

1 **Multi-ancestry GWAS reveals loci linked to human variation in LINE-1- and Alu-**  
2 **insertion numbers**

3

4 Juan I. Bravo<sup>a</sup>, Lucia Zhang<sup>a,b</sup>, Bérénice A. Benayoun<sup>a,c,d,e,f,\*</sup>

5 <sup>a</sup> Leonard Davis School of Gerontology, University of Southern California, Los Angeles,  
6 CA 90089, USA.

7 <sup>b</sup> Quantitative and Computational Biology Department, USC Dornsife College of Letters,  
8 Arts and Sciences, Los Angeles, California, USA.

9 <sup>c</sup> Molecular and Computational Biology Department, USC Dornsife College of Letters,  
10 Arts and Sciences, Los Angeles, CA 90089, USA.

11 <sup>d</sup> Biochemistry and Molecular Medicine Department, USC Keck School of Medicine, Los  
12 Angeles, CA 90089, USA.

13 <sup>e</sup> USC Norris Comprehensive Cancer Center, Epigenetics and Gene Regulation, Los  
14 Angeles, CA 90089, USA.

15 <sup>f</sup> USC Stem Cell Initiative, Los Angeles, CA 90089, USA.

16

17 \* Corresponding author. Leonard Davis School of Gerontology, University of Southern  
18 California, Los Angeles, CA, 90089, USA.

19 *E-mail address:* berenice.benayoun@usc.edu (B.A. Benayoun)

20

21

22 **ABSTRACT**

23       LINE-1 (L1) and Alu are two families of transposable elements (TEs) occupying  
24 ~17% and ~11% of the human genome, respectively. Though only a small fraction of L1  
25 copies is able to produce the machinery to mobilize autonomously, Alu and degenerate  
26 L1s can hijack their functional machinery and mobilize *in trans*. The expression and  
27 subsequent mobilization of L1 and Alu can exert pathological effects on their hosts.  
28 These features have made them promising focus subjects in studies of aging where  
29 they can become active. However, mechanisms regulating TE activity are incompletely  
30 characterized, especially in diverse human populations. To address these gaps, we  
31 leveraged genomic data from the 1000 Genomes Project to carry out a trans-ethnic  
32 GWAS of L1/Alu insertion singletons. These are rare, recently acquired insertions  
33 observed in only one person and which we used as proxies for variation in L1/Alu  
34 insertion numbers. Our approach identified SNVs in genomic regions containing genes  
35 with potential and known TE regulatory properties, and it enriched for SNVs in regions  
36 containing known regulators of L1 expression. Moreover, we identified reference TE  
37 copies and structural variants that associated with L1/Alu singletons, suggesting their  
38 potential contribution to TE insertion number variation. Finally, a transcriptional analysis  
39 of lymphoblastoid cells highlighted potential cell cycle alterations in a subset of samples  
40 harboring L1/Alu singletons. Collectively, our results suggest that known TE regulatory  
41 mechanisms may be active in diverse human populations, expand the list of loci  
42 implicated in TE insertion number variability, and reinforce links between TEs and  
43 disease.

44

45

46 **KEYWORDS:** LINE-1, Alu, transposons, insertions, GWAS, regulators

## 47 1. INTRODUCTION

48 In the human genome, the two most abundant families of transposable elements  
49 (TEs) are Long INterspersed Element-1 (LINE-1; L1) and Alu, which account for ~16–  
50 17% and ~9–11% of the genome, respectively [1, 2]. Full-length L1 elements span ~6  
51 kilobases and produce bicistronic messenger ribonucleic acids (mRNAs) encoding two  
52 polypeptides, ORF1p and ORF2p, necessary for L1 transposition (reviewed in [3]). The  
53 L1 family can be segregated into 3 subfamilies depending on the evolutionary age of the  
54 copy: the L1M (mammalian-wide) lineage is the oldest, the L1P (primate-specific)  
55 lineage is of intermediate age, and the L1PA lineage is the youngest. Importantly, only  
56 the L1PA1/L1Hs subfamily contains ~80-100 actively mobile copies in the average  
57 human genome [4], with the remaining ~500,000 L1 copies being rendered non-  
58 autonomous due to the presence of loss-of-function mutations or truncations [1]. In  
59 contrast to L1 elements, Alu elements are short (~300 bp) non-autonomous  
60 retrotransposons that rely on functional L1 machinery for their mobilization [5-7]. Alu  
61 retrotransposons can also be segregated by evolutionary age into the following  
62 subfamilies: AluJ is the oldest lineage and is likely completely inactive in humans, AluS  
63 is the middle-aged lineage and contains mobile copies, and AluY is the youngest  
64 lineage and contains the largest number of functionally intact elements [8].

65  
66 For an expansion of L1 insertions to occur, L1 must undergo a multi-step and  
67 tightly regulated lifecycle. This lifecycle begins with transcription of an active, full-length  
68 L1 copy and ends with reverse transcription and integration of L1 by a process called  
69 target primed reverse transcription (TPRT) (reviewed in [3]). Importantly, though neither  
70 Alu elements or degenerate L1 copies can mobilize autonomously, they can hijack  
71 proteins from transposition-competent L1s and mobilize *in trans* [6, 7, 9]. Though not  
72 traditionally considered part of the L1/Alu lifecycles, other genome remodeling (i.e.  
73 recombination or DNA repair) mechanisms can further contribute to TE insertion  
74 number variation. This includes, for example, repeat-mediated deletion (RMD) events  
75 whereby two repetitive elements (often Alu elements) on the same chromatid recombine  
76 and cause the deletion of one of the repeats as well as the intervening sequence, which  
77 may include additional repeats [10, 11]. More broadly, non-allelic homologous

78 recombination (NAHR) events ([12-15] and reviewed in [16, 17]) can directly generate  
79 large chromosomal deletions and duplications, which may include repetitive sequences.  
80 Ultimately, TE insertion number variation is shaped by a combination of *de novo*  
81 insertions resulting from their lifecycle and genome remodeling that can expand or  
82 retract the number of insertions.

83

84 Characterizing the mechanisms governing L1 and Alu transcriptional and  
85 mobilization control will be important, given their associations with, and potential  
86 contributions to, aging and aging-associated diseases like cancer (discussed in [18-20]).  
87 Fundamentally, L1 and/or Alu can alter several hallmarks of aging [21], such as  
88 genomic instability, cellular senescence, and inflammation. Though the origin of the  
89 signal is unclear (genomic, extra-chromosomal, or cytosolic), an increase in L1 copies  
90 has been observed with chronological aging [22] and during cellular senescence [23].  
91 Moreover, a key feature of cellular senescence is the senescence-associated secretory  
92 phenotype (SASP) whereby cells secrete an amalgamation of pro-inflammatory factors  
93 [24] that may contribute to chronic, low-grade, sterile inflammation with chronological  
94 age (a phenomenon referred to as “inflamm-aging”) [24, 25]. L1 can induce a  
95 senescent-like state *in vitro* in several cell lines [26, 27] and its cytoplasmic  
96 complementary DNA (cDNA) is implicated in the maturation of the SASP response and  
97 the establishment of deep senescence through the production of interferons [28].  
98 Similarly, Alu RNA can upregulate senescence markers in retinal pigment epithelium  
99 (RPE) cells from human eyes with geographic atrophy [29], and knockdown of Alu  
100 transcripts was reported to promote senescence exit in adult adipose-derived  
101 mesenchymal stem cells [30]. These findings highlight the relevance of L1 and Alu  
102 retrotransposons in pathological, age-associated features and highlight the important of  
103 characterizing TE control mechanisms.

104

105 To maintain homeostasis, it is imperative that host cells tightly regulate TE  
106 activity (reviewed in [31, 32]). These mechanisms, however, remain incompletely  
107 characterized due to the cell-specific and transposon-specific nature of TE regulatory  
108 mechanisms. Indeed, no systematic, genome-wide screen for regulators of Alu

109 expression or mobilization has been carried out thus far, to our knowledge, and our  
110 understanding of L1 control mechanisms remains incomplete, limiting our ability to  
111 understand why they are derepressed during aging. To address these gaps, a number  
112 of *in vitro* and *in silico* approaches have been developed to scan for novel regulators of  
113 TE expression or mobilization. *In vitro* approaches have relied on clustered regularly  
114 interspaced short palindromic repeats (CRISPR)-based and small RNA-based tools to  
115 decipher L1 regulation in several types of cancer cells [33-37]. These approaches,  
116 however, can be technically challenging to implement in non-cancerous cells, like  
117 primary cells, which may not tolerate hyper-elevated transposon activity or that may  
118 resist genetic perturbations. To complement these methods, a number of *in silico*  
119 approaches have been developed that utilize chromatin immunoprecipitation followed  
120 by sequencing (ChIP-seq) [38], gene co-expression networks [39], or insertion number-  
121 expression correlations [40] to explore L1 regulation without external manipulations.  
122 More recently, we screened for candidate regulators of L1 RNA levels in lymphoblastoid  
123 cell lines (LCLs) using trans-expression quantitative trait locus (trans-eQTL) analysis  
124 [41]. These tools highlight the need for, and usefulness of, alternative approaches that  
125 utilize increasingly available, large ‘-omic’ datasets to identify potentially novel  
126 mechanisms of TE control.

127  
128 In this study, we develop a genome-wide association study (GWAS) pipeline to  
129 identify genomic loci associated with variation in global L1 and Alu insertion singletons  
130 in diverse human populations. Global singleton insertions reflect rare, recently acquired  
131 insertions observed uniquely in only one person in a population [42]. Thus, we used  
132 insertion singletons as proxies for L1/Alu insertion number variation, which can arise  
133 through *de novo* transposition events or genome remodeling mechanisms. We  
134 demonstrate that our GWAS approach captures, and enriches, genomic regions  
135 containing known and potential regulators of TE activity. We observe that this approach  
136 also captures reference insertions and polymorphic structural variants that may  
137 influence L1 or Alu insertion number variation. Finally, we note that associated loci fall  
138 into a few genes with clinical relevance, strengthening the association between TEs and  
139 disease.

## 140 2. RESULTS

### 141 2.1 Identification of genomic loci associated with L1/Alu singletons in diverse human 142 populations

143 To unbiasedly identify potential genetic sources of L1 and Alu insertion number  
144 variation in human populations, we leveraged a publicly available human “omic” dataset  
145 with thoroughly characterized genetic information. For this analysis, we utilized 2503  
146 multi-ethnic samples from the 1000 Genomes Project for which both single nucleotide  
147 variant (SNV) and structural variant (SV) data were available. Specifically, this included  
148 individuals from 5 super-populations: 660 African (AFR), 504 East Asian (EAS), 503  
149 European (EUR), 489 South Asian (SAS), and 347 Admixed American (AMR)  
150 individuals who declared themselves healthy at the time of sample collection (**Figure**  
151 **1A**). As a quality control step, we checked whether the combined SNV and SV data  
152 segregated samples by population following principal component analysis (PCA). These  
153 analyses demonstrated that the top four principal components segregated population  
154 groups within each super-population (**Figure S1A-S1E**).

155

156 We then carried out an integration of the available multi-ethnic SNV and SV  
157 genomic data (**Figure 1B**). For the phenotype, we focused on global singleton SVs,  
158 which are rare SVs that are observed exactly once in a single person [42], for L1 and  
159 Alu insertions. The rarity of these insertions is similar to, but more stringent than, the  
160 threshold used in a study of *Arabidopsis* that explored the genetic basis of variable  
161 transposition rate [43]. We employed this more stringent threshold (observed in only 1  
162 individual) to increase the likelihood of capturing a recent mobilization event (likely to  
163 have occurred in the germline or the early embryonic stage). Unique, or private,  
164 insertions are likely to be family-specific rather than population-specific and suggest a  
165 recent relaxation of transposon control. For all these reasons, we hypothesized that  
166 L1/Alu insertion singletons could serve as proxies for elevated L1/Alu insertion number  
167 variation.

168

169 Thus, we first split our samples into cases and controls, depending on whether or  
170 not they contained an L1 and/or an Alu insertion singleton (**Supplementary Table**

171 **S1A)**. Second, we carried out a GWAS within each super-population to identify  
172 common, polymorphic SNVs and SVs associated with case-control status. Third, to  
173 maximize statistical power and identify shared, trans-ethnic sources of TE singleton  
174 number variation, we meta-analyzed our GWAS results across the 5 super-populations  
175 using a random effects statistical model, which allows for increased generalizability  
176 across diverse human cohorts compared to fixed effects models [44, 45]. Interestingly,  
177 several hundred L1/Alu global singleton insertions were detected in each super-  
178 population, ranging from 322 in the American cohort to 866 in the African cohort (**Figure**  
179 **1C**). Though most case samples had 1-3 L1/Alu global singletons, several samples  
180 exhibited much more extreme TE singleton accumulation, especially within the African  
181 super-population (**Figure 1C**). Finally, we note that L1/Alu global singleton insertions  
182 were distributed across all autosomes in all 5 super-populations (**Figure 1C**). Many  
183 L1/Alu singletons resided in intronic (907/3024 ~ 30%) and intergenic (834/3024 ~  
184 27.6%) regions and were enriched for these features compared to background (**Table**  
185 **1**). This is partially consistent with a prior report where *de novo* engineered L1 insertions  
186 were not enriched in genes and endogenous L1 elements were depleted from genic  
187 regions in one cell line [46], though the preferences for introns and exons in that study  
188 are unclear. Theoretically, the deleterious effects of insertions within intronic and  
189 intergenic regions may be blunted compared to the effects of insertions within exonic  
190 regions, and we speculate that this may be part of the reason for L1/Alu insertion site  
191 preferences. In contrast, L1/Alu singletons were depleted in regions already containing  
192 LINE or SINE transposons, compared to their background distribution (**Table 1**). We  
193 speculate that the target motif needed for L1/Alu mobilization may be disrupted in these  
194 regions, restricting L1/Alu integration into these sites.

195

196 As expected, GWAS in each super-population was generally underpowered  
197 (**Figure S2A-S2E**). Though we were able to identify many significant (FDR < 0.05)  
198 variants in the African cohort (**Figure S2A**), we could not identify significant variants in  
199 the East Asian (**Figure S2B**) and European (**Figure S2C**) cohorts, and we were only  
200 able to identify a handful of significant variants in the South Asian (**Figure S2D**) and  
201 Admixed American (**Figure S2E**) cohorts. These observations, and the abundant

202 associations observed in the African super-population in particular, may be related to  
203 the high genetic diversity among African populations [47, 48], including containing more  
204 mobile element variants and more population-specific mobile element variants  
205 compared to other super-populations [49-51]. In contrast to the individual analyses, the  
206 GWAS meta-analysis integrating all super-populations identified 658 significant variants  
207 distributed across all 22 autosomes, though there was especially strong and recurrent  
208 signal on chromosome 21 (**Figure 1D, Table S1B**). Importantly, to ensure that  
209 associations were not completely dependent on the African super-population, we also  
210 ran the meta-analysis using only the four non-African super-populations (**Figure S3A-**  
211 **S3C**). Though this analysis had a smaller sample size, we identified 194 significant  
212 variants that were shared with the complete meta-analysis and 21 variants that were  
213 unique (**Figure S3C, Table S1C**), suggesting that the variants identified in the complete  
214 meta-analysis are robust. Moreover, there was significant ( $p = 5.26E-213$ ) overlap  
215 between the variants identified in the non-African meta-analysis and the African-only  
216 analysis (**Figure S3D**), further demonstrating that i) the genomic architecture underlying  
217 L1/Alu singleton variation is similar between non-African and African super-populations  
218 and ii) a subset of variants can be independently identified in two separate populations.  
219 Since super-population-specific variants were limited, we focused on the results from  
220 the complete meta-analysis. To simplify functional annotation of significant variants and  
221 discard potential false positives, we omitted from downstream analyses significant  
222 variants overlapping the “ENCODE blacklist v2” [52]. After filtering 188 blacklisted  
223 SNVs, we were left with 449 greenlisted SNVs and 21 greenlisted SVs that were  
224 significantly associated with case-control status (**Figure 1D, Table 2, Supplementary**  
225 **Table S1B**).

226

227 To assess the potential functions of greenlisted, significant variants, nearby  
228 genes were assigned to variants and over-representation analysis (ORA) was carried  
229 out (**Figure S4A**). Except for 51 greenlisted SNVs that were not linked to any gene, the  
230 remaining SNVs were linked to 1-2 genes and were found within 1000 kilobases of a  
231 transcriptional start site (TSS) (**Figure S4B**). This observation highlights the association  
232 of intergenic and gene-proximal, rather than distal, genetic variation with L1/Alu



233 insertion number differences. Over-representation analysis of the associated genes  
234 using the Gene Ontology (GO) Biological Process gene set revealed an enrichment of  
235 terms related to heart development (such as ‘regulation of heart growth’, ‘cardiac  
236 chamber morphogenesis’, and ‘positive regulation of cardiac muscle cell proliferation’)  
237 and neuronal function (such as ‘neuron recognition’ and ‘axonal fasciculation’; **Figure**  
238 **S4B, Table S1D**). Interestingly, genes related to ‘reproduction’ were also significantly  
239 over-represented among our list of associated genes (**Supplementary Table S1D**).  
240 Similar to the SNVs, greenlisted SVs were all linked to 1-2 genes and were within 1000  
241 kilobases of an annotated TSS (**Figure S4C**). Likely due to the low number of  
242 greenlisted SVs, and consequently low number of associated genes, we were unable to  
243 identify any significantly enriched GO Biological Process gene sets (**Supplementary**  
244 **Table S1E**). Given the limited number of greenlisted SVs and the unavailability of SV  
245 sequences, we largely focused on greenlisted SNVs in downstream enrichment  
246 analyses.

247  
248 As a complementary approach, we also predicted the functional impact of  
249 significant variants using SnpEff [53] (**Supplementary Table S1F**). Most variants were  
250 assigned a ‘modifier’ impact by SnpEff—this annotation describes non-coding variants  
251 where definitive functional predictions are not straightforward. Nonetheless, we highlight  
252 a few variants which overlapped clinically relevant genes. For example, the most  
253 significant, greenlisted variant we identified was an inversion SV  
254 (INV\_delly\_INV00066128,  $p = 5.67E-23$ , odds ratio = 4.38) residing in an intronic Alu  
255 copy within the *APP* (amyloid beta precursor protein) gene, an important biological  
256 marker for Alzheimer’s disease (AD) [54]. Similarly, we identified several SNVs  
257 (rs61994687,  $p = 4.11E-7$ , odds ratio = 0.54; rs1175403595,  $p = 3.85E-8$ , odds ratio =  
258 0.44; rs1343402870,  $p = 4.49E-6$ , odds ratio = 0.43) in intronic or downstream regions  
259 of *PWRN1* within the Prader-Willi syndrome (PWS) region. To explore non-protein-  
260 coding roles greenlisted SNVs may play, we assigned them to ENCODE candidate cis-  
261 Regulatory Elements (cCREs) [55] (**Figure S4D**). Although about 8% (36/449) of  
262 greenlisted SNVs resided in an ENCODE cCRE, these were significantly depleted (FDR  
263 =  $6.86E-15$ ) in our greenlisted SNVs compared to background SNVs (**Figure S4D**).

264 Ultimately, our results suggest that proximal, intergenic variation is associated with  
265 L1/Alu insertion singleton number variation.

266

267 To further explore the potential functional roles of SNV-associated genes, we  
268 leveraged expression data from the Genotype-Tissue Expression (GTEx) Portal [56-58]  
269 to assess the patterns of tissue expression of SNV-linked genes (**Figure S4E**,  
270 **Supplementary Table S1G**). In particular, expression patterns in the brain and gonads  
271 were of special interest, given that L1 activity tends to be more frequent in those tissues  
272 compared to others (discussed in [59]), and that L1/Alu integration events observed in  
273 our GWAS would have to occur in the germline for transmission across generations.  
274 Thus, we reasoned that if SNV-associated genes played roles in L1/Alu singleton  
275 number variation, they may be more abundantly expressed in the brain and in gonads.  
276 Interestingly, there was a cluster of genes that was very abundantly expressed in testes  
277 but not in other tissue types (including ovarian tissue), suggesting the existence of  
278 potential sex-specific mechanisms of *de novo* L1/Alu insertion transmission.  
279 Furthermore, there was also a cluster of genes that were abundantly expressed across  
280 brain regions and much less abundantly expressed across other tissue types. More  
281 generally, there were many SNV-linked genes that were abundantly expressed in more  
282 than ~50% of tissue types. These results highlight that significant SNV-associated  
283 genes have tissue-specific expression patterns, including some genes that are very  
284 abundant in tissues with documented L1 activity.

285

## 286 *2.2 Significant SNVs are enriched near regulators of L1 expression*

287 One of the primary motivations for carrying out this study was to identify novel,  
288 candidate regulators of L1 and/or Alu mobilization. To determine whether our approach  
289 captured genes with transposon regulatory potential, we assessed whether our list of  
290 greenlisted SNVs was enriched for 1) genes with known TE regulatory capabilities and  
291 2) genes in broader pathways involved in TE regulation (**Figure 2A**). Interestingly, our  
292 greenlisted SNVs were significantly (FDR = 3.59E-3) enriched in regions near known L1  
293 expression regulators compared to the background list of all SNVs (**Figure 2B, 2C, 2E**).  
294 A few examples of these associations included rs1350516110 ( $p = 9.35E-12$ , odds ratio

295 = 0.33) which was upstream of *RHOT1*, rs201619112 ( $p = 6.24E-9$ , odds ratio = 0.33)  
296 which was downstream of *XPR1*, and rs71475866 ( $p = 1.65E-8$ , odds ratio = 2.54)  
297 which was upstream of *PFKP*. Overall, we identified 24 greenlisted SNVs that were  
298 proximal to 10 genes previously annotated as capable of regulating L1 expression. Our  
299 greenlisted SNVs also captured genes involved in regulating L1 transposition, though  
300 there was no significant enrichment ( $FDR = 7.78E-1$ ) (**Figure 2B, 2D, 2E**). A few  
301 examples of these associations included rs75237296 ( $p = 2.14E-9$ , odds ratio = 0.42) in  
302 the *PABPC1* 3'UTR, rs1288384419 ( $p = 3.17E-8$ , odds ratio = 0.44) in an intron of  
303 *RAD51B*, and rs1471205623 ( $p = 1.34E-7$ , odds ratio = 0.42) upstream of *MPHOSPH8*.  
304 Here, we identified 5 greenlisted SNVs near 3 genes capable of regulating L1  
305 transposition. Finally, we checked the abundances of SNVs linked to candidate  
306 regulators of L1 RNA levels in lymphoblastoid cells [41]. However, we were not able to  
307 detect any greenlisted SNVs in regions containing candidate genes (**Figure 2B and 2E**).  
308 Nevertheless, these results highlight the ability of our approach to enrich for genomic  
309 regions containing known regulators of the retrotransposon lifecycle and suggest that  
310 these regulators may play important roles in diverse human populations.

311

312 We next repeated the above analyses using gene sets for broader pathways  
313 involved in TE regulation, including a gene set for “histone methyltransferase activity”  
314 (GO:0042054) and one for “RNA modifications” (GO:0009451). Though neither gene set  
315 was significantly enriched among our greenlisted SNVs ( $FDR = 1$  for methyltransferase  
316 activity and  $FDR = 0.27$  for RNA modification), there was some degree of overlap with  
317 each gene set (**Figure 2E**). We identified 4 greenlisted SNVs that were proximal to 2  
318 genes with histone methyltransferase activity, including *EEF2KMT* and *PRDM7*. We  
319 also identified 8 greenlisted SNVs that were proximal to 3 genes with RNA modification  
320 capabilities, including *A1CF*, *ADARB2*, and *METTL14*. Importantly, ADARs (RNA-  
321 specific adenosine deaminases) are a family of double-stranded RNA (dsRNA)-binding  
322 proteins that modulate A-to-I editing events, including among Alu RNA species, which is  
323 important for preventing aberrant activation of innate immune signaling pathways [60].  
324 Though *ADARB2* cannot catalyze A-to-I editing, it can negatively regulate the editing  
325 functions of other ADARs [60], making it a potential candidate regulator of Alu activity.

326 These results demonstrate that our approach can capture genes implicated, but with  
327 uncharacterized roles, in TE regulation.

328

### 329 *2.3 Significant SNVs are enriched in regions containing features that promote genome* 330 *instability*

331 After scanning for known and potential regulators of TE activity, we next explored  
332 the possibility that significant variants tagged genetically unstable TE loci (**Figure 3A**).  
333 Such loci could theoretically contribute to TE insertion number variation through *de novo*  
334 transposition events. To probe this possibility, we assessed whether greenlisted SNVs  
335 were enriched for Alu and L1 loci belonging to subfamilies that have retained their ability  
336 to mobilize (**Figure 3B**). Interestingly, SNVs overlapping the active AluS subfamily were  
337 significantly (FDR = 4.80E-6) enriched in our SNV greenlist compared to background  
338 (**Figure 3B**). Our enrichment analysis also highlighted a significant (FDR = 7.29E-7)  
339 depletion of inactive-L1M-overlapping SNVs and a significant (FDR = 3.46E-12)  
340 enrichment of L1PA-overlapping SNVs, all compared to background SNVs (**Figure 3B**).  
341 To obtain a higher resolution view of the transposition capabilities of overlapping L1PA  
342 copies, we checked the overlap of our greenlisted SNVs with annotations for putatively  
343 active L1 copies. Surprisingly, active copies, with either an intact, full-length L1 or only  
344 an intact ORF2, were not significantly enriched/depleted in our greenlisted SNV list  
345 (**Figure 3C**). However, non-intact, full-length L1 copies—annotated for their regulatory  
346 potential—were significantly (FDR = 2.41E-6) enriched compared to background. These  
347 results suggest that most greenlisted SNV-overlapping L1 copies are limited in their  
348 ability to generate *de novo* insertions and thus may influence insertion number through  
349 alternative mechanisms (i.e. genome remodeling). This is potentially in contrast to  
350 greenlisted SNV-overlapping AluS copies, which may still be measurably active in the  
351 human genome. The TE enrichments we identified above were consistent with those  
352 identified using the Transposable Element Enrichment Analyzer (TEENA) tool [61]  
353 (**Supplementary Table S1H**).

354

355 We further explored the more general possibility that greenlisted SNVs tagged  
356 genomic regions that may be prone to genome remodeling that may influence TE

357 insertion numbers. Such structural alterations may be facilitated by, but may not require,  
358 the presence of repetitive elements. In particular, extensive homology between  
359 segmental duplications, often in the vicinity of Alu elements [62], can facilitate NAHR  
360 and drive recurrent genomic rearrangements [63] that can help form SV hotspots [64,  
361 65]. Noting the enrichment of AluS copies we observed among our greenlisted SNVs,  
362 we next assessed whether our greenlisted SNVs significantly overlapped regions of  
363 segmental duplication [66, 67] or regions characterized as SV hotspots [65] (**Figure**  
364 **3D**). Consistent with the notion that SNVs may tag regions with potentially elevated  
365 rates of genome instability, our SNVs were very significantly enriched in regions of  
366 segmental duplication (FDR = 2.64E-99), as well as in regions harboring SV hotspots  
367 (FDR = 8.82E-7). These results further link variation in TE insertion numbers to genomic  
368 loci where structural instability may arise through genome remodeling mechanisms.

369

370 Finally, we note that our GWAS analysis identified 21 polymorphic SVs that were  
371 significantly (FDR < 0.05) associated with the presence/absence of L1/Alu insertion  
372 singletons. These polymorphic SVs varied in nature and included inversions, Alu  
373 insertions, an L1 insertion, SINE-VNTR-Alu (SVA) insertions, and a multiallelic copy  
374 number variant (**Figure S5A**). With the exception of the CN0 copy number variant  
375 (YL\_CN\_PEL\_1784,  $p = 3.30E-6$ , odds ratio = 0.28) which was associated with lower  
376 odds of carrying an L1/Alu insertion singleton, all of the other structural variants were  
377 associated with higher odds of carrying an insertion singleton (odds ratio > 1). Since the  
378 sequences for these SVs were not available, it is unclear whether common, polymorphic  
379 L1/Alu insertion SVs may be directly increasing the singleton number through novel  
380 transposition events, or whether any of these polymorphic SVs may be influencing the  
381 L1/Alu singleton number through genome remodeling mechanisms. Indeed,  
382 polymorphic inversions, many of which are often flanked by retrotransposons, are  
383 associated with genetic instability [68]. Ultimately, these results suggest a tight link  
384 between common, polymorphic SVs of different types and L1/Alu singleton SVs,  
385 whereby having the former is generally associated with higher odds of having the latter.

386

387 *2.4 Case samples exhibit elevated cell cycle-related gene expression profiles*

388 To gain insight into the functional differences between controls and cases, we  
389 leveraged publicly available lymphoblastoid cell line mRNA-seq data generated by the  
390 GEUVADIS consortium for a subset of European and African samples in the 1000  
391 Genomes Project [69] (**Figure 4A**). This included 358 European samples from 4  
392 populations (British, Finnish, Tuscan, and Utah residents with European ancestry) and  
393 86 African samples from 1 population (Yoruba), which we recently used to quantify gene  
394 and TE expression profiles [41]. We utilized this expression data to construct consensus  
395 gene co-expression networks for both the European and African samples using the  
396 WGCNA [70] package. This approach led to the identification of 20 consensus modules  
397 and 1 module (MEgrey) containing genes that were not assigned to the consensus  
398 modules (**Supplementary Table S1I**). We then ran a module-trait correlation analysis  
399 comparing the expression of these modules with the case/control status of the  
400 European and African samples (**Figure 4B**). Here, we used a stricter threshold of  $p <$   
401  $0.01$  to call significant correlations. We were not able to identify any significant module-  
402 phenotype correlations using the European network, which is potentially consistent with  
403 our difficulty in calling significant GWAS variants in this super-population at the available  
404 sample sizes (**Figure S2C**). In contrast, the MEroyalblue module was significantly ( $p =$   
405  $4.0E-4$ ) correlated with African case/control status. To combine the results from each  
406 network, we utilized Fisher's method to meta-analyze the p-values for modules  
407 exhibiting correlations in the same direction. By meta-analysis, the MEroyalblue module  
408 was still significantly ( $p = 0.002$ ) and positively correlated with case status. Finally, to  
409 functionally characterize this module, we ran over-representation analysis using the GO  
410 Biological Process gene set collection (**Figure 4C, Supplementary Table S1J**). The top  
411 10 over-represented gene sets were involved in cell cycle-related processes, including  
412 "mitotic cell cycle", "cell division", and "sister chromatid segregation". These findings are  
413 consistent with the biology underlying an expansion of TE insertions. Though L1 can  
414 mobilize in non-dividing cells [71, 72], L1 retrotransposition exhibits a cell cycle bias and  
415 peaks during the S phase [73]. Alternatively, chromosome segregation errors during  
416 mitosis or meiosis can generate cells with abnormal ploidy and either increased or  
417 decreased dosages of both genic and transposon content [74]. These results implicate

418 cell cycle differences in cells from individuals with unique L1/Alu insertion singleton  
419 variation.

420

421

## 422 3. DISCUSSION

### 423 3.1 A new approach to identify loci implicated in L1 and Alu insertion number variation

424 In this work, we developed a pipeline to computationally identify candidate loci  
425 involved in L1/Alu singleton number variation by GWAS analysis. Importantly, our study  
426 incorporates natural human genetic variation present in populations of different  
427 geographic origin via trans-ethnic GWAS meta-analysis to identify shared, candidate  
428 regulatory loci. Though several studies have begun to screen for regulators and  
429 potential regulators of L1 expression or transposition in cell culture models or across  
430 tissues [33-40], these can be limited in that the generalizability of these findings to  
431 different ethnic populations is unclear. Moreover, no systematic, genome-wide screen  
432 for candidate regulators of Alu activity has been carried out thus far, to our knowledge.  
433 To address these gaps, we previously utilized trans-eQTL analysis to identify potential  
434 regulators of L1 RNA levels in European and African populations [41]. Here, we utilized  
435 genomic data from samples originating from 5 super-populations to identify candidate  
436 loci modulating L1/Alu insertion singleton numbers.

437

438 TE insertion number variation can arise through *de novo* transposition events or  
439 through genome remodeling mechanisms that can generate large deletions or  
440 duplications. We were particularly interested in identifying new candidate regulators of  
441 L1/Alu transposition. Consistent with the notion that greenlisted SNVs may play roles in  
442 the retrotransposon lifecycle, our approach enriched genomic regions containing genes  
443 that can regulate L1 expression levels. Though other known regulators of TE activity,  
444 and pathways involved in TE control, were not significantly enriched among our  
445 greenlisted SNVs, we nonetheless identified many SNVs in genomic regions containing  
446 these genes. This included, for example, *MPHOSPH8*—a component of the HUSH  
447 complex important for L1 repression, regulating both L1 expression and transposition  
448 [33, 37, 75-77]. It also included *PABPC1* which is a poly(A) binding protein that attaches  
449 to the poly(A) tail of L1, is important for the formation of L1 ribonucleoprotein particles  
450 (RNPs), and modulates L1 and Alu transposition [78-80]. As a third example, we identified  
451 variants near *ADARB2*, a negative regulator of RNA A-to-I editing, including among Alu  
452 RNAs [60]. These results suggest that SNV-associated genes identified in this study



453 hold TE regulatory potential and it may therefore be informative to (i) test and validate  
454 these in future studies or (ii) use these to prioritize future, targeted studies of TE  
455 regulators.

456

457 Our approach also identified an enrichment of greenlisted SNVs in regions  
458 containing reference TE insertions, including AluS and full-length, non-intact L1PA  
459 copies. Though neither of these can mobilize autonomously, they can hijack machinery  
460 from transposition-competent L1s and mobilize *in trans* [6, 7, 9]. Thus, it is possible that  
461 greenlisted SNVs tag reference insertions contributing to L1/Alu singleton variation  
462 through transposition-dependent mechanisms. Of course, an alternative possibility is  
463 that these repetitive elements are directly involved in genomic remodeling involving  
464 mechanisms like repeat-mediated deletions or NAHR. We also note that greenlisted  
465 SNVs were enriched in regions containing segmental duplications and structural  
466 variation hotspots where recombination-based mechanisms, including NAHR, may lead  
467 to duplications or deletions of the local genomic architecture. Thus, it is also possible  
468 that greenlisted SNVs tag genomic regions prone to structural variation that can alter  
469 the L1/Alu insertion number through recombination-dependent mechanisms.

470

471 Importantly, we also suggest the possibility that genome remodeling mechanisms  
472 (including recombination) may interact with gene-based mechanisms of TE regulation.  
473 Indeed, genes such as *BRCA1* are known to regulate L1 transposition [33] and are also  
474 known to undergo Alu-Alu recombination events that can give rise to new mutations in  
475 the gene [81-83]. Observations such as these highlight the possibility that TE insertions  
476 may modulate structural variation in genomic regions containing genes regulating  
477 retrotransposon lifecycles, which may facilitate an expansion of TE insertions through  
478 transposition-based mechanisms, which may influence further structural variation  
479 driving this whole process. This possibility is consistent with the enrichment of  
480 greenlisted SNVs in regions containing L1 expression regulators, Alu and L1 repeats,  
481 and other genomically unstable features like segmental duplications and structural  
482 variation hotspots. In the future, it may be informative to experimentally assess this

483 possibility in contexts where genome instability is a hallmark feature that is coupled with  
484 TE de-repression, such as aging [21] or aging-associated diseases like cancer [84, 85].

485  
486 Finally, our approach also identified several common, polymorphic SVs that were  
487 significantly associated with L1/Alu insertion singletons. Overwhelmingly, the presence  
488 of polymorphic SVs of different types—inversions, Alu insertions, an L1 insertion, and  
489 SVA insertions—was associated with increased odds of a global L1/Alu insertion  
490 singleton. The exception to this was a multi-allelic copy number variant  
491 (YL\_CN\_PEL\_1784,  $p = 3.30E-6$ , odds ratio = 0.28) where 0 copies were present; this  
492 SV was associated with decreased odds of a global L1/Alu insertion singleton. Based  
493 on these results, we speculate that specific polymorphic SVs (i) may directly drive  
494 genome instability that can facilitate the acquisition of L1/Alu copies or (ii) may serve as  
495 markers for elevated risk of indirectly acquiring additional L1/Alu copies. Indeed, active  
496 donor L1 copies that can mobilize and generate *de novo* insertions are usually highly  
497 polymorphic in human populations (reviewed in [59]), and polymorphic inversions, many  
498 of which are often flanked by retrotransposons, are also associated with genetic  
499 instability and genomic disorders [68].

500  
501 In summation, this study provides a list of variants that are associated with L1/Alu  
502 insertion singletons and includes (i) SNVs in regions containing regulators of TE activity,  
503 (ii) SNVs in regions containing features associated with genome instability, including  
504 retrotransposons, that may influence TE insertion number variation through  
505 transposition-dependent or transposition-independent mechanisms, and (iii) common,  
506 polymorphic SVs that may also influence TE insertion number variation through  
507 transposition-dependent or transposition-independent mechanisms.

508

509

### 510 *3.2 L1/Alu insertion singleton-associated loci contain genes of clinical relevance*

511 The most significant, greenlisted variant we identified was a polymorphic  
512 inversion SV (INV\_delly\_INV00066128, chr21:26001780,  $p = 5.67E-23$ , odds ratio =  
513 4.38) residing in an intronic Alu copy within the *APP* gene, an important marker of

514 Alzheimer's disease (AD). AD is characterized by (i) the accumulation of amyloid  $\beta$  ( $A\beta$ )  
515 plaques derived from amyloidogenic *APP* processing and (ii) neurofibrillary tangles of  
516 hyperphosphorylated tau [86]. Importantly, tau protein can induce TE expression and  
517 there is speculation that TEs may mobilize in tauopathies [86]. Whether APP protein or  
518  $A\beta$  plaques can similarly modulate TE expression or mobilization is an interesting area  
519 of potential future research; indeed, our findings are consistent with the possibility that  
520 *APP* products may act as regulators of TE activity. Of course, another possibility is that  
521 the genomic region containing *APP* may be a source of L1/Alu insertion number  
522 variation independent of the functional properties of APP protein (i.e. genomic instability  
523 at that locus may be the driver of TE insertion number variation). Indeed, repeat-based  
524 recombination events have been documented in AD [87]. Nevertheless, our results offer  
525 another connection between transposable element regulation and Alzheimer's disease.

526

527 We also identified several SNVs proximal to *PWRN1*, which resides within the  
528 Prader-Willi syndrome region and is thought to play a role in PWS [88]. Prader-Willi is  
529 an imprinting disorder where genes in the chromosome 15q11-q13 region are  
530 maternally imprinted and paternal copies are not expressed [89]. This lack of paternal  
531 gene expression is predominantly caused by *de novo* paternally inherited deletions of  
532 the 15q11-q13 region, though, less frequently, inheritance of two maternal chromosome  
533 15 copies is the cause [89]. Importantly, a feature of genomic disorders like PWS is the  
534 presence of segmental duplications that can serve as substrates for NAHR [90, 91].  
535 Thus, we hypothesize that this particular region might be more prone to L1/Alu insertion  
536 number variation as a consequence of recombination-based chromosomal alterations.  
537 Nevertheless, it is unclear (i) whether L1/Alu repeats are differentially active in PWS  
538 compared to healthy controls or, more specifically, (ii) whether genes like *PWRN1* can  
539 differentially express or mobilize L1/Alu transposons. Ultimately, the associations  
540 between L1/Alu singletons and both *APP* and *PWRN1* further implicate  
541 retrotransposons in disease.

542

543 Though we did not identify aging-specific signatures in our analyses, these  
544 results should nonetheless be relevant for carrying out a comprehensive analysis of

545 aging-associated genetic risk factors. Though further validation work is needed, it is  
546 hypothesized that L1 insertions increase with chronological aging [22] and during  
547 cellular senescence [23]. Moreover, both L1 and Alu can promote features of cellular  
548 senescence [26-30]. Though these associations are increasingly being characterized,  
549 the study of transposable elements has generally been limited, and their relationship  
550 with aging processes remain incompletely characterized. Thus, this study may serve to  
551 accelerate identification of these relationships by providing an initial set of high-  
552 confidence variants that can be incorporated into future genetic scans for aging-  
553 associated risk factors. Given the ever-increasing availability of large-cohort-based  
554 association studies, these analyses should provide immediate, human-relevant, and  
555 novel aging molecular targets.

556

557

### 558 *3.3 Limitations and future considerations*

559 In this study, we sought to identify new, candidate genes implicated in Alu and L1  
560 insertion number control. One specific mechanism of interest by which this can occur is  
561 target-primed reverse transcription (TPRT)-mediated transposition. Since insertion  
562 sequences for 1000 Genomes Project samples were not available, to our knowledge, it  
563 is difficult to assess to what degree TPRT is driving the associations we identified, since  
564 we cannot check for sequence features of TPRT. In studies with larger cohorts where  
565 insertion sequences are available and insertions with TPRT features can be identified,  
566 our approach could theoretically be applied to explore the genetic basis of TPRT-  
567 specific insertion number variation. Of course, our approach is generally restricted to  
568 identifying genomic loci where variation is common across human populations. We  
569 note, however, that significant variants were enriched in regions containing genes  
570 involved in L1 expression regulation. Since expression is one of the early steps of the  
571 L1 life cycle, our approach captured variants and genes with potential significance to  
572 TPRT-mediated transposition.

573

574 We also note that there are several variables that we are unable to control for in  
575 this study. To protect patient privacy, biological covariates such as chronological age

576 were not available and therefore could not be corrected for in our analysis. Since  
577 increases in L1 expression and L1 copies have been observed with chronological aging  
578 [22], differences in insertion number may reflect age differences rather than genetic  
579 differences. Importantly, however, samples were considered healthy at the time of  
580 sample collection, potentially mitigating health-related effects on insertion number. The  
581 origin and developmental timing of the rare Alu and L1 insertions we utilized is also  
582 unclear. That is, it is unclear whether global singletons used in this study arose *de novo*  
583 in the germline, arose during life through somatic mutation, or even whether they just  
584 arose during the cultivation of the lymphoblastoid cells used to amplify each sample's  
585 DNA. Depending on when these insertions were acquired, the associations identified in  
586 this study may be relevant for either germ cell or somatic cell TE biology. A potential  
587 avenue of future research to address this question would be the incorporation of trio  
588 parent-child genome sequencing [92] and multi-generational genome sequencing to  
589 help identify *bona fide de novo* insertions and their developmental timing. Additionally,  
590 to further strengthen reliability and generalizability of our findings, it will be important to  
591 assess whether significant variants identified in this study are also identified in other  
592 independent cohorts, such as those in the UK Biobank.

593

594 More broadly, though we employed GWAS to scan most of the genome for  
595 variants associated with differences in genomic TE content, we note that the X and Y  
596 chromosomes were omitted from our scan for several reasons. These included: (1) the  
597 necessary data was of poor quality, with over 80% of the X chromosome not genotyped  
598 by the 1000 Genomes Project, (2) the necessary data from the Y chromosome was not  
599 available, to our knowledge, and (3) the analyses require additional analytical  
600 considerations (i.e. male hemizyosity and X chromosome inactivation) that warrant  
601 their own study design [93-96]. Though the data availability and data quality limitations  
602 have been partially addressed by additional sequencing and variant calling in 1000  
603 Genomes Project samples [92], incorporation of this data and of additional analytical  
604 approaches would be beyond the scope of this paper. Thus, the contributions of the X  
605 and Y chromosomes to differences in TE genomic load are a ripe area for future  
606 investigation.

607

608           Ultimately, future studies modulating genes identified with our approach will need  
609 to be carried out to validate any causal contributions of the identified variants to L1/Alu  
610 genomic load. This validation may take on multiple forms, depending on the hypothesis  
611 being tested. For example, if variants tag genes regulating TE mobilization, the tagged  
612 genes can be knocked-down or overexpressed and TE mobilization can be assessed *in*  
613 *vitro* using a retrotransposition assay [97]. Alternatively, if the variants tag genomic  
614 regions prone to structural remodeling (i.e. contribute to NAHR or segmental  
615 duplications), these regions can be sequenced prior to and after a challenge, to assess  
616 whether the number of TE copies has been altered. The underlying mechanisms driving  
617 differences in genomic TE load will likely be multifactorial.

618

619

### 620 *3.4 Conclusions*

621           In this study, we employed GWAS across human populations of different  
622 geographic origin to computationally identify genomic loci associated with variation in  
623 L1/Alu insertion numbers. Our approach enriches for SNVs in genomic regions  
624 containing known regulators of L1 expression. This observation suggests that the TE-  
625 regulatory properties of these genes may extend beyond isogenic cell culture models to  
626 more diverse human populations. Moreover, this observation also suggests that our list  
627 of associated genes likely contains novel regulators of L1 or Alu activity that may be  
628 prioritized in future, validation studies. Our approach also identified reference insertions  
629 and non-reference, polymorphic SVs that may modulate L1/Alu insertion numbers  
630 through transposition-dependent or transposition-independent mechanisms. Finally, the  
631 observation that some significant variants reside in genes of clinical relevance, like *APP*  
632 and *PWRN1*, reinforce accumulating evidence of biological associations between TE  
633 regulation and disease. Ultimately, our approach adds to the analytical toolkit that can  
634 be used to study the regulation of TE activity.

635

636

## 637 **4. METHODS**

### 638 **4.1 Publicly available genomic data acquisition**

639 The multi-ethnic GWAS analysis was carried out on 2503 individuals derived  
640 from 5 super-populations (African, East Asian, European, South Asian, and American)  
641 and for which paired single nucleotide variant and structural variant data were available  
642 from Phase 3 of the 1000 Genomes Project [98-100]. Specifically, Phase 3 autosomal  
643 SNVs called on the GRCh38 reference genome were obtained from The International  
644 Genome Sample Resource (IGSR) FTP site  
645 ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000\\_genomes\\_project/rele](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20190312_biallelic_SNV_and_INDEL/)  
646 [ase/20190312\\_biallelic\\_SNV\\_and\\_INDEL/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20190312_biallelic_SNV_and_INDEL/)). Structural variants, called against the  
647 GRCh37 reference genome and then lifted over to GRCh38, were also obtained from  
648 the IGSR FTP site  
649 ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated\\_sv\\_map/supporting/GRCh](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/supporting/GRCh38_positions/)  
650 [38\\_positions/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/supporting/GRCh38_positions/)).

651  
652 Human gene expression data across 54 non-diseased tissue sites was obtained  
653 from the GTEx Analysis v8 on the GTEx Portal [56-58]. Specifically, we downloaded the  
654 matrix containing the median gene-level transcripts per million (TPMs) by tissue, and we  
655 extracted the expression values for significant SNV-associated genes. After filtering out  
656 genes with no detectable expression (0 TPMs), we generated heatmaps using the  
657 pheatmap v1.0.12 package in R. Gene expression values were centered and scaled  
658 across tissues to visualize and compare the relative expression levels across tissues.

659

### 660 **4.2 Aggregating and pre-processing genotype data for GWAS analysis**

661 To define the phenotype of interest for GWAS, we first extracted global singleton  
662 Alu and L1 insertions. We utilized VCFtools v0.1.17 [101] to extract all autosomal SVs  
663 with no missing data (--max-missing 1) and an allele count of 1 across all samples (--  
664 non-ref-ac 1 --max-non-ref-ac 1), i.e. global singletons. From these, we extracted Alu-  
665 and L1-specific insertions using BCFtools v1.10.2 [102] to keep entries annotated with  
666 the 'SVTYPE="LINE1"' and 'SVTYPE="ALU"' tags. We calculated the enrichment of  
667 various genomic features overlapping L1/Alu singletons using the annotatePeaks script

668 within the Hypergeometric Optimization of Motif EnRichment (HOMER) v4.11.1 software  
669 suite [103] using hg38 v6.4 annotations. Finally, VCF files containing global singleton L1  
670 or Alu insertions were converted to PLINK BED format using PLINK v1.90b6.17 [104].  
671 We note that SVs on sex chromosomes were not included in any part of the analysis  
672 since (1) the necessary data was of poor quality, with over 80% of the X chromosome  
673 not genotyped by the 1000 Genomes Project, (2) the necessary data from the Y  
674 chromosome was not available, to our knowledge, and (3) the analyses require  
675 additional analytical considerations (i.e. male hemizyosity and X chromosome  
676 inactivation) that warrant their own study design [93-96].

677

678 Secondly, we prepared polymorphic SVs for inclusion in the GWAS analysis.  
679 VCFtools was used to isolate SVs with the following properties in each individual super-  
680 population: possessed a minimum and maximum of two alleles (biallelic), possessed a  
681 minor allele frequency (MAF) of at least 1%, passed Hardy-Weinberg equilibrium  
682 thresholding at  $p < 1e-6$ , had no missing samples, and was located on an autosome. To  
683 focus on shared, trans-ethnic sources of genetic variation, we used BCFtools to identify  
684 and subset SVs that were shared across all 5 super-populations.

685

686 Third, we prepared polymorphic SNVs for inclusion in the GWAS analysis. All  
687 SNVs were first annotated with rsIDs from dbSNP build 155 using BCFtools. Within  
688 each super-population, VCFtools was used to remove indels and keep autosomal SNVs  
689 with the same parameters as the polymorphic SVs. We note that for similar reasons as  
690 with the polymorphic SVs, sex chromosome SNVs were also omitted from all analyses.  
691 We then used BCFtools to identify and subset SNVs that were shared across all 5  
692 super-populations. Finally, we used BCFtools to generate the final genotype matrices by  
693 combining shared, polymorphic SNVs with shared, polymorphic SVs. VCF files  
694 containing the combined SNVs and SVs were then converted to PLINK BED format  
695 using PLINK, keeping the allele order. PLINK was also used to prune the combined  
696 SNV and SV matrices (--indep-pairwise 50 10 0.1) and to generate principal  
697 components (PCs) from the pruned genotypes, for inclusion as covariates in the GWAS.

698



### 699 **4.3 Super-population-specific and trans-ethnic GWAS**

700 We began by running GWAS within each super-population using PLINK  
701 v1.90b6.17 [104]. For the phenotype, we added the number of Alu and L1 global  
702 singleton insertions for each sample and segregated samples into cases and controls,  
703 depending on whether they contained or did not contain a global singleton. We ran  
704 GWAS analyses using a logistic regression model that included the following covariates:  
705 biological sex and the top 4 principal components from the pruned SNV and SV  
706 genotype matrices. Individual results from each super-population were combined via  
707 meta-analysis using PLINK. To help call significant variants, we generated a null  
708 distribution of p-values for each super-population by running 20 instances of the GWAS  
709 where the case-control status for each sample was randomly shuffled with the case-  
710 control status of a different sample. Each set of permutation results was meta-analyzed  
711 across super-populations to similarly obtain 20 random distributions of meta-analysis p-  
712 values. For the meta-analysis, we focused on the p-values and odds ratios generated  
713 using a random effects statistical model, as opposed to a fixed effects model, since 1)  
714 there may be heterogeneity across super-populations in response to different genetic  
715 variants, and 2) we were interested in enhancing the generalizability of our findings to  
716 facilitate future follow-up studies.

717  
718 To limit false positives, the Benjamini-Hochberg (BH) false discovery rate (FDR)  
719 was calculated in each analysis, and we used the p-value corresponding to a BH FDR <  
720 0.05 as the threshold for GWAS significance. As a secondary threshold, we used the  
721 permutation data to identify p-values corresponding to an average empirical FDR <  
722 0.05. To note, we calculated the average empirical FDR at a given p-value  $p_i$  by (i)  
723 counting the total number of null points with  $p \leq p_i$ , (ii) dividing by the number of  
724 permutations, to obtain an average number of null points with  $p \leq p_i$ , and (iii) dividing the  
725 average number of null points with  $p \leq p_i$  by the number of real points with  $p \leq p_i$ . GWAS  
726 variants were considered significant if they passed the stricter of the two thresholds in  
727 each analysis.

728

### 729 **4.4 Annotation of variants and annotation enrichment analyses**

730 We obtained BED files containing annotated genomic regions from various  
731 sources. We obtained the ENCODE blacklist v2 [52] for hg38 from  
732 <https://github.com/Boyle-Lab/Blacklist/tree/master/lists>. This curated “blacklist”  
733 represents a set of genomic regions with anomalous signal in next-generation  
734 sequencing experiments independent of cell type and individual experiment [52]. SNVs  
735 overlapping these regions were significantly enriched in our initial, unfiltered significant  
736 SNV list compared to the background SNV list (Fisher’s exact test, FDR = 1.33E-302).  
737 Segmental duplications [66, 67] and RepeatMasker annotations using the Repbase  
738 library [105] were obtained from the UCSC Genome Browser [106]. We obtained SV  
739 hotspot coordinates on hg19 from [65], and we used the online UCSC LiftOver tool to  
740 map coordinates to the hg38 genome assembly using the default settings. The BED  
741 tracks for full-length and intact L1s, only ORF2-intact L1s, and full-length non-intact L1s  
742 were obtained from L1Base v2—a dedicated database of putatively active L1 insertions  
743 [107]. We obtained the Registry (version 4) of candidate cis-Regulatory Elements  
744 (cCREs) for hg38 from the Search Candidate Regulatory Elements by  
745 ENCODE (SCREEN) web interface [55] (<http://screen-beta.wenglab.org>). We used the  
746 ‘intersect’ command in BEDTools v2.31.1 [108] to assign genomic region annotations to  
747 all overlapping variants used in this study. We were also interested in assessing  
748 whether variants were linked to specific regulatory gene annotations. All variants used  
749 in the study were submitted to the GREAT v4.0.4 [109, 110] online platform with the  
750 default settings (basal plus extension, proximal with 5 kb upstream and 1 kb  
751 downstream, plus distal up to 1000 kb) to assign gene annotations to each variant.  
752 These gene annotations were then used to assess the number of variants linked to  
753 genes in several TE regulatory lists—including a list of genes that regulated L1  
754 expression in a CRISPR screen using cancer cells [37], a list of genes that regulated L1  
755 transposition in an independent CRISPR screen using cancer cells [33], a list of  
756 candidate genes influencing intronic, intergenic, or exonic L1 RNA levels in  
757 lymphoblastoid cell lines [41], a list of genes with histone methyltransferase activity  
758 (GO:0042054), and a list of genes with RNA modification activity (GO:0009451). The  
759 two GO gene sets were obtained on 2024-07-29 from the Molecular Signatures  
760 Database (MSigDB) v2023.2.Hs [111, 112], corresponding to GO release 2023-07-27.

761

762           Given the limited number of significant SVs and the unavailability of SV  
763 sequences, we largely focused on blacklist-filtered, significant SNVs in downstream  
764 enrichment analyses. For each of the above annotations, we assessed whether  
765 greenlisted SNVs were significantly enriched or depleted for that annotation compared  
766 to background SNVs—all SNVs that were tested in the GWAS. The numbers of  
767 background and greenlisted SNVs overlapping or not overlapping a set of annotations  
768 were placed into a contingency table, and statistical significance was assessed using  
769 Fisher's exact test (with the `fisher.test` function in R v4.3.3). After all tests were carried  
770 out, p-values were FDR corrected using the `p.adjust` function in R. All  
771 enrichments/depletions with an FDR < 0.05 were considered significant. For the repeat  
772 enrichment analyses, we also analyzed our greenlisted SNVs using the TEENA web  
773 server [61] (on August 8, 2024), specifying the hg38 assembly and using all other  
774 default options.

775

#### 776 **4.5 RNA-seq and gene co-expression network analyses**

777           For lymphoblastoid cell line transcriptional analyses, mRNA-sequencing was  
778 initially carried out by the GEUVADIS consortium [69] on LCLs from a small subset of  
779 European and African (Yoruban, specifically) samples from the 1000 Genomes Project.  
780 Recently, we re-processed this data to quantify gene and transposon expression levels  
781 [41]. Briefly, reads were trimmed using `fastp v0.20.1` [113], trimmed reads were aligned  
782 to the GRCh38 human genome assembly using `STAR v2.7.3a` [114], and the  
783 `TEtranscripts v2.1.4` [115] package was used to obtain gene and TE counts, using the  
784 GENCODE release 33 [116] annotations and a repeat GTF file provided on the  
785 Hammell lab website. To note, the EBV genome (GenBank ID V01555.2) was included  
786 as an additional contig in our reference genome, since LCLs are generated by infecting  
787 B-cells with Epstein-Barr virus (EBV).

788

789           Using these gene/TE count matrices, lowly expressed genes were filtered out if  
790 50% of European or Yoruban samples did not have over 0.44 counts per million (cpm)  
791 or 0.43 cpm, respectively, which correspond to 10 reads in each cohort's median-length

792 library. Since we were interested in building consensus co-expression networks  
793 between the European and Yoruban samples, we also removed genes that were not  
794 expressed in both groups. After, the filtered counts underwent a variance stabilizing  
795 transformation (vst) using DESeq2 v1.42.1 [117] and the following covariates were  
796 regressed out using the 'removeBatchEffect' function in Limma v3.58.1 [118]: biological  
797 sex, sequencing lab, population category, principal components 1-2 of the pruned  
798 genotype matrices containing both SNVs and SVs, and EBV expression levels. The  
799 population category variable was omitted in the Yoruban batch correction since that did  
800 not vary.

801

802 The batch-corrected VST matrices were then used to perform weighted gene co-  
803 expression network analysis (WGCNA) [70] using the WGCNA v1.72-5 R package. We  
804 used the 'blockwiseConsensusModules' function to automate consensus network  
805 construction for both the European and African expression data, specifying these  
806 parameters: corType = "bicor", power = 12, networkType = "signed", maxPOutliers =  
807 0.05, mergeCutHeight = 0.25, deepSplit = 2, minKMEtoStay = 0, pamRespectsDendro =  
808 FALSE, minModuleSize = 30, and randomSeed = 90280. Phenotype-module  
809 correlations, and the corresponding p-values, were calculated using WGCNA's 'cor' and  
810 'corPvalueFisher' functions, respectively. The p-values for the European and Yoruban  
811 correlations were meta-analyzed using Fisher's method. For visualization purposes  
812 only, to show a correlation direction in the meta-analysis, we took the average of the  
813 European and Yoruban correlations. Modules showing opposite correlations across the  
814 two consensus networks were disregarded in the meta-analysis. Correlations with a p-  
815 value < 0.01 were considered significant.

816

#### 817 **4.6 Functional enrichment analyses**

818 We used the over-representation analysis (ORA) paradigm as implemented in  
819 the R package clusterProfiler v4.10.1 [119]. Gene Ontology Biological Process gene  
820 sets were obtained from the R package msigdb v7.5.1, an Ensembl ID-mapped  
821 collection of gene sets from the Molecular Signatures Database [111, 112]. For ORA  
822 with genes linked to greenlisted SNVs and SVs, we used the genes linked to the

823 background SNVs and SVs, respectively, for the universe background to compute  
824 enrichment significance. For ORA analysis of co-expression network modules, we used  
825 all genes in the network for the universe background. All gene sets with an FDR < 0.05  
826 were considered significant, and the top 10 significant gene sets, at most, were plotted.  
827 All enrichments results are included in **Supplementary Table S1D, S1E, and S1J**.  
828  
829

830 **DECLARATIONS**

831 **Ethics approval and consent to participate**

832 Not applicable.

833

834 **Consent for publication**

835 Not applicable.

836

837 **Availability of data and materials**

838 All code is available on the Benayoun lab GitHub  
839 ([https://github.com/BenayounLaboratory/TE\\_GWAS](https://github.com/BenayounLaboratory/TE_GWAS)). Analyses were conducted using R  
840 version 4.3.3. Code was re-run independently on R version 4.3.0 to check for  
841 reproducibility.

842

843 **Competing interests**

844 The authors declare that they have no competing interests.

845

846 **Funding**

847 This work was supported by the National Science Foundation  
848 [<https://www.nsf.gov/>] Graduate Research Fellowship Program (NSF GRFP) DGE-  
849 1842487 (J.I.B.), the National Institute on Aging [<https://www.nia.nih.gov/>] T32  
850 AG052374 (J.I.B.), the University of Southern California with a Provost Fellowship  
851 (J.I.B.), and the National Institute of General Medical Sciences  
852 [<https://www.nigms.nih.gov/>] R35 GM142395 (to B.A.B).

853 The funders had no role in study design, data collection and analysis, decision to  
854 publish, or preparation of the manuscript.

855

856 **Authors' contributions**

857 **Juan I Bravo:** Conceptualization, formal analysis, investigation, methodology,  
858 visualization, writing - original draft preparation, writing - review and editing. **Lucia**  
859 **Zhang:** Validation, formal analysis, writing - review and editing. **B er enice A.**

860 **Benayoun:** Conceptualization, formal analysis, funding acquisition, methodology,  
861 supervision, visualization, writing - original draft preparation, writing - review and editing.

862

### 863 **Acknowledgements**

864 We would like to thank Prof. Rachel Brem at the University of California Berkeley  
865 for her feedback and insights on the GWAS analyses. We would also like to thank Dr.  
866 Heather C. Mefford at St. Jude Children’s Research Hospital for referring us to  
867 publications that helped shape the analyses linking our variants to regions potentially  
868 prone to genomic instability and involved in disease. Finally, we would like to thank Dr.  
869 Minhoo Kim, Mr. Aaron Lemus, and Ms. Eyael Habtemariam for their comments on the  
870 manuscript.

871

## 872 REFERENCES

- 873 1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K,  
874 Doyle M, FitzHugh W, et al: **Initial sequencing and analysis of the human genome.**  
875 *Nature* 2001, **409**:860-921.
- 876 2. Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, Smith HO, Yandell M,  
877 Evans CA, Holt RA, et al: **The Sequence of the Human Genome.** *Science* 2001, **291**:1304-  
878 1351.
- 879 3. Ostertag EM, Kazazian HH, Jr.: **Biology of mammalian L1 retrotransposons.** *Annu Rev*  
880 *Genet* 2001, **35**:501-538.
- 881 4. Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran JV, Kazazian HH: **Hot**  
882 **L1s account for the bulk of retrotransposition in the human population.** *Proceedings of*  
883 *the National Academy of Sciences* 2003, **100**:5280-5285.
- 884 5. Rubin CM, Houck CM, Deininger PL, Friedmann T, Schmid CW: **Partial nucleotide**  
885 **sequence of the 300-nucleotide interspersed repeated human DNA sequences.** *Nature*  
886 1980, **284**:372-374.
- 887 6. Dewannieux M, Esnault C, Heidmann T: **LINE-mediated retrotransposition of marked**  
888 **Alu sequences.** *Nature Genetics* 2003, **35**:41-48.
- 889 7. Wallace N, Wagstaff BJ, Deininger PL, Roy-Engel AM: **LINE-1 ORF1 protein enhances Alu**  
890 **SINE retrotransposition.** *Gene* 2008, **419**:1-6.
- 891 8. Bennett EA, Keller H, Mills RE, Schmidt S, Moran JV, Weichenrieder O, Devine SE: **Active**  
892 **Alu retrotransposons in the human genome.** *Genome Research* 2008, **18**:1875-1883.
- 893 9. Wei W, Gilbert N, Ooi SL, Lawler JF, Ostertag EM, Kazazian HH, Boeke JD, Moran JV:  
894 **Human L1 Retrotransposition: cisPreference versus trans Complementation.** *Molecular*  
895 *and Cellular Biology* 2001, **21**:1429-1439.
- 896 10. Mendez-Dorantes C, Bhargava R, Stark JM: **Repeat-mediated deletions can be induced**  
897 **by a chromosomal break far from a repeat, but multiple pathways suppress such**  
898 **rearrangements.** *Genes Dev* 2018, **32**:524-536.
- 899 11. Mendez-Dorantes C, Tsai LJ, Jahanshir E, Lopezcolorado FW, Stark JM: **BLM has Contrary**  
900 **Effects on Repeat-Mediated Deletions, based on the Distance of DNA DSBs to a Repeat**  
901 **and Repeat Divergence.** *Cell Reports* 2020, **30**:1342-1357.e1344.
- 902 12. Robberecht C, Voet T, Esteki MZ, Nowakowska BA, Vermeesch JR: **Nonallelic**  
903 **homologous recombination between retrotransposable elements is a driver of de**  
904 **novo unbalanced translocations.** *Genome Research* 2013, **23**:411-418.
- 905 13. Boone Philip M, Yuan B, Campbell Ian M, Scull Jennifer C, Withers Marjorie A, Baggett  
906 Brett C, Beck Christine R, Shaw Christine J, Stankiewicz P, Moretti P, et al: **The**  
907 **<em>Alu</em>-Rich Genomic Architecture of <em>SPAST</em> Predisposes to**  
908 **Diverse and Functionally Distinct Disease-Associated CNV Alleles.** *The American*  
909 *Journal of Human Genetics* 2014, **95**:143-161.
- 910 14. Campbell IM, Gambin T, Dittwald P, Beck CR, Shuvarikov A, Hixson P, Patel A, Gambin A,  
911 Shaw CA, Rosenfeld JA, Stankiewicz P: **Human endogenous retroviral elements**  
912 **promote genome instability via non-allelic homologous recombination.** *BMC Biology*  
913 2014, **12**:74.



- 914 15. Startek M, Szafranski P, Gambin T, Campbell IM, Hixson P, Shaw CA, Stankiewicz P,  
915 Gambin A: **Genome-wide analyses of LINE–LINE-mediated nonallelic homologous**  
916 **recombination**. *Nucleic Acids Research* 2015, **43**:2188-2198.
- 917 16. Belancio VP, Deininger PL, Roy-Engel AM: **LINE dancing in the human genome:**  
918 **transposable elements and disease**. *Genome Medicine* 2009, **1**:97.
- 919 17. Cordaux R, Batzer MA: **The impact of retrotransposons on human genome evolution**.  
920 *Nature Reviews Genetics* 2009, **10**:691-703.
- 921 18. Kolomietz E, Meyn MS, Pandita A, Squire JA: **The role of Alu repeat clusters as**  
922 **mediators of recurrent chromosomal aberrations in tumors**. *Genes, Chromosomes and*  
923 *Cancer* 2002, **35**:97-112.
- 924 19. Bravo JI, Nozownik S, Danthi PS, Benayoun BA: **Transposable elements, circular RNAs**  
925 **and mitochondrial transcription in age-related genomic regulation**. *Development* 2020,  
926 **147**.
- 927 20. Zhang X, Zhang R, Yu J: **New Understanding of the Relevant Role of LINE-1**  
928 **Retrotransposition in Human Disease and Immune Modulation**. *Frontiers in Cell and*  
929 *Developmental Biology* 2020, **8**.
- 930 21. López-Otín C, Blasco MA, Partridge L, Serrano M, Kroemer G: **Hallmarks of aging: An**  
931 **expanding universe**. *Cell* 2023, **186**:243-278.
- 932 22. De Cecco M, Criscione SW, Peterson AL, Neretti N, Sedivy JM, Kreiling JA: **Transposable**  
933 **elements become active and mobile in the genomes of aging mammalian somatic**  
934 **tissues**. *Aging (Albany NY)* 2013, **5**:867-883.
- 935 23. De Cecco M, Criscione SW, Peckham EJ, Hillenmeyer S, Hamm EA, Manivannan J,  
936 Peterson AL, Kreiling JA, Neretti N, Sedivy JM: **Genomes of replicatively senescent cells**  
937 **undergo global epigenetic changes leading to gene silencing and activation of**  
938 **transposable elements**. *Aging Cell* 2013, **12**:247-256.
- 939 24. Campisi J: **Aging, Cellular Senescence, and Cancer**. *Annual Review of Physiology* 2013,  
940 **75**:685-705.
- 941 25. Franceschi C, Garagnani P, Parini P, Giuliani C, Santoro A: **Inflammaging: a new**  
942 **immune–metabolic viewpoint for age-related diseases**. *Nature Reviews Endocrinology*  
943 2018, **14**:576-590.
- 944 26. Wallace NA, Belancio VP, Deininger PL: **L1 mobile element expression causes multiple**  
945 **types of toxicity**. *Gene* 2008, **419**:75-81.
- 946 27. Belancio VP, Roy-Engel AM, Pochampally RR, Deininger P: **Somatic expression of LINE-1**  
947 **elements in human tissues**. *Nucleic acids research* 2010, **38**:3909-3922.
- 948 28. De Cecco M, Ito T, Petrashen AP, Elias AE, Skvir NJ, Criscione SW, Caligiana A, Broccoli G,  
949 Adney EM, Boeke JD, et al: **L1 drives IFN in senescent cells and promotes age-**  
950 **associated inflammation**. *Nature* 2019, **566**:73-78.
- 951 29. Yamada K, Kaneko H, Shimizu H, Suzumura A, Namba R, Takayama K, Ito S, Sugimoto M,  
952 Terasaki H: **Lamivudine Inhibits Alu RNA-induced Retinal Pigment Epithelium**  
953 **Degeneration via Anti-inflammatory and Anti-senescence Activities**. *Translational*  
954 *Vision Science & Technology* 2020, **9**:1-1.
- 955 30. Wang J, Geesman GJ, Hostikka SL, Atallah M, Blackwell B, Lee E, Cook PJ, Pasaniuc B,  
956 Shariat G, Halperin E, et al: **Inhibition of activated pericentromeric SINE/Alu repeat**

- 957 transcription in senescent human adult stem cells reinstates self-renewal. *Cell Cycle*  
958 2011, **10**:3016-3030.
- 959 31. Levin HL, Moran JV: **Dynamic interactions between transposable elements and their**  
960 **hosts.** *Nature Reviews Genetics* 2011, **12**:615-627.
- 961 32. Rebollo R, Romanish MT, Mager DL: **Transposable Elements: An Abundant and Natural**  
962 **Source of Regulatory Sequences for Host Genes.** *Annual Review of Genetics* 2012,  
963 **46**:21-42.
- 964 33. Liu N, Lee CH, Swigut T, Grow E, Gu B, Bassik MC, Wysocka J: **Selective silencing of**  
965 **euchromatic L1s revealed by genome-wide screens for L1 regulators.** *Nature* 2018,  
966 **553**:228-232.
- 967 34. Tristan-Ramos P, Morell S, Sanchez L, Toledo B, Garcia-Perez JL, Heras SR: **sRNA/L1**  
968 **retrotransposition: using siRNAs and miRNAs to expand the applications of the cell**  
969 **culture-based LINE-1 retrotransposition assay.** *Philosophical Transactions of the Royal*  
970 *Society B: Biological Sciences* 2020, **375**:20190346.
- 971 35. Mita P, Sun X, Fenyö D, Kahler DJ, Li D, Agmon N, Wudzinska A, Keegan S, Bader JS, Yun  
972 C, Boeke JD: **BRCA1 and S phase DNA repair pathways restrict LINE-1**  
973 **retrotransposition in human cells.** *Nat Struct Mol Biol* 2020, **27**:179-191.
- 974 36. Briggs EM, Mita P, Sun X, Ha S, Vasilyev N, Leopold ZR, Nudler E, Boeke JD, Logan SK:  
975 **Unbiased proteomic mapping of the LINE-1 promoter using CRISPR Cas9.** *Mobile DNA*  
976 2021, **12**:21.
- 977 37. Li X, Bie L, Wang Y, Hong Y, Zhou Z, Fan Y, Yan X, Tao Y, Huang C, Zhang Y, et al: **LINE-1**  
978 **transcription activates long-range gene expression.** *Nature Genetics* 2024, **56**:1494-  
979 1502.
- 980 38. Sun X, Wang X, Tang Z, Grivainis M, Kahler D, Yun C, Mita P, Fenyö D, Boeke JD:  
981 **Transcription factor profiling reveals molecular choreography and key regulators of**  
982 **human retrotransposon expression.** *Proc Natl Acad Sci U S A* 2018, **115**:E5526-e5535.
- 983 39. Chung N, Jonaid GM, Quinton S, Ross A, Sexton CE, Alberto A, Clymer C, Churchill D,  
984 Navarro Leija O, Han MV: **Transcriptome analyses of tumor-adjacent somatic tissues**  
985 **reveal genes co-expressed with transposable elements.** *Mobile DNA* 2019, **10**:39.
- 986 40. Tristán-Ramos P, Rubio-Roldan A, Peris G, Sánchez L, Amador-Cubero S, Viollet S,  
987 Cristofari G, Heras SR: **The tumor suppressor microRNA let-7 inhibits human LINE-1**  
988 **retrotransposition.** *Nature Communications* 2020, **11**:5712.
- 989 41. Bravo JI, Mizrahi CR, Kim S, Zhang L, Suh Y, Benayoun BA: **An eQTL-based approach**  
990 **reveals candidate regulators of LINE-1 RNA levels in lymphoblastoid cells.** *PLOS*  
991 *Genetics* 2024, **20**:e1011311.
- 992 42. Johnston HR, Hu Y, Cutler DJ: **Population Genetics Identifies Challenges in Analyzing**  
993 **Rare Variants.** *Genetic Epidemiology* 2015, **39**:145-148.
- 994 43. Baduel P, Leduque B, Ignace A, Gy I, Gil J, Loudet O, Colot V, Quadrana L: **Genetic and**  
995 **environmental modulation of transposition shapes the evolutionary potential of**  
996 ***Arabidopsis thaliana*.** *Genome Biology* 2021, **22**:138.
- 997 44. Evangelou E, Ioannidis JPA: **Meta-analysis methods for genome-wide association**  
998 **studies and beyond.** *Nature Reviews Genetics* 2013, **14**:379-389.
- 999 45. Peterson RE, Kuchenbaecker K, Walters RK, Chen CY, Popejoy AB, Periyasamy S, Lam M,  
1000 Iyegbe C, Strawbridge RJ, Brick L, et al: **Genome-wide Association Studies in Ancestrally**

- 1001 **Diverse Populations: Opportunities, Methods, Pitfalls, and Recommendations.** *Cell*  
1002 2019, **179**:589-603.
- 1003 46. Sultana T, van Essen D, Siol O, Bailly-Bechet M, Philippe C, Zine El Aabidine A, Pioger L,  
1004 Nigumann P, Sacconi S, Andrau J-C, et al: **The Landscape of L1 Retrotransposons in the**  
1005 **Human Genome Is Shaped by Pre-insertion Sequence Biases and Post-insertion**  
1006 **Selection.** *Molecular Cell* 2019, **74**:555-570.e557.
- 1007 47. Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera AV, Lowther C,  
1008 Gauthier LD, Wang H, et al: **A structural variation reference for medical and population**  
1009 **genetics.** *Nature* 2020, **581**:444-451.
- 1010 48. Pereira L, Mutesa L, Tindana P, Ramsay M: **African genetic diversity and adaptation**  
1011 **inform a precision medicine agenda.** *Nature Reviews Genetics* 2021, **22**:284-306.
- 1012 49. Ewing AD, Kazazian HH: **Whole-genome resequencing allows detection of many rare**  
1013 **LINE-1 insertion alleles in humans.** *Genome Research* 2011, **21**:985-990.
- 1014 50. Wong JS, Jadhav T, Young E, Wang Y, Xiao M: **Characterization of full-length LINE-1**  
1015 **insertions in 154 genomes.** *Genomics* 2021, **113**:3804-3810.
- 1016 51. Kojima S, Koyama S, Ka M, Saito Y, Parrish EH, Endo M, Takata S, Mizukoshi M, Hikino K,  
1017 Takeda A, et al: **Mobile element variation contributes to population-specific genome**  
1018 **diversification, gene regulation and disease risk.** *Nature Genetics* 2023, **55**:939-951.
- 1019 52. Amemiya HM, Kundaje A, Boyle AP: **The ENCODE Blacklist: Identification of Problematic**  
1020 **Regions of the Genome.** *Scientific Reports* 2019, **9**:9354.
- 1021 53. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM: **A**  
1022 **program for annotating and predicting the effects of single nucleotide polymorphisms,**  
1023 **SnPEff.** *Fly* 2012, **6**:80-92.
- 1024 54. Hampel H, Hardy J, Blennow K, Chen C, Perry G, Kim SH, Villemagne VL, Aisen P,  
1025 Vendruscolo M, Iwatsubo T, et al: **The Amyloid- $\beta$  Pathway in Alzheimer's Disease.**  
1026 *Molecular Psychiatry* 2021, **26**:5481-5503.
- 1027 55. Abascal F, Acosta R, Addleman NJ, Adrian J, Afzal V, Ai R, Aken B, Akiyama JA, Jammal  
1028 OA, Amrhein H, et al: **Expanded encyclopaedias of DNA elements in the human and**  
1029 **mouse genomes.** *Nature* 2020, **583**:699-710.
- 1030 56. Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F,  
1031 Young N, et al: **The Genotype-Tissue Expression (GTEx) project.** *Nature Genetics* 2013,  
1032 **45**:580-585.
- 1033 57. Carithers LJ, Ardlie K, Barcus M, Branton PA, Britton A, Buia SA, Compton CC, DeLuca DS,  
1034 Peter-Demchok J, Gelfand ET, et al: **A Novel Approach to High-Quality Postmortem**  
1035 **Tissue Procurement: The GTEx Project.** *Biopreservation and Biobanking* 2015, **13**:311-  
1036 319.
- 1037 58. Consortium TG, Aguet F, Anand S, Ardlie KG, Gabriel S, Getz GA, Graubert A, Hadley K,  
1038 Handsaker RE, Huang KH, et al: **The GTEx Consortium atlas of genetic regulatory effects**  
1039 **across human tissues.** *Science* 2020, **369**:1318-1330.
- 1040 59. Faulkner GJ, Billon V: **L1 retrotransposition in the soma: a field jumping ahead.** *Mobile*  
1041 *DNA* 2018, **9**:22.
- 1042 60. Raghava Kurup R, Oakes EK, Manning AC, Mukherjee P, Vadlamani P, Hundley HA: **RNA**  
1043 **binding by ADAR3 inhibits adenosine-to-inosine editing and promotes expression of**  
1044 **immune response protein MAVS.** *Journal of Biological Chemistry* 2022, **298**.

- 1045 61. Li Y, Lyu R, Chen S, Wang Y, Sun M-a: **TEENA: an integrated web server for transposable**  
1046 **element enrichment analysis in various model and non-model organisms.** *Nucleic Acids*  
1047 *Research* 2024, **52**:W126-W131.
- 1048 62. Bailey JA, Liu G, Eichler EE: **An *Alu* Transposition Model for the Origin and**  
1049 **Expansion of Human Segmental Duplications.** *The American Journal of Human Genetics*  
1050 2003, **73**:823-834.
- 1051 63. Sharp AJ, Locke DP, McGrath SD, Cheng Z, Bailey JA, Vallente RU, Pertz LM, Clark RA,  
1052 Schwartz S, Segraves R, et al: **Segmental Duplications and Copy-Number Variation in**  
1053 **the Human Genome.** *The American Journal of Human Genetics* 2005, **77**:78-88.
- 1054 64. Perry GH, Tchinda J, McGrath SD, Zhang J, Picker SR, Cáceres AM, Iafrate AJ, Tyler-Smith  
1055 C, Scherer SW, Eichler EE, et al: **Hotspots for copy number variation in chimpanzees**  
1056 **and humans.** *Proc Natl Acad Sci U S A* 2006, **103**:8006-8011.
- 1057 65. Lin Y-L, Gokcumen O: **Fine-Scale Characterization of Genomic Structural Variation in**  
1058 **the Human Genome Reveals Adaptive and Biomedically Relevant Hotspots.** *Genome*  
1059 *Biology and Evolution* 2019, **11**:1136-1151.
- 1060 66. Bailey JA, Yavor AM, Massa HF, Trask BJ, Eichler EE: **Segmental Duplications:**  
1061 **Organization and Impact Within the Current Human Genome Project Assembly.**  
1062 *Genome Research* 2001, **11**:1005-1017.
- 1063 67. Bailey JA, Gu Z, Clark RA, Reinert K, Samonte RV, Schwartz S, Adams MD, Myers EW, Li  
1064 PW, Eichler EE: **Recent Segmental Duplications in the Human Genome.** *Science* 2002,  
1065 **297**:1003-1007.
- 1066 68. Porubsky D, Höps W, Ashraf H, Hsieh P, Rodriguez-Martin B, Yilmaz F, Ebler J, Hallast P,  
1067 Maria Maggiolini FA, Harvey WT, et al: **Recurrent inversion polymorphisms in humans**  
1068 **associate with genetic instability and genomic disorders.** *Cell* 2022, **185**:1986-  
1069 2005.e1926.
- 1070 69. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas MA,  
1071 González-Porta M, Kurbatova N, Griebel T, Ferreira PG, et al: **Transcriptome and**  
1072 **genome sequencing uncovers functional variation in humans.** *Nature* 2013, **501**:506-  
1073 511.
- 1074 70. Langfelder P, Horvath S: **WGCNA: an R package for weighted correlation network**  
1075 **analysis.** *BMC Bioinformatics* 2008, **9**:559.
- 1076 71. Kubo S, Seleme MdC, Soifer HS, Perez JLG, Moran JV, Kazazian HH, Kasahara N: **L1**  
1077 **retrotransposition in nondividing and primary human somatic cells.** *Proceedings of the*  
1078 *National Academy of Sciences* 2006, **103**:8036-8041.
- 1079 72. Macia A, Widmann TJ, Heras SR, Ayllon V, Sanchez L, Benkaddour-Boumzaouad M,  
1080 Muñoz-Lopez M, Rubio A, Amador-Cubero S, Blanco-Jimenez E, et al: **Engineered LINE-1**  
1081 **retrotransposition in nondividing human neurons.** *Genome Research* 2017, **27**:335-348.
- 1082 73. Mita P, Wudzinska A, Sun X, Andrade J, Nayak S, Kahler DJ, Badri S, LaCava J, Ueberheide  
1083 B, Yun CY, et al: **LINE-1 protein localization and functional dynamics during the cell**  
1084 **cycle.** *eLife* 2018, **7**:e30058.
- 1085 74. Potapova T, Gorbsky GJ: **The Consequences of Chromosome Segregation Errors in**  
1086 **Mitosis and Meiosis.** *Biology* 2017, **6**:12.

- 1087 75. Gu Z, Liu Y, Zhang Y, Cao H, Lyu J, Wang X, Wylie A, Newkirk SJ, Jones AE, Lee M, et al:  
1088 **Silencing of LINE-1 retrotransposons is a selective dependency of myeloid leukemia.**  
1089 *Nature Genetics* 2021, **53**:672-682.
- 1090 76. Müller I, Moroni AS, Shlyueva D, Sahadevan S, Schoof EM, Radzisheuskaya A, Højfeldt  
1091 JW, Tatar T, Koche RP, Huang C, Helin K: **MPP8 is essential for sustaining self-renewal**  
1092 **of ground-state pluripotent stem cells.** *Nature Communications* 2021, **12**:3034.
- 1093 77. Danac JMC, Matthews RE, Gungi A, Qin C, Parsons H, Antrobus R, Timms RT,  
1094 Tchasovnikarova IA: **Competition between two HUSH complexes orchestrates the**  
1095 **immune response to retroelement invasion.** *Molecular Cell* 2024, **84**:2870-2881.e2875.
- 1096 78. Dai L, Taylor MS, O'Donnell KA, Boeke JD: **Poly(A) binding protein C1 is essential for**  
1097 **efficient L1 retrotransposition and affects L1 RNP formation.** *Mol Cell Biol* 2012,  
1098 **32**:4323-4336.
- 1099 79. Taylor Martin S, LaCava J, Mita P, Molloy Kelly R, Huang Cheng Ran L, Li D, Adney  
1100 Emily M, Jiang H, Burns Kathleen H, Chait Brian T, et al: **Affinity Proteomics Reveals**  
1101 **Human Host Factors Implicated in Discrete Stages of LINE-1 Retrotransposition.** *Cell*  
1102 2013, **155**:1034-1048.
- 1103 80. Goodier JL, Cheung LE, Kazazian HH, Jr: **Mapping the LINE1 ORF1 protein interactome**  
1104 **reveals associated inhibitors of human retrotransposition.** *Nucleic Acids Research*  
1105 2013, **41**:7401-7419.
- 1106 81. Mazoyer S: **Genomic rearrangements in the BRCA1 and BRCA2 genes.** *Human Mutation*  
1107 2005, **25**:415-422.
- 1108 82. Peixoto A, Pinheiro M, Massena L, Santos C, Pinto P, Rocha P, Pinto C, Teixeira MR:  
1109 **Genomic characterization of two large Alu-mediated rearrangements of the BRCA1**  
1110 **gene.** *Journal of Human Genetics* 2013, **58**:78-83.
- 1111 83. Wang Y, Bernhardt AJ, Nacson J, Kraiss JJ, Tan Y-F, Nicolas E, Radke MR, Handorf E, Llop-  
1112 Guevara A, Balmaña J, et al: **BRCA1 intronic Alu elements drive gene rearrangements**  
1113 **and PARP inhibitor resistance.** *Nature Communications* 2019, **10**:5661.
- 1114 84. Hanahan D, Weinberg Robert A: **Hallmarks of Cancer: The Next Generation.** *Cell* 2011,  
1115 **144**:646-674.
- 1116 85. Rodić N, Sharma R, Sharma R, Zampella J, Dai L, Taylor MS, Hruban RH, Iacobuzio-  
1117 Donahue CA, Maitra A, Torbenson MS, et al: **Long Interspersed Element-1 Protein**  
1118 **Expression Is a Hallmark of Many Human Cancers.** *The American Journal of Pathology*  
1119 2014, **184**:1280-1286.
- 1120 86. Evering TH, Marston JL, Gan L, Nixon DF: **Transposable elements and Alzheimer's**  
1121 **disease pathogenesis.** *Trends Neurosci* 2023, **46**:170-172.
- 1122 87. Pascarella G, Hon CC, Hashimoto K, Busch A, Luginbühl J, Parr C, Hin Yip W, Abe K, Kratz  
1123 A, Bonetti A, et al: **Recombination of repeat elements generates somatic complexity in**  
1124 **human genomes.** *Cell* 2022, **185**:3025-3040.e3026.
- 1125 88. Buiting K, Nazlican H, Galetzka D, Wawrzik M, Groß S, Horsthemke B: **C15orf2 and a**  
1126 **novel noncoding transcript from the Prader-Willi/Angelman syndrome region show**  
1127 **monoallelic expression in fetal brain.** *Genomics* 2007, **89**:588-595.
- 1128 89. Angulo MA, Butler MG, Cataletto ME: **Prader-Willi syndrome: a review of clinical,**  
1129 **genetic, and endocrine findings.** *J Endocrinol Invest* 2015, **38**:1249-1263.

- 1130 90. Makoff AJ, Flomen RH: **Detailed analysis of 15q11-q14 sequence corrects errors and**  
1131 **gaps in the public access sequence to fully reveal large segmental duplications at**  
1132 **breakpoints for Prader-Willi, Angelman, and inv dup(15) syndromes.** *Genome Biology*  
1133 2007, **8**:R114.
- 1134 91. Mefford HC, Eichler EE: **Duplication hotspots, rare genomic disorders, and common**  
1135 **disease.** *Current Opinion in Genetics & Development* 2009, **19**:196-204.
- 1136 92. Byrska-Bishop M, Evani US, Zhao X, Basile AO, Abel HJ, Regier AA, Corvelo A, Clarke WE,  
1137 Musunuri R, Nagulapalli K, et al: **High-coverage whole-genome sequencing of the**  
1138 **expanded 1000 Genomes Project cohort including 602 trios.** *Cell* 2022, **185**:3426-  
1139 3440.e3419.
- 1140 93. Gao F, Chang D, Biddanda A, Ma L, Guo Y, Zhou Z, Keinan A: **XWAS: A Software Toolset**  
1141 **for Genetic Data Analysis and Association Studies of the X Chromosome.** *J Hered* 2015,  
1142 **106**:666-671.
- 1143 94. Keur N, Ricaño-Ponce I, Kumar V, Matzaraki V: **A systematic review of analytical**  
1144 **methods used in genetic association analysis of the X-chromosome.** *Briefings in*  
1145 *Bioinformatics* 2022, **23**.
- 1146 95. Belloy ME, Le Guen Y, Stewart I, Williams K, Herz J, Sherva R, Zhang R, Merritt V,  
1147 Panizzon MS, Hauger RL, et al: **Role of the X Chromosome in Alzheimer Disease**  
1148 **Genetics.** *JAMA Neurology* 2024, **81**:1032-1042.
- 1149 96. Simmonds E, Leonenko G, Yaman U, Bellou E, Myers A, Morgan K, Brookes K, Hardy J,  
1150 Salih D, Escott-Price V: **Chromosome X-wide association study in case control studies of**  
1151 **pathologically confirmed Alzheimer's disease in a European population.** *Translational*  
1152 *Psychiatry* 2024, **14**:358.
- 1153 97. Rangwala SH, Kazazian HH: **The L1 retrotransposition assay: A retrospective and**  
1154 **toolkit.** *Methods* 2009, **49**:219-226.
- 1155 98. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A,  
1156 Clark AG, Donnelly P, Eichler EE, et al: **A global reference for human genetic variation.**  
1157 *Nature* 2015, **526**:68-74.
- 1158 99. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye  
1159 K, Jun G, Hsi-Yang Fritz M, et al: **An integrated map of structural variation in 2,504**  
1160 **human genomes.** *Nature* 2015, **526**:75-81.
- 1161 100. Lowy-Gallego E, Fairley S, Zheng-Bradley X, Ruffier M, Clarke L, Flicek P, null n: **Variant**  
1162 **calling on the GRCh38 assembly with the data from phase three of the 1000 Genomes**  
1163 **Project [version 2; peer review: 2 approved].** *Wellcome Open Research* 2019, **4**.
- 1164 101. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter  
1165 G, Marth GT, Sherry ST, et al: **The variant call format and VCFtools.** *Bioinformatics* 2011,  
1166 **27**:2156-2158.
- 1167 102. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,  
1168 McCarthy SA, Davies RM, Li H: **Twelve years of SAMtools and BCFtools.** *GigaScience*  
1169 2021, **10**.
- 1170 103. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H,  
1171 Glass CK: **Simple combinations of lineage-determining transcription factors prime cis-**  
1172 **regulatory elements required for macrophage and B cell identities.** *Mol Cell* 2010,  
1173 **38**:576-589.

- 1174 104. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar P,  
1175 de Bakker PIW, Daly MJ, Sham PC: **PLINK: A Tool Set for Whole-Genome Association**  
1176 **and Population-Based Linkage Analyses.** *The American Journal of Human Genetics*  
1177 2007, **81**:559-575.
- 1178 105. Jurka J: **Rebase Update: a database and an electronic journal of repetitive elements.**  
1179 *Trends in Genetics* 2000, **16**:418-420.
- 1180 106. Nassar LR, Barber GP, Benet-Pagès A, Casper J, Clawson H, Diekhans M, Fischer C,  
1181 Gonzalez JN, Hinrichs Angie S, Lee Brian T, et al: **The UCSC Genome Browser database:**  
1182 **2023 update.** *Nucleic Acids Research* 2022, **51**:D1188-D1195.
- 1183 107. Penzkofer T, Jäger M, Figlerowicz M, Badge R, Mundlos S, Robinson PN, Zemojtel T:  
1184 **L1Base 2: more retrotransposition-active LINE-1s, more mammalian genomes.** *Nucleic*  
1185 *Acids Research* 2016, **45**:D68-D73.
- 1186 108. Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic**  
1187 **features.** *Bioinformatics* 2010, **26**:841-842.
- 1188 109. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G:  
1189 **GREAT improves functional interpretation of cis-regulatory regions.** *Nature*  
1190 *Biotechnology* 2010, **28**:495-501.
- 1191 110. Tanigawa Y, Dyer ES, Bejerano G: **WhichTF is functionally important in your open**  
1192 **chromatin data?** *PLOS Computational Biology* 2022, **18**:e1010378.
- 1193 111. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,  
1194 Pomeroy SL, Golub TR, Lander ES, Mesirov JP: **Gene set enrichment analysis: A**  
1195 **knowledge-based approach for interpreting genome-wide expression profiles.**  
1196 *Proceedings of the National Academy of Sciences* 2005, **102**:15545-15550.
- 1197 112. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP:  
1198 **Molecular signatures database (MSigDB) 3.0.** *Bioinformatics* 2011, **27**:1739-1740.
- 1199 113. Chen S, Zhou Y, Chen Y, Gu J: **fastp: an ultra-fast all-in-one FASTQ preprocessor.**  
1200 *Bioinformatics* 2018, **34**:i884-i890.
- 1201 114. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,  
1202 Gingeras TR: **STAR: ultrafast universal RNA-seq aligner.** *Bioinformatics* 2012, **29**:15-21.
- 1203 115. Jin Y, Tam OH, Paniagua E, Hammell M: **TEtranscripts: a package for including**  
1204 **transposable elements in differential expression analysis of RNA-seq datasets.**  
1205 *Bioinformatics* 2015, **31**:3593-3599.
- 1206 116. Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, Mudge JM, Sisu  
1207 C, Wright J, Armstrong J, et al: **GENCODE reference annotation for the human and**  
1208 **mouse genomes.** *Nucleic Acids Research* 2018, **47**:D766-D773.
- 1209 117. Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for**  
1210 **RNA-seq data with DESeq2.** *Genome Biology* 2014, **15**:550.
- 1211 118. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK: **limma powers**  
1212 **differential expression analyses for RNA-sequencing and microarray studies.** *Nucleic*  
1213 *Acids Research* 2015, **43**:e47-e47.
- 1214 119. Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, Feng T, Zhou L, Tang W, Zhan L, et al:  
1215 **clusterProfiler 4.0: A universal enrichment tool for interpreting omics data.** *The*  
1216 *Innovation* 2021, **2**:100141.
- 1217

1218



1219 **Legends to Figures**

1220

1221 **Fig 1. Overview of the pipeline to scan for genetic variants associated with L1/Alu**  
1222 **global singletons.**

1223 **(A)** An illustration of the samples used in this study. SNV and SV genetic data was  
1224 available for 2503 individuals from 5 super-populations, including 660 Africans (AFR),  
1225 504 East Asians (EAS), 503 Europeans (EUR), 489 South Asians (SAS), and 347  
1226 Admixed Americans (AMR). Males and females were approximately equally  
1227 represented, with male-to-female ratios (M/F ratios) ranging from 0.91 to 1.14. **(B)** A  
1228 schematic illustrating the trans-ethnic integration of available SNV and SV data to  
1229 identify variants associated with L1/Alu insertion global singletons. Within each super-  
1230 population, samples were segregated into cases and controls depending on whether or  
1231 not they harbored a global Alu or L1 insertion singleton. GWAS was carried out within  
1232 each super-population to identify polymorphic SNVs and SVs associated with case-  
1233 control status. Finally, GWAS results from all 5 super-populations were meta-analyzed  
1234 using a random effects statistical model, yielding a summary meta-analysis odds ratio  
1235 and p-value for each variant. **(C)** The frequency of Alu and L1 insertion singletons in  
1236 each super-population (*left panel*) or among cases within each super-population (*middle*  
1237 *panel*). The distribution of insertion singletons across autosomes is also shown (*right*  
1238 *panel*). **(D)** A Manhattan plot for the trans-ethnic GWAS meta-analysis. The dashed line  
1239 at  $p = 1.40E-5$  corresponds to an average empirical FDR  $< 0.05$ , based on 20 random  
1240 permutations. One such permutation is shown in the bottom panel for illustrative  
1241 purposes. The solid line at  $p = 6.00E-6$  corresponds to a Benjamini-Hochberg FDR  $<$   
1242  $0.05$ . The stricter of the two thresholds,  $p = 6.00E-6$ , was used to define significant  
1243 SNVs and SVs. Significant variants overlapping regions in the ENCODE blacklist v2 are  
1244 shown in blue and were omitted from downstream analyses. FDR: False Discovery  
1245 Rate.

1246

1247 **Fig 2. Significant SNVs lie in genomic regions containing genes involved in**  
1248 **transposon control.**

1249 **(A)** Scheme for assessing whether greenlisted SNVs were enriched in regions  
1250 containing genes with TE regulatory potential. For a given gene set with regulatory  
1251 potential (regulatory set A or B), the proportion of SNVs near genes in that gene set  
1252 were calculated for the background and significant SNV lists, and statistical significance  
1253 was assessed using Fisher's exact test. **(B)** Enrichment analysis of greenlisted SNVs  
1254 near genes previously implicated in L1 expression control [37] or L1 transposition  
1255 control [33] by CRISPR screening in cancer cell lines. Three specific examples of  
1256 greenlisted SNVs near **(C)** genes controlling L1 expression and **(D)** genes controlling L1  
1257 transposition are shown. All GWAS associations have an FDR < 0.05, but, as  
1258 customary in the GWAS field, we report the raw p-values. **(E)** A summary of the  
1259 associations we identified with various TE regulatory gene sets, highlighting the number  
1260 of associated SNVs and the regulatory genes those SNVs were proximal to. Though not  
1261 exclusive regulators of TE activity *per se*, we included in our analysis gene sets involved  
1262 in "histone methyltransferase activity" and "RNA modification" functions, since those  
1263 processes have been implicated in transposon control. FDR: False Discovery Rate.

1264

1265 **Fig 3. Significant SNVs are enriched in genomic regions containing features**  
1266 **associated with genome instability.**

1267 **(A)** Scheme for assessing whether greenlisted SNVs are enriched in regions containing  
1268 elements known for promoting genome instability. For a given set of potentially  
1269 genetically unstable regions (unstable element set A or B), the proportion of SNVs  
1270 overlapping regions in each set are calculated for the background and significant SNV  
1271 lists, and statistical significance is assessed using Fisher's exact test. **(B)** Enrichment  
1272 analysis of greenlisted SNVs overlapping evolutionary age-stratified Alu (*left column*)  
1273 and L1 (*right column*) copies. **(C)** Enrichment analysis of greenlisted SNVs overlapping  
1274 curated L1 loci in L1Base v2 [107]. This database contains putatively active L1 copies  
1275 (with either full-length, fully intact L1 copies or L1 copies with only ORF2 intact), as well  
1276 as non-autonomous, full-length, non-intact L1 copies with regulatory potential. **(D)**  
1277 Enrichment analysis of greenlisted SNVs overlapping genomic regions containing  
1278 segmental duplications [66, 67] or structural variation hotspots [65]. FDR: False  
1279 Discovery Rate.

1280

1281 **Fig 4. Alterations in the cell cycle are positively correlated with case status.**

1282 **(A)** Scheme for characterizing transcriptomic differences between case and control  
1283 samples. Gene expression profiles were quantified using mRNA-sequencing data from  
1284 lymphoblastoid cells belonging to 358 European and 86 African individuals. To note, all  
1285 African individuals here were from the Yoruban population. These gene expression  
1286 profiles were used to construct consensus gene co-expression networks with WGCNA.  
1287 We then quantified the correlations between each module in the network and the case-  
1288 control status of all samples (encoded as 0 for controls and 1 for cases). Finally, over-  
1289 representation analysis (ORA) using the Gene Ontology (GO) Biological Process gene  
1290 set collection was used to assign functions to significantly correlated modules. **(B)** The  
1291 correlations and correlation p-values between consensus network modules and case-  
1292 control status in the European and African cohorts. Boxes were color-coded according  
1293 to the strength of the correlation. A meta-analysis was also carried out to summarize  
1294 statistical results by combining European and African correlation p-values using Fisher's  
1295 method. For visualization purposes only, the average of the European and African  
1296 correlations was assigned to the meta-analysis. Correlations with opposite trends in the  
1297 two cohorts were disregarded in the meta-analysis and colored grey. Correlations with  $p$   
1298  $< 0.01$  were considered statistically significant and highlighted in bold. **(C)** The top 10  
1299 ORA results for the MEroyalblue module using the GO Biological Process gene set  
1300 collection. The colors represent the gene ratio (i.e. the fraction of module genes from  
1301 the listed gene set) and the sizes of the dots represent the Benjamini-Hochberg FDR.  
1302 NA: Not Applicable. FDR: False Discovery Rate.

1303

1304

1305 **Legends to Tables**

1306

1307 **Table 1. Enrichment of genomic features overlapping global L1/Alu singletons.**

1308

HOMER Annotation	Number of peaks	Log <sub>2</sub> Ratio (obs/exp)	LogP enrichment (+values depleted)
Intron	907	0.451	-56.155
Intergenic	834	0.157	-8.634
DNA	133	0.426	-7.538
Unknown	4	2.474	-5.062
snRNA	2	2.681	-3.229
3UTR	35	0.371	-2.529
Low_complexity	8	0.563	-1.715
TTS	37	0.184	-1.429
LTR	261	0.042	-1.142
snoRNA	0	-0.001	0
SINE?	0	-0.004	0.003
RC?	0	-0.082	0.059
tRNA	0	-0.127	0.092
miRNA	0	-0.139	0.101
RNA	0	-0.149	0.109
scRNA	0	-0.166	0.122
rRNA	0	-0.265	0.201
srpRNA	0	-0.339	0.265
RC	0	-0.449	0.365
DNA?	0	-0.506	0.42
pseudo	2	-0.088	0.442
ncRNA	7	-0.033	0.554
Retroposon	3	-0.427	0.852
LTR?	0	-1.098	1.141
5UTR	1	-1.382	1.324
Promoter	30	-0.267	1.746
CpG-Island	3	-1.585	3.86
Simple_repeat	22	-0.615	3.863
LINE	550	-0.139	5.176
Exon	18	-1.057	8.089
Satellite	26	-1.493	22.78
SINE	141	-1.431	112.799

1309

1310

1311 **Table 2. The top 5 most significant, greenlisted variants with an odds ratio > 1**  
 1312 **(top) or < 1 (bottom).**

1313

Variant	Chr	Position	Nearby_Genes	Odds_Ratio	P_Meta	P_AFR	P_AMR	P_EAS	P_EUR	P_SAS
INV_delly_INV00066128	21	26001780	ATP5J, APP	4.3813	5.67E-23	2.16E-14	9.90E-13	3.27E-06	3.14E-06	3.38E-10
ALU_umary_ALU_7919	10	42173057	ZNF33B	3.1524	7.75E-14	2.42E-07	0.000214	0.000959	0.000269	1.61E-08
L1_umary_LINE1_1066	5	20834958	CDH18	2.2789	7.64E-13	2.41E-06	0.03278	0.01845	0.03176	1.86E-05
SVA_umary_SVA_602	14	64980975	RAB15, FNTB	2.221	3.07E-10	1.08E-06	0.008247	0.0006175	0.05464	0.001859
ALU_umary_ALU_3176	4	29731532	PCDH7	3.3541	5.17E-10	2.76E-10	0.0006186	0.0001465	0.0003296	4.82E-16
rs199980215	1	68007166	GNG12, DIRAS3	0.3235	4.91E-20	1.50E-08	5.24E-07	0.002965	0.03624	3.54E-05
rs1317852644	13	64517674		0.2899	2.87E-19	0.02896	5.08E-07	0.0002641	0.00253	3.15E-08
rs201291405	6	47190021	ADGRF1, TNFRSF21	0.2706	2.62E-18	3.54E-06	3.31E-05	0.02904	0.0103	1.36E-07
rs1056816611	14	31879290	ARHGAP5, NUBPL	0.3132	6.39E-18	1.86E-06	0.004238	0.01085	0.0001221	4.84E-07
rs1255393055	14	31879311	ARHGAP5, NUBPL	0.3132	6.39E-18	1.86E-06	0.004238	0.01085	0.0001221	4.84E-07

1314

1315

1316 **Legends to Supplementary Figures**

1317

1318 **S1 Fig. Quality control of combined SNV and SV 1000 Genomes Project data.**

1319 PCA plots for pruned SNV and SV genotype data from **(A)** African, **(B)** East Asian, **(C)**  
1320 European, **(D)** South Asian, and **(E)** Admixed American samples. Colors and shapes  
1321 represent different populations within each super-population.

1322

1323 **S2 Fig. GWAS in individual super-populations is underpowered.**

1324 Manhattan plots for GWAS results in individual super-populations, including for the **(A)**  
1325 African, **(B)** East Asian, **(C)** European, **(D)** South Asian, and **(E)** Admixed American  
1326 cohorts. Solid lines correspond to a Benjamini-Hochberg FDR < 0.05 and dashed lines  
1327 correspond to an average empirical FDR < 0.05, based on 20 random permutations.  
1328 The Benjamini-Hochberg FDR and average empirical FDR, respectively, corresponded  
1329 to the following p-values in each super-population:  $p = 1.18E-6$  and  $p = 4.61E-6$  in the  
1330 African cohort,  $p = 9.06E-7$  and  $p = 8.46E-7$  in the South Asian cohort, and  $p = 3.53E-7$   
1331 and  $p = 1.07E-6$  in the Admixed American cohort. The stricter of the two thresholds in  
1332 each super-population was used to define significant SNVs and SVs. No variant at an  
1333 FDR < 0.05 was identified in the East Asian and European cohorts. Significant variants  
1334 overlapping regions in the ENCODE blacklist v2 are shown in blue. FDR: False  
1335 Discovery Rate.

1336

1337 **S3 Fig. GWAS meta-analysis associations are conserved in non-African super-**  
1338 **populations.**

1339 **(B)** A schematic illustrating the trans-ethnic integration of available SNV and SV data to  
1340 identify variants associated with L1/Alu insertion global singletons in non-African super-  
1341 populations. Within each non-African super-population, samples were segregated into  
1342 cases and controls depending on whether or not they harbored a global Alu or L1  
1343 insertion singleton. GWAS was carried out within each super-population to identify  
1344 polymorphic SNVs and SVs associated with case-control status. Finally, GWAS results  
1345 from all 4 non-African super-populations were meta-analyzed using a random effects  
1346 statistical model, yielding a summary meta-analysis odds ratio and p-value for each

1347 variant. **(B)** A Manhattan plot for the trans-ethnic GWAS meta-analysis omitting the  
1348 African super-population. The dashed line at  $p = 6.60E-6$  corresponds to an average  
1349 empirical FDR  $< 0.05$ , based on 20 random permutations. One such permutation is  
1350 shown in the bottom panel for illustrative purposes. The solid line at  $p = 1.92E-6$   
1351 corresponds to a Benjamini-Hochberg FDR  $< 0.05$ . The stricter of the two thresholds,  $p$   
1352  $= 1.92E-6$ , was used to define significant SNVs and SVs. Significant variants  
1353 overlapping regions in the ENCODE blacklist v2 are shown in blue and were omitted  
1354 from downstream analyses. **(C)** A Venn diagram comparing the number of significant  
1355 variants identified using either all five super-populations or all four non-African super-  
1356 populations. The statistical significance of the overlap was calculated using a one-sided  
1357 Fisher's exact test. **(D)** A Venn diagram comparing the number of significant variants  
1358 identified using either all four non-African super-populations or only the African super-  
1359 population. The statistical significance of the overlap was calculated using a one-sided  
1360 Fisher's exact test. FDR: False Discovery Rate.

1361

#### 1362 **S4 Fig. Functional annotation of significant variants.**

1363 **(A)** Scheme for predicting functions of genomic regions containing significant variants.  
1364 All SNVs and SVs used in this study were assigned genes using the GREAT [109]  
1365 online platform. Significant SNV- and SV-associated genes were then tested for  
1366 functional gene set enrichment by over-representation analysis (ORA) using  
1367 clusterProfiler [119], specifying the corresponding background-associated genes as the  
1368 universe. **(B)** The number of genes associated to greenlisted SNVs (*left*), the distance  
1369 between greenlisted SNVs and the transcriptional start sites (TSS) of associated genes  
1370 (*middle*), and the top 10 ORA results for associated genes using the GO Biological  
1371 Process gene set collection (*right*). The colors represent the gene ratio (i.e. the fraction  
1372 of significant SNV-associated genes from the listed gene set) and the sizes of the dots  
1373 represent the Benjamini-Hochberg FDR. **(C)** The number of genes associated to  
1374 greenlisted SVs (*left*) and the distance between greenlisted SVs and the transcriptional  
1375 start sites (TSS) of associated genes (*middle*). **(D)** Enrichment analysis of greenlisted  
1376 SNVs overlapping genomic regions with candidate cis-Regulatory Elements (cCREs)  
1377 from the ENCODE Registry v4 [55]. **(E)** Heatmap comparing the median expression

1378 levels of significant SNV-associated genes in each tissue included in the GTEx Analysis  
1379 v8. Values in the heatmap represent z-scores of gene expression across tissues.  
1380 Negative values are lower than average tissue expression and positive values are  
1381 higher than average tissue expression. FDR: False Discovery Rate.

1382

1383 **S5 Fig. Polymorphic SVs of different classes are associated with L1/Alu insertion**  
1384 **singletons.**

1385 **(A)** One example of each type of significant, polymorphic SV that was associated with  
1386 L1/Alu singletons. These classes included inversions, Alu insertions, an L1 insertion,  
1387 SINE-VNTR-Alu (SVA) insertions, and a multiallelic copy number variant. All GWAS  
1388 associations have an FDR < 0.05, but, as customary in the GWAS field, we report the  
1389 raw p-values. FDR: False Discovery Rate.

1390

1391



1392 **Inventory of Supplementary Tables**

1393

1394 **Supplementary Table S1. Sample singleton counts, significant variant**  
1395 **annotations, and gene co-expression network results.**

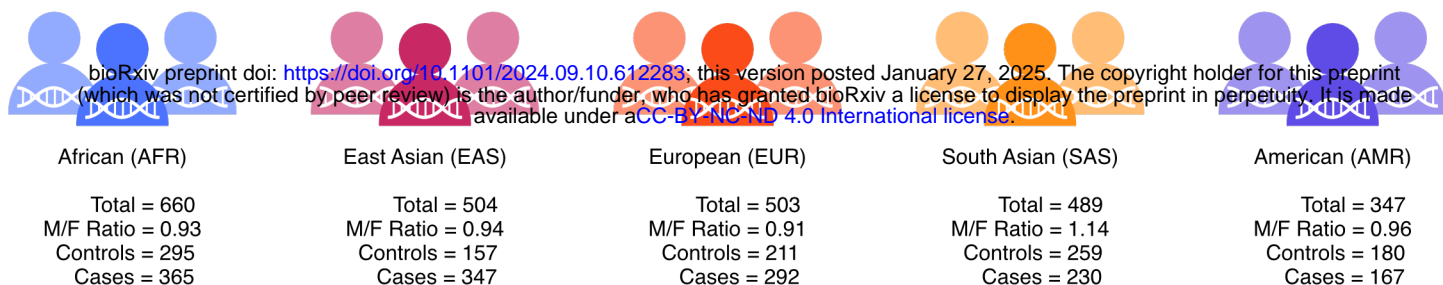
1396 **(A)** Number of singletons for each sample. **(B)** All variants passing  $FDR < 0.05$  in the  
1397 GWAS meta-analysis. **(C)** Significant variants uniquely identified in the non-African  
1398 meta-analysis and absent in the complete meta-analysis. **(D)** Over-representation  
1399 analysis of genes associated to greenlisted, significant SNVs using GO Biological  
1400 Process gene sets. **(E)** Over-representation analysis of genes associated to greenlisted,  
1401 significant SVs using GO Biological Process gene sets. **(F)** SnpEff annotations of  
1402 significant variants. **(G)** Expression levels (median tissue-specific TPMs) of significant  
1403 SNV-associated genes in the GTEx Analysis v8. The cluster of brain-associated genes  
1404 is in blue, and the cluster of testes-associated genes is in orange. **(H)** TE enrichment  
1405 analysis with TEENA. **(I)** Lymphoblastoid cell line WGCNA network gene-module  
1406 assignments. **(J)** Over-representation analysis of MEroyalblue genes using GO  
1407 Biological Process gene sets.

1408

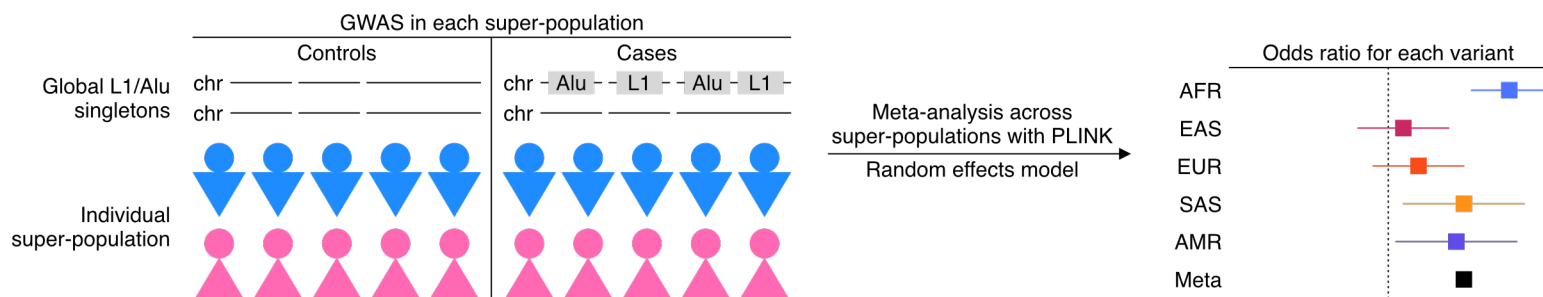
1409

**Figure 1**

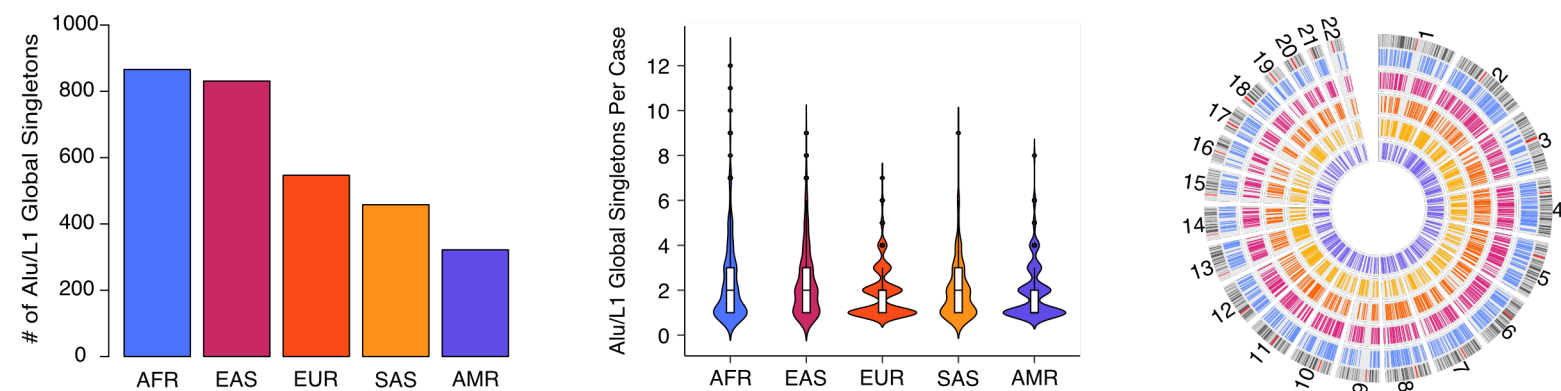
**A** The 2503 samples from the 1000Genomes Project used in this study



**B** Case-control study design for trans-ethnic GWAS



**C** Frequency and genomic distribution of L1/Alu global singletons across superpopulations



**D** Manhattan plot for the trans-ethnic GWAS meta-analysis associations

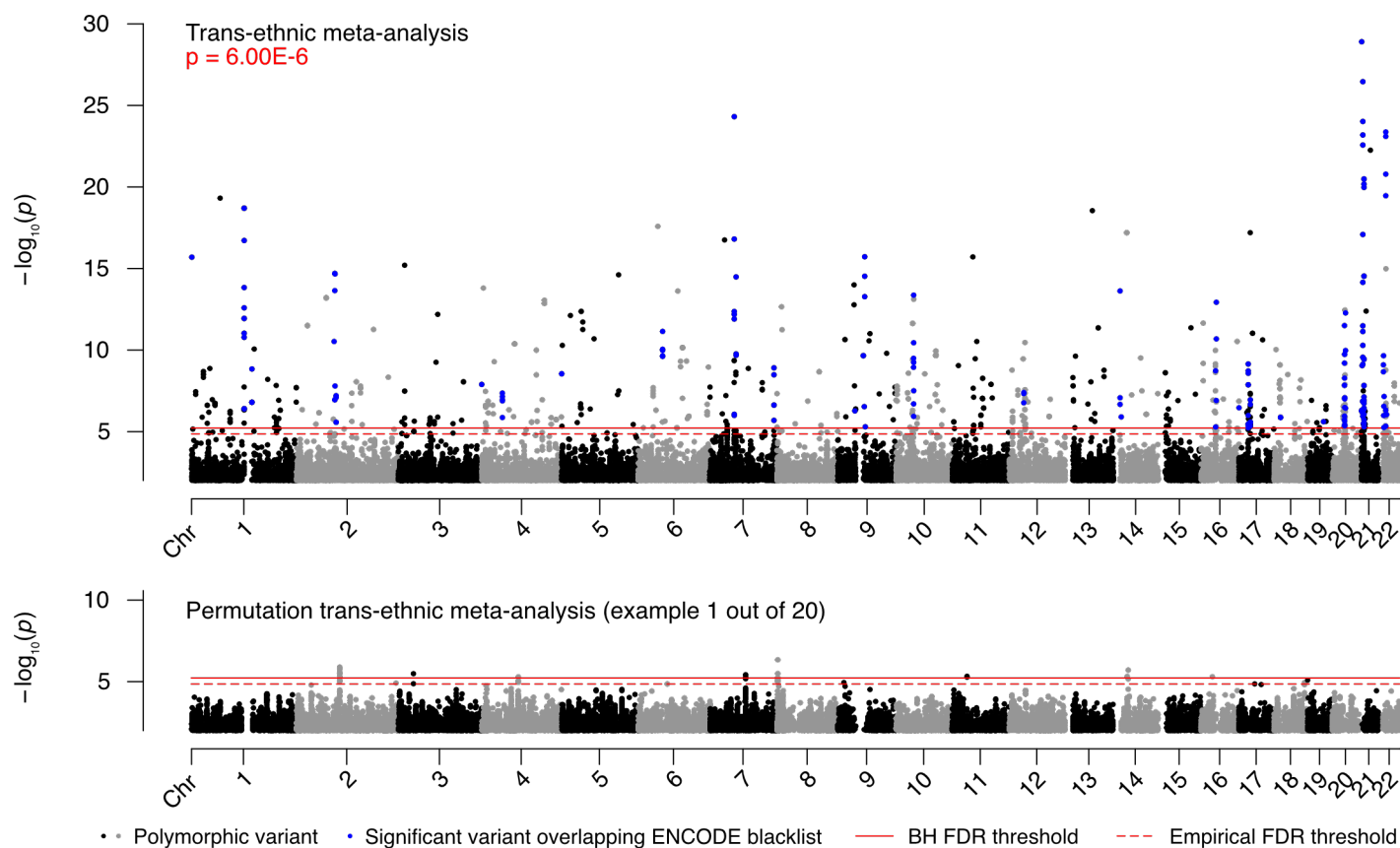
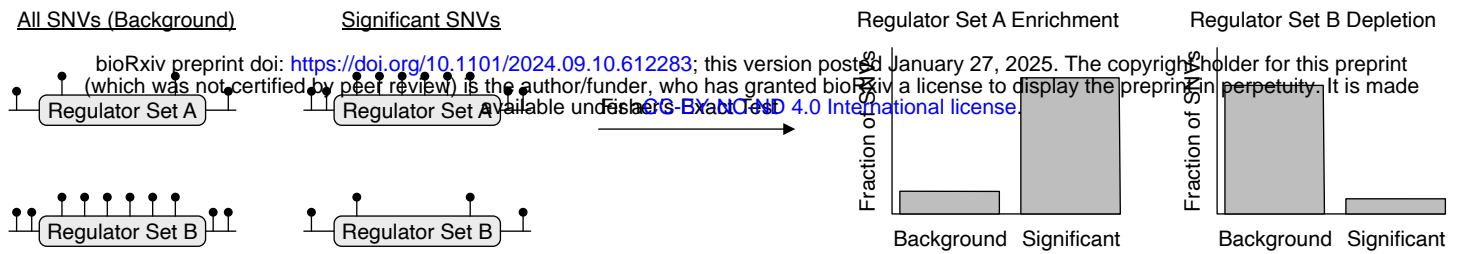
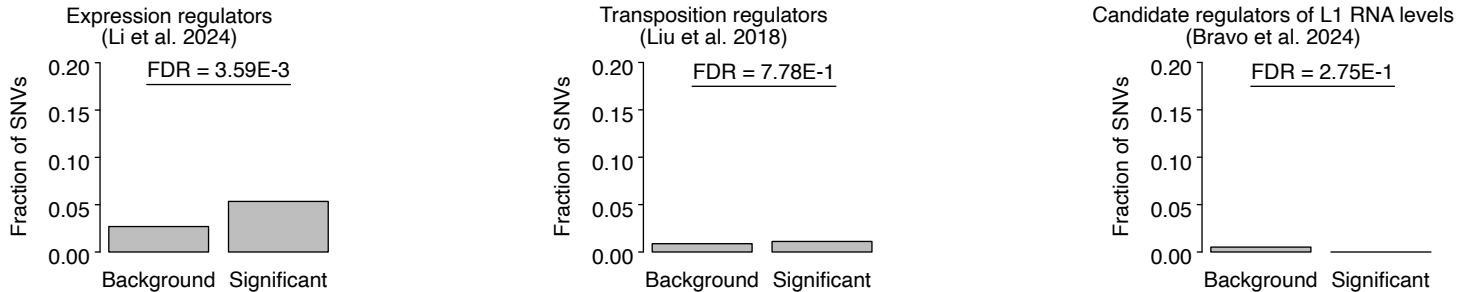


Figure 2

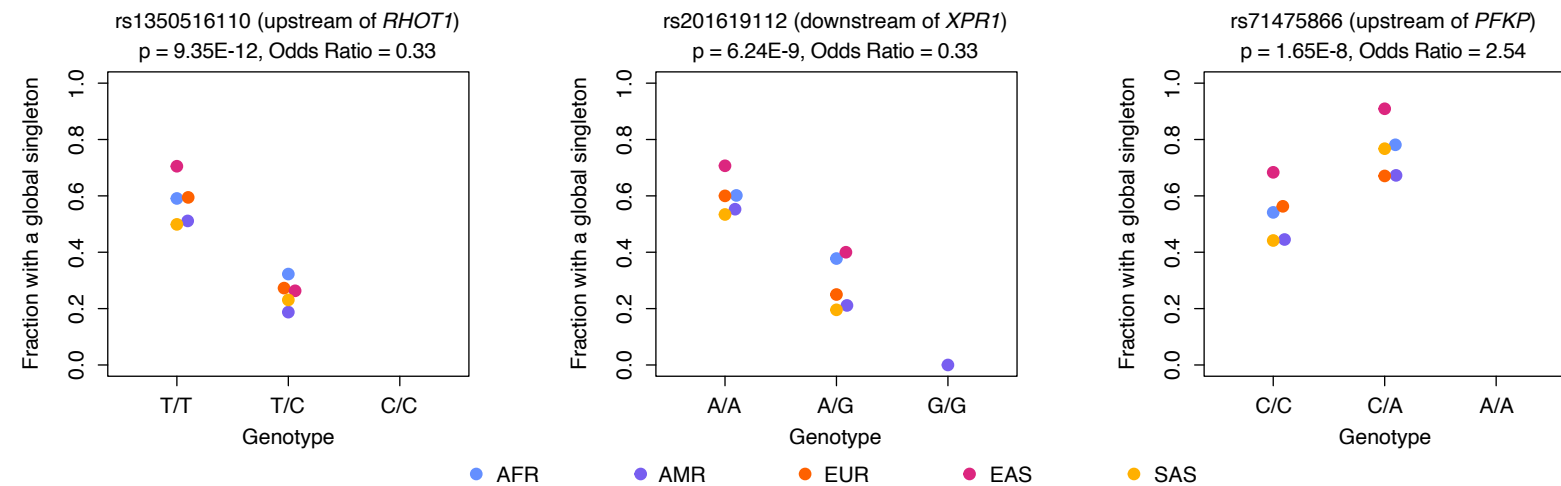
**A** Scheme for linking significant SNVs with potential transposon regulators



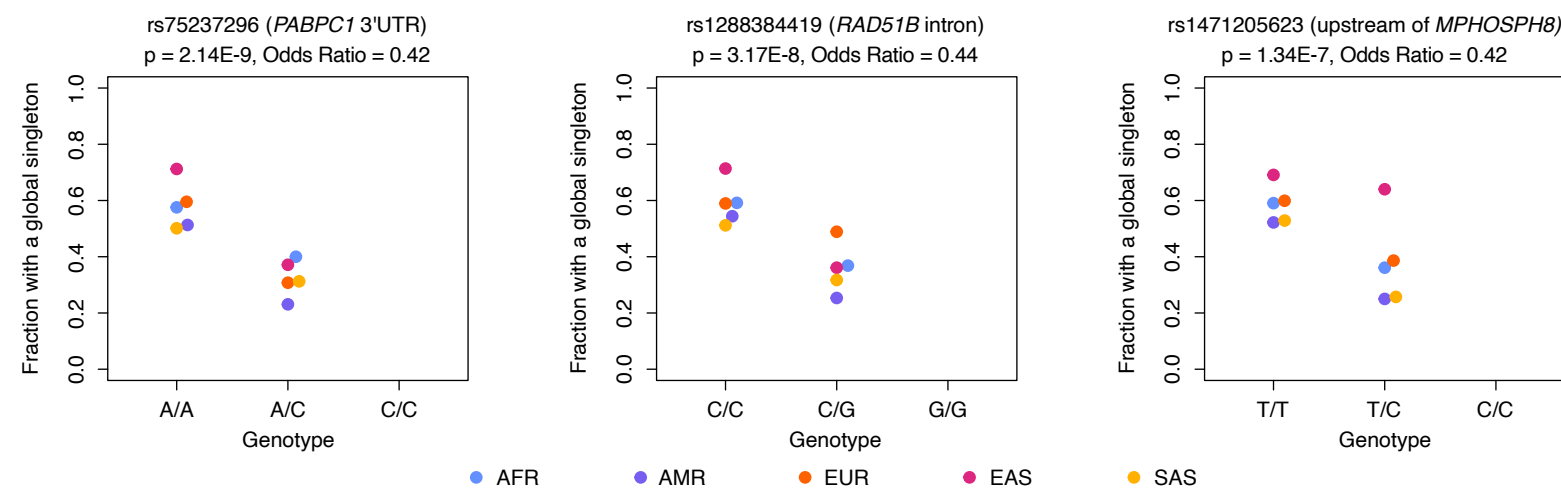
**B** SNVs near known L1 regulators



**C** Example SNVs near L1 expression regulators



**D** Example SNVs near L1 transposition regulators

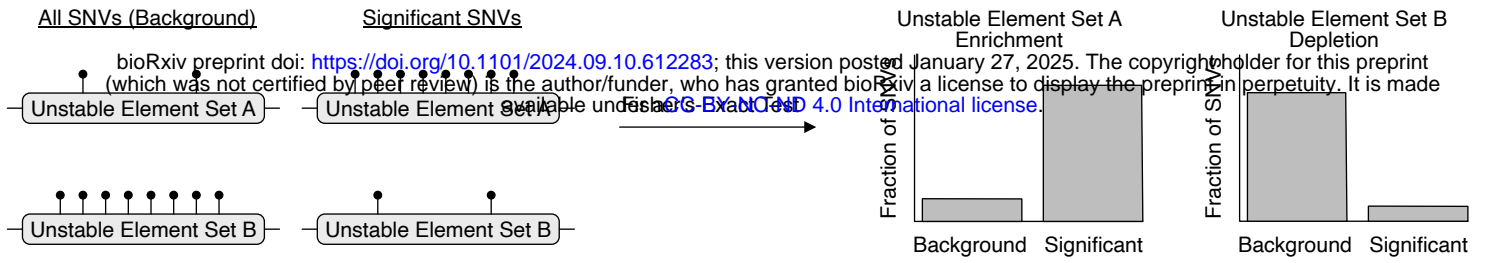


**E** Summary of SNVs near known TE regulators and genes with functions involved in TE regulation

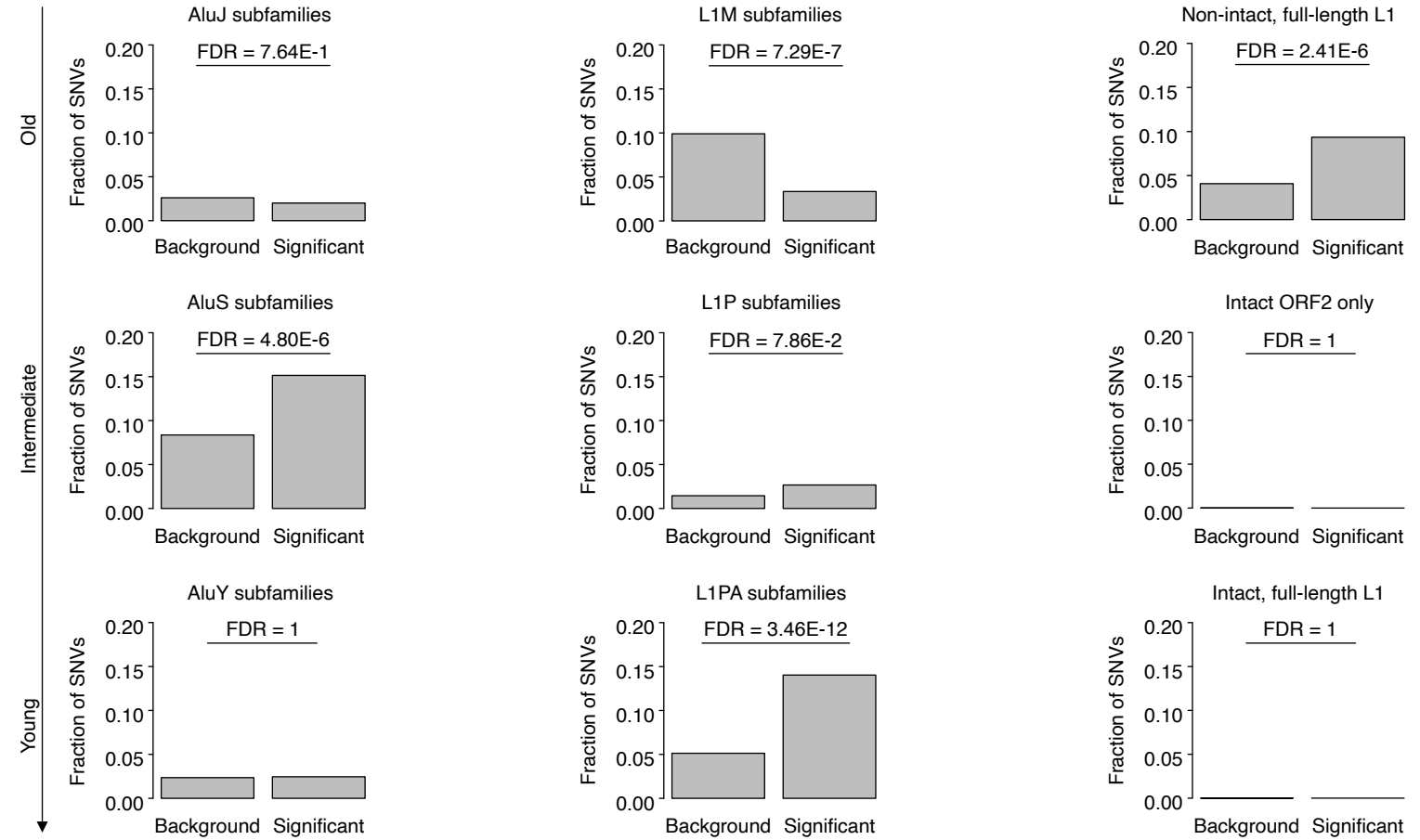
Regulatory gene set	Significant SNVs near genes	Genes near significant SNVs
L1 expression regulators (Li et al. 2024)	24	<i>IPO9, MAPK14, METTL14, MPHOSPH8, PFKP, PHF3, RBPJ, RHOT1, RRAGA, XPR1</i>
L1 transposition regulators (Liu et al. 2018)	5	<i>MPHOSPH8, PABPC1, RAD51B</i>
Candidate regulators of L1 RNA levels (Bravo et al. 2024)	0	
Histone methyltransferase activity GO:0042054	4	<i>EEF2KMT, PRDM7</i>
RNA modification GO:0009451	8	<i>A1CF, ADARB2, METTL14</i>

**Figure 3**

**A** Scheme for linking significant SNVs with transposons and other regions of potential genomic instability



**B** SNV overlap with evolutionary age-stratified Alu and L1 copies



**C** SNV overlap with L1Base v2 annotations

**D** SNV overlap with regions of potential genomic instability

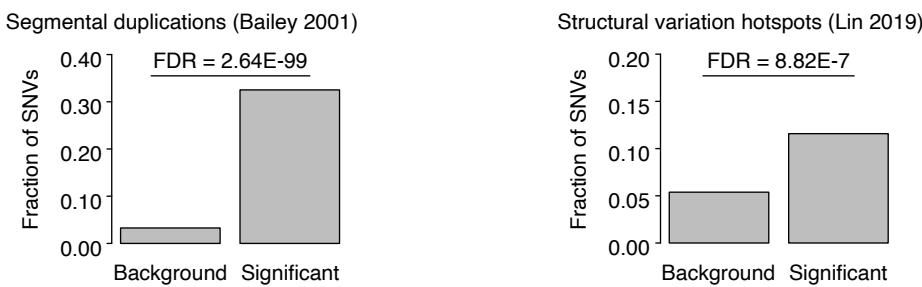
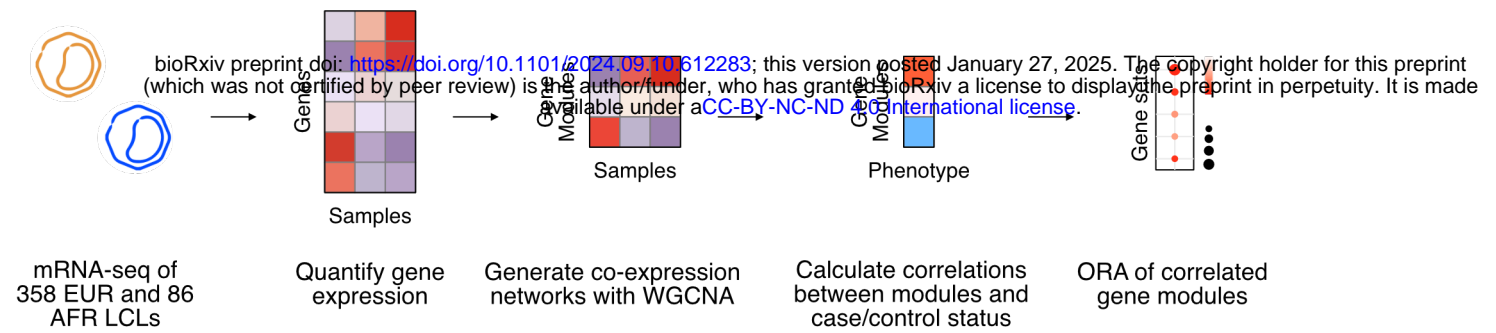


Figure 4

**A** Scheme for assessing transcriptomic differences between cases and controls



**B** Network module correlations with case/control status

Module	EUR case status	AFR case status	Meta-analysis
MEblack	0.077 (0.1)	0.056 (0.6)	0.066 (0.3)
MEpurple	0.083 (0.1)	0.057 (0.6)	0.07 (0.3)
MEgreenyellow	0.039 (0.5)	0.09 (0.4)	0.064 (0.5)
MEmidnightblue	-0.0024 (1)	0.12 (0.3)	NA (NA)
MEsalmon	-0.12 (0.03)	0.09 (0.4)	NA (NA)
MElightyellow	0.084 (0.1)	0.04 (0.7)	0.062 (0.3)
MEcyan	-0.039 (0.5)	0.11 (0.3)	NA (NA)
MEred	-0.057 (0.3)	0.23 (0.03)	NA (NA)
MElightcyan	0.017 (0.8)	-0.13 (0.2)	NA (NA)
MEroyalblue	0.035 (0.5)	<b>0.37 (4E-4)</b>	<b>0.2 (0.002)</b>
MEblue	0.036 (0.5)	0.18 (0.1)	0.11 (0.2)
MEbrown	-0.021 (0.7)	0.15 (0.2)	NA (NA)
MEturquoise	-0.087 (0.1)	-0.12 (0.3)	-0.1 (0.1)
MEgrey60	-0.022 (0.7)	0.08 (0.5)	NA (NA)
MEmagenta	-0.03 (0.6)	-0.14 (0.2)	-0.087 (0.3)
MEgreen	0.0054 (0.9)	-0.14 (0.2)	NA (NA)
MEtan	0.046 (0.4)	-0.0054 (1)	NA (NA)
MEyellow	0.048 (0.4)	-0.15 (0.2)	NA (NA)
MElightgreen	0.049 (0.4)	-0.048 (0.7)	NA (NA)
MEpink	0.074 (0.2)	-0.15 (0.2)	NA (NA)
MEgrey	-0.052 (0.3)	0.014 (0.9)	NA (NA)

Correlation (p-value) scale: -1 to 1

**C** ORA of the MEroyalblue module with the GO Biological Process gene sets

