

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

3

4 Juan I. Bravo^a, Lucia Zhang^{a,b}, Bérénice A. Benayoun^{a,c,d,e,f,*}

5 ^a Leonard Davis School of Gerontology, University of Southern California, Los Angeles,
6 CA 90089, USA.

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

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

11 ^d Biochemistry and Molecular Medicine Department, USC Keck School of Medicine, Los
12 Angeles, CA 90089, USA.

13 ^e USC Norris Comprehensive Cancer Center, Epigenetics and Gene Regulation, Los
14 Angeles, CA 90089, USA.

15 ^f 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 Long INterspersed Element-1 (LINE-1; L1) and Alu are two families of
24 transposable elements (TEs) occupying ~17% and ~11% of the human genome,
25 respectively. Though only a small fraction of L1 copies is able to produce the machinery
26 to mobilize autonomously, Alu elements and degenerate L1 copies can hijack their
27 functional machinery and mobilize *in trans*. The expression and subsequent copy
28 number expansion of L1 and Alu can exert pathological effects on their hosts, promoting
29 genome instability, inflammation, and cell cycle alterations. These features have made
30 L1 and Alu promising focus subjects in studies of aging and aging diseases where they
31 can become active. However, the mechanisms regulating variation in their expression
32 and copy number remain incompletely characterized. Moreover, the relevance of known
33 mechanisms to diverse human populations remains unclear, as mechanisms are often
34 characterized in isogenic cell culture models. To address these gaps, we leveraged
35 genomic data from the 1000 Genomes Project to carry out a trans-ethnic GWAS of L1
36 and Alu insertion global singletons. These singletons are rare insertions observed only
37 once in a population, potentially reflecting recently acquired L1 and Alu integrants or
38 structural variants, and which we used as proxies for L1/Alu-associated copy number
39 variation. Our computational approach identified single nucleotide variants in genomic
40 regions containing genes with potential and known TE regulatory properties, and it
41 enriched for single nucleotide variants in regions containing known regulators of L1
42 expression. Moreover, we identified many reference TE copies and polymorphic
43 structural variants that were associated with L1/Alu singletons, suggesting their potential
44 contribution to TE copy number variation through transposition-dependent or
45 transposition-independent mechanisms. Finally, a transcriptional analysis of
46 lymphoblastoid cells highlighted potential cell cycle alterations in a subset of samples
47 harboring L1/Alu singletons. Collectively, our results (i) suggest that known TE
48 regulatory mechanisms may also play regulatory roles in diverse human populations, (ii)
49 expand the list of genic and repetitive genomic loci implicated in TE copy number
50 variation, and (iii) reinforce the links between TEs and disease.

51

52 **KEYWORDS:** LINE-1, Alu, transposons, GWAS, regulators, copy number

53 1. INTRODUCTION

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

71
72 For a transposition-dependent expansion of L1 copy number to occur, L1 must
73 undergo a multi-step lifecycle. This lifecycle includes (i) transcription of an active, full-
74 length L1 copy, (ii) potential RNA processing of the L1 transcript, (iii) nuclear export of
75 the transcript to the cytoplasm, (iv) translation of the two open reading frames (ORFs),
76 (v) potential post-translational modifications of ORF1p/ORF2p and binding of those
77 proteins to the transcripts that produced them (*cis* preference) to form ribonucleoprotein
78 (RNP) complexes, (vi) entry into the nucleus, and finally (vii) reverse transcription and
79 integration by a process called target primed reverse transcription (TPRT) (reviewed in
80 [3]). Importantly, though neither Alu elements or degenerate L1 copies can mobilize
81 autonomously, they can hijack proteins from transposition-competent L1s and mobilize
82 *in trans* [6, 7, 9]. Though not traditionally considered part of the L1/Alu lifecycles, other
83 transposition-independent, but homology-dependent, mechanisms can further

84 contribute to TE copy number variation. This includes repeat-mediated deletion (RMD)
85 events whereby two repetitive elements (often Alu elements) on the same chromatid
86 can recombine to cause a deletion rearrangement, potentially resulting in the deletion of
87 one of the repeats as well as the intervening sequence, which may include additional
88 repeats [10, 11]. More broadly, TE-mediated and TE-independent non-allelic
89 homologous recombination (NAHR) events ([12-15] and reviewed in [16, 17]) can
90 directly generate large chromosomal deletions and duplications, which may include
91 repetitive sequences. From an evolutionary perspective, the burst of Alu transposition
92 in primates is hypothesized to have sensitized the ancestral genome to Alu-mediated
93 recombination events, which may have played a role in the emergence and expansion
94 of segmental duplications, which would provide additional substrates for NAHR [18].
95 Ultimately, TE copy number is shaped by a combination of *de novo* insertions resulting
96 from their lifecycle and genomic structural remodeling that can expand or retract the
97 copy number.

98
99 Characterizing the mechanisms governing L1 and Alu transcriptional and copy
100 number control will be important, given their associations with, and potential
101 contributions to, aging and aging-associated diseases like cancer (discussed in [19-21]).
102 Fundamentally, L1 and/or Alu can alter several hallmarks of aging [22], such as
103 genomic instability, cellular senescence, and inflammation. Though the origin of the
104 signal is unclear (genomic, extra-chromosomal, or cytosolic), an increase in L1 copy
105 number has been observed with chronological aging [23] and during cellular
106 senescence [24]. Moreover, a key feature of cellular senescence is the senescence-
107 associated secretory phenotype (SASP) whereby cells secrete an amalgamation of pro-
108 inflammatory factors [25] that may contribute to chronic, low-grade, sterile inflammation
109 with chronological age (a phenomenon referred to as “inflamm-aging”) [25, 26]. L1 can
110 induce a senescent-like state *in vitro* in several cell lines [27, 28] and its cytoplasmic
111 complementary DNA (cDNA) is implicated in the maturation of the SASP response and
112 the establishment of deep senescence through the production of interferons [29].
113 Similarly, Alu RNA can upregulate senescence markers in retinal pigment epithelium
114 (RPE) cells from human eyes with geographic atrophy [30], and knockdown of Alu

115 transcripts were reported to promote senescence exit in adult adipose-derived
116 mesenchymal stem cells [31]. These findings highlight the relevance of L1 and Alu
117 retrotransposons in pathological, age-associated features.

118
119 To maintain homeostasis, it is imperative that host cells tightly regulate TE
120 activity (reviewed in [32, 33]). These mechanisms, however, remain incompletely
121 characterized due to the cell-specific and transposon-specific nature of TE regulatory
122 mechanisms. Indeed, no systematic, genome-wide screen for regulators of Alu
123 expression or copy number has been carried out thus far, to our knowledge. To address
124 these gaps, a number of *in vitro* and *in silico* approaches have been developed to scan
125 for novel regulators of TE expression or copy number. *In vitro* approaches have relied
126 on clustered regularly interspaced short palindromic repeats (CRISPR)-based and small
127 RNA-based tools to decipher L1 regulation in several types of cancer cells [34-38].
128 These approaches, however, can be technically challenging to implement in non-
129 cancerous cells, like primary cells, which may not tolerate hyper-elevated transposon
130 activity or that may resist genetic perturbations. To complement these methods, a
131 number of *in silico* approaches have been developed that utilize chromatin
132 immunoprecipitation followed by sequencing (ChIP-seq) [39], gene co-expression
133 networks [40], or copy number-expression correlations [41] to explore L1 regulation
134 without external manipulations. More recently, we screened for candidate regulators of
135 L1 RNA levels in lymphoblastoid cell lines (LCLs) using trans-expression quantitative
136 trait locus (trans-eQTL) analysis [42]. These tools highlight the need for, and usefulness
137 of, alternative approaches that utilize increasingly available, large ‘-omic’ datasets to
138 identify potentially novel mechanisms of TE control.

139
140 In this study, we develop a genome-wide association study (GWAS) pipeline to
141 identify genomic loci associated with global L1 and Alu insertion singletons in diverse
142 human populations. Global singleton insertions are rare insertions observed only once
143 in a population [43], potentially reflecting recently acquired L1 and Alu integrants or
144 structural variants. Thus, we used insertion singletons as proxies for L1/Alu-associated
145 copy number variation, which can arise through *de novo* transposition events or

146 alternative mechanisms. We demonstrate that our GWAS approach captures, and
147 enriches, genomic regions containing known and potential regulators of TE activity. We
148 observe that this approach also captures reference insertions and polymorphic
149 structural variants that may influence L1 or Alu copy number variation through
150 transposition-dependent or -independent mechanisms. Finally, we note that associated
151 loci fall into a few genes with clinical relevance, strengthening the association between
152 TEs and disease.

153 2. RESULTS

154 2.1 Identification of genomic loci associated with L1/Alu singletons in diverse human 155 populations

156 To unbiasedly identify potential genetic sources of L1 and Alu copy number
157 variation in human populations, we leveraged a publicly available human “omic” dataset
158 with thoroughly characterized genetic information. For this analysis, we utilized 2503
159 multi-ethnic samples from the 1000 Genomes Project for which both single nucleotide
160 variant (SNV) and structural variant (SV) data were available. Specifically, this included
161 individuals from 5 super-populations: 660 African (AFR), 504 East Asian (EAS), 503
162 European (EUR), 489 South Asian (SAS), and 347 Admixed American (AMR)
163 individuals who declared themselves healthy at the time of sample collection (**Figure**
164 **1A**). As a quality control step, we checked whether the combined SNV and SV data
165 segregated samples by population following principal component analysis (PCA). These
166 analyses demonstrated that the top four principal components segregated population
167 groups within each super-population (**Figure S1A-S1E**).

168

169 We then carried out an integration of the available multi-ethnic SNV and SV
170 genomic data (**Figure 1B**). For the phenotype, we focused on global singleton SVs,
171 which are rare SVs that occur exactly once in a study [43], for L1 and Alu insertions.
172 Given their rarity, these insertion singletons may reflect recently acquired L1 and Alu
173 integrants or structural variants. Thus, we hypothesized that these global singletons
174 could serve as proxies for elevated L1/Alu-associated copy number variation, which
175 may arise from *de novo* transposition events or alternative mechanisms. Importantly,
176 Alu is dependent on L1 machinery for its mobilization [6, 7] and both L1 and Alu can
177 contribute to structural variation through recombination-based mechanisms [12, 14, 15],
178 so variation in L1 and Alu copy numbers are likely to have overlapping regulatory
179 mechanisms. Thus, we first split our samples into cases and controls, depending on
180 whether or not they contained an L1 and/or an Alu insertion global singleton (**Table**
181 **S1A**). Second, we carried out a GWAS within each super-population to identify
182 common, polymorphic SNVs and SVs associated with case-control status. Third, to
183 maximize statistical power and identify shared, trans-ethnic sources of TE singleton

184 number variation, we meta-analyzed our GWAS results across the 5 super-populations
185 using a random effects statistical model [44, 45]. Interestingly, several hundred L1/Alu
186 global singleton insertions were detected in each super-population, ranging from 322 in
187 the American cohort to 866 in the African cohort (**Figure 1C**). Though most case
188 samples had 1-3 L1/Alu global singletons, several samples exhibited much more
189 extreme TE singleton accumulation, especially within the African super-population
190 (**Figure 1C**). Finally, we note that L1/Alu global singleton insertions were distributed
191 across all autosomes in all 5 super-populations (**Figure 1C**).

192
193 As expected, GWAS in each super-population was generally underpowered
194 (**Figure S2A-S2E**). Though we were able to identify many significant (FDR < 0.05)
195 variants in the African cohort (**Figure S2A**), we could not identify significant variants in
196 the East Asian (**Figure S2B**) and European (**Figure S2C**) cohorts, and we were only
197 able to identify a handful of significant variants in the South Asian (**Figure S2D**) and
198 Admixed American (**Figure S2E**) cohorts. In contrast, the GWAS meta-analysis
199 integrating all super-populations identified 658 significant variants distributed across all
200 22 autosomes, though there was especially strong and recurrent signal on chromosome
201 21 (**Figure 1D, Table S1B**). To simplify functional annotation of significant variants and
202 discard potential false positives, we omitted from downstream analyses significant
203 variants overlapping the “ENCODE blacklist v2” [46]. This curated “blacklist” represents
204 a set of genomic regions with anomalous signal in next-generation sequencing
205 experiments independent of cell type and individual experiment [46], and SNVs
206 overlapping these regions were significantly enriched in our significant SNV list
207 compared to the background SNV list (Fisher’s exact test, FDR = 1.33E-302). After
208 filtering 188 blacklisted SNVs, we were left with 449 greenlisted SNVs and 21
209 greenlisted SVs that were significantly associated with case-control status (**Figure 1D,**
210 **Table S1B**).

211
212 To assess the potential functions of greenlisted, significant variants, nearby
213 genes were assigned to variants using the Genomic Regions Enrichment of Annotations
214 Tool (GREAT) [47, 48] and over-representation analysis (ORA) was carried out using

215 clusterProfiler [49] (**Figure S3A**). Except for 51 greenlisted SNVs that were not linked to
216 any gene, the remaining SNVs were linked to 1-2 genes and were found predominantly
217 within 500 kilobases of a transcriptional start site (TSS) (**Figure S3B**). This observation
218 highlights the association of intergenic and gene-proximal, rather than distal, genetic
219 variation with L1/Alu copy number differences. Over-representation analysis of the
220 associated genes using the Gene Ontology (GO) Biological Process gene set revealed
221 an enrichment of terms related to heart development (such as ‘regulation of heart
222 growth’, ‘cardiac chamber morphogenesis’, and ‘positive regulation of cardiac muscle
223 cell proliferation’) and neuronal function (such as ‘neuron recognition’ and ‘axonal
224 fasciculation’; **Figure S3B, Table S1C**). Interestingly, genes related to ‘reproduction’
225 were also significantly over-represented among our list of associated genes (**Table**
226 **S1C**). Similar to the SNVs, greenlisted SVs were all linked to 1-2 genes and were
227 mostly within 500 kilobases of an annotated TSS (**Figure S3C**). Likely due to the low
228 number of greenlisted SVs, and consequently low number of associated genes, we
229 were unable to identify any significantly enriched GO Biological Process gene sets
230 (**Table S1D**). Given the limited number of greenlisted SVs and the unavailability of SV
231 sequences, we largely focused on greenlisted SNVs in downstream enrichment
232 analyses.

233
234 As a complementary approach to GREAT, we also predicted the functional
235 impact of significant variants using SnpEff [50] (**Table S1E**). Most variants were
236 assigned a ‘modifier’ impact by SnpEff—this annotation describes non-coding variants
237 where definitive functional predictions are not straightforward. One exception to this
238 included an SNV (rs367696690) introducing a synonymous substitution (p.Asp1605Asp)
239 with a low predicted impact in *IGFN1* (immunoglobulin like and fibronectin type III
240 domain containing 1). It is worth noting that synonymous substitutions can still impact
241 mRNA levels and protein function because of host preferences for specific codons (i.e.
242 codon usage bias; reviewed in [51]). Another exception was a missense variant
243 (rs1406034099) introducing a nonsynonymous substitution (p.Leu13639Phe or
244 p.Leu478Phe) with a moderate predicted impact in *MUC16* (mucin 16, cell surface
245 associated). Nonetheless, we highlight a few variants which overlapped clinically

246 relevant genes. For example, the most significant, greenlisted variant we identified was
247 an inversion SV (INV_delly_INV00066128) residing in an intronic Alu copy within the
248 *APP* (amyloid beta precursor protein) gene, an important biological marker for
249 Alzheimer's disease (AD) [52]. Similarly, we identified several SNVs (rs61994687,
250 rs1175403595, rs1343402870) in intronic or downstream regions of *PWRN1* within the
251 Prader-Willi syndrome (PWS) region. To explore non-protein-coding roles greenlisted
252 SNVs may play, we assigned them to ENCODE candidate cis-Regulatory Elements
253 (cCREs) [53] (**Figure S3D**). Although about 8% (36/449) of greenlisted SNVs resided in
254 an ENCODE cCRE, these were significantly depleted (FDR = 6.86E-15) in our
255 greenlisted SNVs compared to background SNVs (**Figure S3D**). Ultimately, our results
256 suggest that proximal, intergenic variation is associated with L1/Alu insertion singleton
257 number variation.

258

259 To further explore the potential functional roles of SNV-associated genes, we
260 leveraged expression data from the Genotype-Tissue Expression (GTEx) Portal [54-56]
261 to assess the pattern of expression of genes linked to greenlisted SNVs across tissues
262 (**Figure S3E, Table S1F**). In particular, expression patterns in the brain and gonads
263 were of special interest, given that L1 activity tends to be more frequent in those tissues
264 compared to others (discussed in [57]), and that L1/Alu integration events observed in
265 our GWAS would have to occur in the germline for transmission across generations.
266 Thus, we reasoned that if SNV-associated genes played roles in L1/Alu singleton
267 number variation, they may be more abundantly expressed in the brain and in gonads.
268 Interestingly, there was a cluster of genes that were very abundantly expressed in
269 testes but not in other tissue types (including ovarian tissue), suggesting the existence
270 of potential sex-specific mechanisms of *de novo* L1/Alu insertion transmission.
271 Furthermore, there was also a cluster of genes that were abundantly expressed across
272 brain regions and much less abundantly expressed across other tissue types. More
273 generally, there were many SNV-linked genes that were abundantly expressed in more
274 than ~50% of tissue types. These results highlight that significant SNV-associated
275 genes have tissue-specific expression patterns, including some genes that are very
276 abundant in tissues with documented L1 activity.

277

278 *2.2 Significant SNVs are enriched near regulators of L1 expression*

279 One of the primary motivations for carrying out this study was to identify novel,
280 candidate regulators of L1 and/or Alu copy number. In particular, there is a gap in our
281 understanding of Alu regulation, as, to our knowledge, no systematic screen for
282 regulators of Alu activity has previously been carried out. To determine whether our
283 approach captured genes with transposon regulatory potential, we assessed whether
284 our list of greenlisted SNVs was enriched for 1) genes with known TE regulatory
285 capabilities and 2) genes in broader pathways involved in TE regulation (**Figure 2A**).
286 Previously, two CRISPR-based screens for regulators of L1 expression [38] and L1
287 transposition [34] were carried out in cancer cell lines. In addition, we also recently
288 carried out an eQTL-based computational screen for candidate regulators of L1 RNA
289 levels in lymphoblastoid cell lines [42]. Interestingly, our greenlisted SNVs were
290 significantly (FDR = 3.59E-3) enriched in regions near known L1 expression regulators
291 compared to the background list of all SNVs (**Figure 2B, 2C, 2E**). A few examples of
292 these associations included rs1350516110 which was upstream of *RHOT1*,
293 rs201619112 which was downstream of *XPR1*, and rs71475866 which was upstream of
294 *PFKP*. Overall, we identified 24 greenlisted SNVs that were proximal to 10 genes
295 previously annotated as capable of regulating L1 expression. Our greenlisted SNVs also
296 captured genes involved in regulating L1 transposition, though there was no significant
297 enrichment (FDR = 7.78E-1) (**Figure 2B, 2D, 2E**). A few examples of these associations
298 included rs75237296 in the *PABPC1* 3'UTR, rs1288384419 in an intron of *RAD51B*, and
299 rs1471205623 upstream of *MPHOSPH8*. Here, we identified 5 greenlisted SNVs near 3
300 genes capable of regulating L1 transposition. Importantly, *PABPC1* is a poly(A) binding
301 protein that attaches to the poly(A) tail of L1, is important for the formation of L1
302 ribonucleoprotein particles (RNPs), and modules L1 and Alu transposition [58-60].
303 Additionally, *MPHOSPH8* is a component of the human silencing hub (HUSH) complex
304 and is important for L1 repression, regulating both L1 expression and L1 transposition
305 [34, 38, 61-63]. Finally, we checked the abundances of SNVs linked to candidate
306 regulators of L1 RNA levels in lymphoblastoid cells [42]. 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 [64].
324 Though *ADARB2* cannot catalyze A-to-I editing, it can negatively regulate the editing
325 functions of other ADARs [64], 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 copy 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**). Specifically, Alu retrotransposons can be segregated into the
337 old and inactive AluJ lineage, the middle-aged and active AluS lineage, and the young
338 and active AluY lineage [8]. Interestingly, though there was no significant enrichment of

339 either AluJ- or AluY-overlapping SNVs (FDR = 0.76 and FDR = 1, respectively), SNVs
340 overlapping AluS copies were significantly (FDR = 4.80E-6) enriched in our SNV
341 greenlist compared to background (**Figure 3B**). Similarly, L1 retrotransposons can be
342 segregated into the old L1M lineage, the middle-aged L1P lineage, and the young L1PA
343 lineage. Importantly, the L1PA1/L1Hs subfamily within the L1PA lineage is the only
344 actively mobile and autonomous subfamily within the human genome. Our enrichment
345 analysis highlighted a significant (FDR = 7.29E-7) depletion of L1M-overlapping SNVs,
346 a trending (FDR = 0.0786) enrichment of L1P-overlapping SNVs, and a significant (FDR
347 = 3.46E-12) enrichment of L1PA-overlapping SNVs, all compared to background SNVs
348 (**Figure 3B**). To obtain a higher resolution view of the transposition capabilities of
349 overlapping L1PA copies, we checked the overlap of our greenlisted SNVs with L1
350 annotations on L1Base v2—a dedicated database of putatively active L1 insertions [65].
351 Surprisingly, active copies, with either an intact, full-length L1 or only an intact ORF2,
352 were not significantly enriched/depleted in our greenlisted SNV list (**Figure 3C**).
353 However, non-intact, full-length L1 copies—annotated for their regulatory potential—
354 were significantly (FDR = 2.41E-6) enriched in our greenlisted SNV list compared to
355 background. Though mutant L1s can be mobilized at very low frequencies by
356 transposition-competent L1s *in trans* [9], these results suggest that most greenlisted
357 SNV-overlapping L1 copies are limited in their ability to generate *de novo* insertions.
358 This is potentially in contrast to greenlisted SNV-overlapping AluS copies, which may
359 still be measurably active in the human genome. We must also consider the alternative
360 interpretation that overlapping transposons may be directly involved in copy number
361 variation through transposition-independent genomic remodeling, through processes
362 such as repeat-mediated deletions and NAHR [12-15]. More generally, these results are
363 consistent with the possibility that greenlisted SNVs tag reference Alu and L1 insertions
364 that may contribute to TE copy number variation through transposition-dependent or
365 transposition-independent mechanisms.

366

367 The TE enrichments we identified above were consistent with those identified
368 using the Transposable Element Enrichment Analyzer (TEENA) [66] (**Table S1G**),
369 which has the added advantage of analyzing other TE families, in addition to Alu and

370 L1, at subfamily-level resolution. The most significantly (FDR < 0.05) enriched Alu
371 subfamilies included AluSg, AluSx3, AluSc, AluSx4, and AluSz6. Likewise, the most
372 significantly enriched L1 subfamilies included L1PA3 and L1PA4. Unexpectedly, the
373 most significantly enriched TE subfamily was HERVH-int of the ERV1 family.
374 Interestingly, human specific endogenous retrovirus H (HERVH) is essential for
375 maintenance of pluripotency in human stem cells (discussed in [67]). Finally, we
376 observed a significant enrichment of other ERV1 subfamilies (PABL_A-int), ERVL
377 subfamilies (HERVL-int), ERVK subfamilies (HERVK9-int, LTR13), and ERVL-MaLR
378 subfamilies (THE1B-int). These results re-iterate the association between specific
379 reference TE loci and variation in the number of L1/Alu singleton insertions.

380

381 We further explored the more general possibility that greenlisted SNVs tagged
382 genomic regions that may be prone to transposition-independent structural alterations
383 that may influence TE copy numbers. Such structural alterations may be facilitated by,
384 but may not require, the presence of repetitive elements. In particular, extensive
385 homology between segmental duplications, often in the vicinity of Alu elements [18], can
386 facilitate NAHR and drive recurrent genomic rearrangements [68] that can help form SV
387 hotspots [69, 70]. Noting the enrichment of AluS copies we observed among our
388 greenlisted SNVs, we next assessed whether our greenlisted SNVs significantly
389 overlapped regions of segmental duplication [71, 72] or regions characterized as SV
390 hotspots [70] (**Figure 3D**). Consistent with the notion that SNVs may tag regions with
391 potentially elevated rates of genome instability, our SNVs were very significantly
392 enriched in regions of segmental duplication (FDR = 2.64E-99), as well as in regions
393 harboring SV hotspots (FDR = 8.82E-7). These results further link variation in TE copy
394 number to genomic loci where structural instability may arise through transposition-
395 independent mechanisms.

396

397 Finally, we note that our GWAS analysis identified 21 polymorphic SVs that were
398 significantly (FDR < 0.05) associated with the presence/absence of L1/Alu insertion
399 singletons. These polymorphic SVs varied in nature and included inversions, Alu
400 insertions, an L1 insertion, SINE-VNTR-Alu (SVA) insertions, and a multiallelic copy

401 number variant (**Figure S4A**). With the exception of the CN0 copy number variant which
402 was associated with lower odds of carrying an L1/Alu insertion singleton (odds ratio =
403 0.28), all of the other structural variants were associated with higher odds of carrying an
404 insertion singleton (odds ratio > 1). Since the sequences for these SVs were not
405 available, it is unclear whether common, polymorphic L1/Alu insertion SVs may be
406 directly increasing the singleton number through novel transposition events, or whether
407 any of these polymorphic SVs may be influencing the L1/Alu singleton number through
408 transposition-independent mechanisms. Indeed, polymorphic inversions, many of which
409 are often flanked by retrotransposons, are associated with genetic instability [73].
410 Ultimately, these results suggest a tight link between common, polymorphic SVs of
411 different types and L1/Alu singleton SVs, whereby having the former is generally
412 associated with higher odds of having the latter.

413

414 *2.4 Case samples exhibit elevated cell cycle-related gene expression profiles*

415 To gain insight into the functional differences between controls and cases, we
416 leveraged publicly available lymphoblastoid cell line mRNA-seq data generated by the
417 GEUVADIS consortium for a subset of European and African samples in the 1000
418 Genomes Project [74] (**Figure 4A**). This included 358 European samples from 4
419 populations (British, Finnish, Tuscan, and Utah residents with European ancestry) and
420 86 African samples from 1 population (Yoruba), which we recently used to quantify gene
421 and TE expression profiles [42]. We utilized this expression data to construct consensus
422 gene co-expression networks for both the European and African samples using the
423 WGCNA [75] package. This approach led to the identification of 20 consensus modules
424 and 1 module (MEgrey) containing genes that were not assigned to the consensus
425 modules (**Table S1H**). We then ran a module-trait correlation analysis comparing the
426 expression of these modules with the case/control status of the European and African
427 samples (**Figure 4B**). Here, we used a stricter threshold of $p < 0.01$ to call significant
428 correlations. We were not able to identify any significant module-phenotype correlations
429 using the European network, which is potentially consistent with our difficulty in calling
430 significant GWAS variants in this super-population at the available sample sizes (**Figure**
431 **S2C**). In contrast, the MEroyalblue module was significantly ($p = 4.0E-4$) correlated with

432 African case/control status. To combine the results from each network, we utilized
433 Fisher's method to meta-analyze the p-values for modules exhibiting correlations in the
434 same direction. By meta-analysis, the MEroyalblue module was still significantly ($p =$
435 0.002) and positively correlated with case status. Finally, to functionally characterize this
436 module, we ran over-representation analysis using the GO Biological Process gene set
437 collection (**Figure 4C, Table S1I**). The top 10 over-represented gene sets were involved
438 in cell cycle-related processes, including "mitotic cell cycle", "cell division", and "sister
439 chromatid segregation". These findings are consistent with the biology of TE copy
440 number expansion. Though L1 can mobilize in non-dividing cells [76, 77], L1
441 retrotransposition exhibits a cell cycle bias and peaks during the S phase [78].
442 Alternatively, chromosome segregation errors during mitosis or meiosis can generate
443 cells with abnormal ploidy and either increased or decreased dosages of both genic and
444 transposon content [79]. These results implicate cell cycle differences in cells from
445 individuals with unique L1/Alu insertion singleton variation.

446

447

448 3. DISCUSSION

449 3.1 A new approach to identify loci implicated in L1 and Alu copy number variation

450 In this work, we developed a pipeline to computationally identify candidate loci
451 involved in L1/Alu singleton number variation by GWAS analysis. Importantly, our study
452 incorporates natural human genetic variation present in populations of different
453 geographic origin via trans-ethnic GWAS meta-analysis to identify shared, candidate
454 regulatory loci. Though several studies have begun to screen for regulators and
455 potential regulators of L1 expression or transposition in cell culture models or across
456 tissues [34-41], these can be limited in that the generalizability of these findings to
457 different ethnic populations is unclear. Moreover, no systematic, genome-wide screen
458 for candidate regulators of Alu activity has been carried out thus far, to our knowledge.
459 To address these gaps, we previously utilized trans-eQTL analysis to identify potential
460 regulators of L1 RNA levels in European and African populations [42]. Here, we utilized
461 genomic data from samples originating from 5 super-populations to identify candidate
462 loci modulating L1/Alu insertion singleton numbers.

463
464 TE copy number variation can arise through *de novo* transposition events or
465 through transposition-independent mechanisms, including recombination-based
466 mechanisms that can generate large deletions or duplications. We were particularly
467 interested in identifying new candidate regulators of L1/Alu transposition. Consistent
468 with the notion that greenlisted SNVs may play roles in the retrotransposon lifecycle, our
469 approach enriched genomic regions containing genes that can regulate L1 expression
470 levels. Though other known regulators of TE activity, and pathways involved in TE
471 control, were not significantly enriched among our greenlisted SNVs, we nonetheless
472 identified many SNVs in genomic regions containing these genes. This included, for
473 example, *MPHOSPH8*—a component of the HUSH complex important for L1
474 repression, regulating both L1 expression and transposition [34, 38, 61-63]. As another
475 example, we identified variants near *ADARB2*, a negative regulator of RNA A-to-I
476 editing, including among Alu RNAs [64]. These results suggest that SNV-associated
477 genes identified in this study hold TE regulatory potential and it may therefore be

478 informative to (i) test and validate these in future studies or (ii) use these to prioritize
479 future, targeted studies of TE regulators.

480

481 Our approach also identified an enrichment of greenlisted SNVs in regions
482 containing reference TE insertions, including AluS and full-length, non-intact L1PA
483 copies. Though neither of these can mobilize autonomously, they can hijack machinery
484 from transposition-competent L1s and mobilize *in trans* [6, 7, 9]. Thus, it is possible that
485 greenlisted SNVs tag reference insertions contributing to L1/Alu singleton variation
486 through transposition-dependent mechanisms. Of course, an alternative possibility is
487 that these repetitive elements are directly involved in genomic remodeling involving
488 transposition-independent mechanisms like repeat-mediated deletions or NAHR. We
489 also note that greenlisted SNVs were enriched in regions containing segmental
490 duplications and structural variation hotspots where recombination-based mechanisms,
491 including NAHR, may lead to duplications or deletions of the local genomic architecture.
492 Thus, it is also possible that greenlisted SNVs tag genomic regions prone to structural
493 variation that can alter the L1/Alu copy number through recombination-dependent
494 mechanisms.

495

496 Importantly, we also suggest the possibility that genome remodeling mechanisms
497 (including recombination) may interact with gene-based mechanisms of TE regulation.
498 Indeed, genes such as *BRCA1* are known to regulate L1 transposition [34] and are also
499 known to undergo Alu-Alu recombination events that can give rise to new mutations in
500 the gene [80-82]. Observations such as these highlight the possibility that TE insertions
501 may modulate structural variation in genomic regions containing genes regulating
502 retrotransposon lifecycles, which may facilitate an expansion of TE copy numbers
503 through transposition-based mechanisms, which may influence further structural
504 variation driving this whole process. This possibility is consistent with the enrichment of
505 greenlisted SNVs in regions containing L1 expression regulators, Alu and L1 repeats,
506 and other genomically unstable features like segmental duplications and structural
507 variation hotspots. In the future, it may be informative to experimentally assess this

508 possibility in contexts where genome instability is a hallmark feature that is coupled with
509 TE de-repression, such as aging [22] or aging-associated diseases like cancer [83, 84].

510

511 Finally, our approach also identified several common, polymorphic SVs that were
512 significantly associated with L1/Alu insertion singletons. Overwhelmingly, the presence
513 of polymorphic SVs of different types—inversions, Alu insertions, an L1 insertion, and
514 SVA insertions—was associated with increased odds of a global L1/Alu insertion
515 singleton. The exception to this was a multi-allelic copy number variant where 0 copies
516 were present; this SV was associated with decreased odds of a global L1/Alu insertion
517 singleton. Based on these results, we speculate that specific polymorphic SVs (i) may
518 directly drive genome instability that can facilitate the acquisition of L1/Alu copies or (ii)
519 may serve as markers for elevated risk of indirectly acquiring additional L1/Alu copies.
520 Indeed, active donor L1 copies that can mobilize and generate *de novo* insertions are
521 usually highly polymorphic in human populations (reviewed in [57]), and polymorphic
522 inversions, many of which are often flanked by retrotransposons, are also associated
523 with genetic instability and genomic disorders [73].

524

525 In summation, this study provides a list of variants that are associated with L1/Alu
526 insertion singletons and includes (i) SNVs in regions containing regulators of TE activity,
527 (ii) SNVs in regions containing features associated with genome instability, including
528 retrotransposons, that may influence TE copy number variation through transposition-
529 dependent or transposition-independent mechanisms, and (iii) common, polymorphic
530 SVs that may also influence TE copy number variation through transposition-dependent
531 or transposition-independent mechanisms.

532

533

534 *3.2 L1/Alu insertion singleton-associated loci contain genes of clinical relevance*

535 The most significant, greenlisted variant we identified was a polymorphic
536 inversion SV (INV_delly_INV00066128, odds ratio = 4.38, FDR = 3.42E-17,
537 chr21:26001780) residing in an intronic Alu copy within the *APP* gene, an important
538 marker of Alzheimer's disease (AD). AD is characterized by (i) the accumulation of

539 amyloid β ($A\beta$) plaques derived from amyloidogenic *APP* processing and (ii)
540 neurofibrillary tangles of hyperphosphorylated tau [85]. Importantly, tau protein can
541 induce TE expression and there is speculation that TEs may mobilize in tauopathies
542 [85]. Whether *APP* protein or $A\beta$ plaques can similarly modulate TE expression or
543 mobilization is an interesting area of potential future research; indeed, our findings are
544 consistent with the possibility that *APP* products may act as regulators of TE activity. Of
545 course, another possibility is that the genomic region containing *APP* may be a source
546 of L1/Alu copy number variation independent of the functional properties of *APP* protein
547 (i.e. genomic instability at that locus may be the driver of TE copy number variation).
548 Nevertheless, our results offer another connection between transposable element
549 regulation and Alzheimer's disease.

550

551 We also identified several SNVs proximal to *PWRN1*, which resides within the
552 Prader-Willi syndrome region and is thought to play a role in PWS [86]. Prader-Willi is
553 an imprinting disorder where genes in the chromosome 15q11-q13 region are
554 maternally imprinted and paternal copies are not expressed [87]. This lack of paternal
555 gene expression is predominantly caused by *de novo* paternally inherited deletions of
556 the 15q11-q13 region, though, less frequently, inheritance of two maternal chromosome
557 15 copies is the cause [87]. Importantly, a feature of genomic disorders like PWS is the
558 presence of segmental duplications that can serve as substrates for NAHR [88, 89].
559 Thus, we hypothesize that this particular region might be more prone to L1/Alu copy
560 number variation as a consequence of recombination-based chromosomal alterations.
561 Nevertheless, it is unclear (i) whether L1/Alu repeats are differentially active in PWS
562 compared to healthy controls or, more specifically, (ii) whether genes like *PWRN1* can
563 differentially express or mobilize L1/Alu transposons. Ultimately, the associations
564 between L1/Alu singletons and both *APP* and *PWRN1* further implicate
565 retrotransposons in disease.

566

567

568 *3.3 Limitations and future considerations*

569 In this study, we sought to identify new, candidate genes implicated in Alu and L1
570 copy number control. One specific mechanism of interest by which this can occur is
571 target-primed reverse transcription (TPRT)-mediated transposition. L1 insertions
572 generated by this method usually carry specific features, including short target site
573 duplications (TSDs), a polyadenine (polyA) tail, and are found integrated at an L1
574 endonuclease motif (reviewed in [57]). Since insertion sequences for 1000 Genomes
575 Project samples were not available, to our knowledge, it is difficult to assess to what
576 degree TPRT is driving the associations we identified. In studies with larger cohorts
577 where insertion sequences are available and insertions with TPRT features can be
578 identified, our approach could theoretically be applied to explore the genetic basis of
579 TPRT-specific copy number variation. Of course, our approach is generally restricted to
580 identifying genomic loci where variation is common across human populations. We
581 note, however, that significant variants were enriched in regions containing genes
582 involved in L1 expression regulation. Since expression is one of the early steps of the
583 L1 life cycle, our approach captured variants and genes with potential significance to
584 TPRT-mediated transposition.

585
586 We also note that there are several variables that we are unable to control for in
587 this study. To protect patient privacy, biological covariates such as chronological age
588 were not available and therefore could not be corrected for in our analysis. Since
589 increases in L1 expression and copy number have been observed with chronological
590 aging [23], differences in copy number may reflect age differences rather than genetic
591 differences. Importantly, however, samples were considered healthy at the time of
592 sample collection, potentially mitigating health-related effects on copy number. The
593 origin and developmental timing of the rare Alu and L1 insertions we utilized is also
594 unclear. That is, it is unclear whether global singletons used in this study arose *de novo*
595 in the germline, arose during life through somatic mutation, or even whether they just
596 arose during the cultivation of the lymphoblastoid cells used to amplify each sample's
597 DNA. Depending on when these insertions were acquired, the associations identified in
598 this study may be relevant for either germ cell or somatic cell TE biology. A potential
599 avenue of future research to address this question would be the incorporation of trio

600 parent-child genome sequencing and multi-generational genome sequencing to help
601 identify *bona fide de novo* insertions and their developmental timing. Ultimately, future
602 studies modulating genes identified with our approach will need to be carried out to
603 validate any causal contributions to L1/Alu regulation.

604

605

606 *3.4 Conclusions*

607 In this study, we employed GWAS across human populations of different
608 geographic origin to computationally identify genomic loci associated with variation in
609 L1/Alu insertion singleton number, specifically, and L1/Alu-associated copy number
610 variation, more generally. Our approach enriches for SNVs in genomic regions
611 containing known regulators of L1 expression. This observation suggests that the TE-
612 regulatory properties of these genes may extend beyond isogenic cell culture models to
613 more diverse human populations. Moreover, this observation also suggests that our list
614 of associated genes likely contains novel regulators of L1 or Alu activity that may be
615 prioritized in future, validation studies. Our approach also identified reference insertions
616 and non-reference, polymorphic SVs that may modulate L1/Alu copy numbers through
617 transposition-dependent or transposition-independent mechanisms. Finally, the
618 observation that some significant variants reside in genes of clinical relevance, like *APP*
619 and *PWRN1*, reinforce accumulating evidence of biological associations between TE
620 regulation and disease. Ultimately, our approach adds to the analytical toolkit that can
621 be used to study the regulation of TE activity.

622

623

624 **4. METHODS**

625 **4.1 Publicly available genomic data acquisition**

626 The multi-ethnic GWAS analysis was carried out on 2503 individuals derived
627 from 5 super-populations (African, East Asian, European, South Asian, and American)
628 and for which paired single nucleotide variant and structural variant data were available
629 from Phase 3 of the 1000 Genomes Project [90-92]. Specifically, Phase 3 autosomal
630 SNVs called on the GRCh38 reference genome were obtained from The International
631 Genome Sample Resource (IGSR) FTP site
632 ([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/)
633 [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
634 GRCh37 reference genome and then lifted over to GRCh38, were also obtained from
635 the IGSR FTP site
636 ([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/)
637 [38_positions/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/supporting/GRCh38_positions/)).

638
639 Human gene expression data across 54 non-diseased tissue sites was obtained
640 from the GTEx Analysis v8 on the GTEx Portal [54-56]. Specifically, we downloaded the
641 matrix containing the median gene-level transcripts per million (TPMs) by tissue, and we
642 extracted the expression values for significant SNV-associated genes. After filtering out
643 genes with no detectable expression (0 TPMs), we generated heatmaps using the
644 pheatmap v1.0.12 package in R. Gene expression values were centered and scaled
645 across tissues to visualize and compare the relative expression levels across tissues.

646

647 **4.2 Aggregating and pre-processing genotype data for GWAS analysis**

648 To define the phenotype of interest for GWAS, we first extracted global singleton
649 Alu and L1 insertions. We utilized VCFtools v0.1.17 [93] to extract all autosomal SVs
650 with no missing data (--max-missing 1) and an allele count of 1 across all samples (--
651 non-ref-ac 1 --max-non-ref-ac 1), i.e. global singletons. From these, we extracted Alu-
652 and L1-specific insertions using BCFtools v1.10.2 [94] to keep entries annotated with
653 the 'SVTYPE="LINE1"' and 'SVTYPE="ALU"' tags. Finally, VCF files containing global

654 singleton L1 or Alu insertions were converted to PLINK BED format using PLINK
655 v1.90b6.17 [95].

656

657 We note that SVs on sex chromosomes were not included in any part of the
658 analysis since (i) Y chromosome SVs were not available, (ii) male genotypes on
659 chromosome X were unknown, and (iii) association studies with X chromosome variants
660 require unique algorithms and cannot easily be incorporated into traditional association
661 pipelines [96, 97].

662

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

670

671 Third, we prepared polymorphic SNVs for inclusion in the GWAS analysis. All
672 SNVs were first annotated with rsIDs from dbSNP build 155 using BCFtools. Within
673 each super-population, VCFtools was used to remove indels and keep autosomal SNVs
674 with the same parameters as the polymorphic SVs. We note that for similar reasons as
675 with the polymorphic SVs, sex chromosome SNVs were also omitted from all analyses.
676 We then used BCFtools to identify and subset SNVs that were shared across all 5
677 super-populations. Finally, we used BCFtools to generate the final genotype matrices by
678 combining shared, polymorphic SNVs with shared, polymorphic SVs. VCF files
679 containing the combined SNVs and SVs were then converted to PLINK BED format
680 using PLINK, keeping the allele order. PLINK was also used to prune the combined
681 SNV and SV matrices (--indep-pairwise 50 10 0.1) and to generate principal
682 components (PCs) from the pruned genotypes, for inclusion as covariates in the GWAS.

683

684 **4.3 Super-population-specific and trans-ethnic GWAS**

685 We began by running GWAS within each super-population using PLINK
686 v1.90b6.17 [95]. For the phenotype, we added the number of Alu and L1 global
687 singleton insertions for each sample and segregated samples into cases and controls,
688 depending on whether they contained or did not contain a global singleton. We ran
689 GWAS analyses using a logistic regression model that included the following covariates:
690 biological sex and the top 4 principal components from the pruned SNV and SV
691 genotype matrices. Individual results from each super-population were combined via
692 meta-analysis using PLINK. To help call significant variants, we generated a null
693 distribution of p-values for each super-population by running 20 instances of the GWAS
694 where the case-control status for each sample was randomly shuffled with the case-
695 control status of a different sample. Each set of permutation results was meta-analyzed
696 across super-populations to similarly obtain 20 random distributions of meta-analysis p-
697 values. For the meta-analysis, we focused on the p-values and odds ratios generated
698 using a random effects statistical model, as opposed to a fixed effects model, since 1)
699 there may be heterogeneity across super-populations in response to different genetic
700 variants, and 2) we were interested in enhancing the generalizability of our findings to
701 facilitate future follow-up studies.

702
703 To limit false positives, the Benjamini-Hochberg (BH) false discovery rate (FDR)
704 was calculated in each analysis, and we used the p-value corresponding to a BH FDR <
705 0.05 as the threshold for GWAS significance. As a secondary threshold, we used the
706 permutation data to identify p-values corresponding to an average empirical FDR <
707 0.05. To note, we calculated the average empirical FDR at a given p-value p_i by (i)
708 counting the total number of null points with $p \leq p_i$, (ii) dividing by the number of
709 permutations, to obtain an average number of null points with $p \leq p_i$, and (iii) dividing the
710 average number of null points with $p \leq p_i$ by the number of real points with $p \leq p_i$. GWAS
711 variants were considered significant if they passed the stricter of the two thresholds in
712 each analysis.

713

714 **4.4 Annotation of variants and annotation enrichment analyses**

715 We obtained BED files containing annotated genomic regions from various
716 sources. We obtained the ENCODE blacklist v2 [46] for hg38 from
717 <https://github.com/Boyle-Lab/Blacklist/tree/master/lists>. Segmental duplications [71, 72]
718 and RepeatMasker annotations using the Rebase library [98] were obtained from the
719 UCSC Genome Browser [99]. We obtained SV hotspot coordinates on hg19 from [70],
720 and we used the online UCSC LiftOver tool to map coordinates to the hg38 genome
721 assembly using the default settings. The BED tracks for full-length and intact L1s, only
722 ORF2-intact L1s, and full-length non-intact L1s were obtained from L1Base v2 [65]. We
723 obtained the Registry (version 4) of candidate cis-Regulatory Elements (cCREs) for
724 hg38 from the Search Candidate Regulatory Elements by ENCODE (SCREEN) web
725 interface [53] (<http://screen-beta.wenglab.org>). We used the ‘intersect’ command in
726 BEDTools v2.31.1 [100] to assign genomic region annotations to all overlapping
727 variants used in this study. We were also interested in assessing whether variants were
728 linked to specific regulatory gene annotations. All variants used in the study were
729 submitted to the GREAT v4.0.4 [47, 48] online platform with the default settings (basal
730 plus extension, proximal with 5 kb upstream and 1 kb downstream, plus distal up to
731 1000 kb) to assign gene annotations to each variant. These gene annotations were then
732 used to assess the number of variants linked to genes in several TE regulatory lists—
733 including a list of genes that regulated L1 expression in a CRISPR screen using cancer
734 cells [38], a list of genes that regulated L1 transposition in an independent CRISPR
735 screen using cancer cells [34], a list of candidate genes influencing intronic, intergenic,
736 or exonic L1 RNA levels in lymphoblastoid cell lines [42], a list of genes with histone
737 methyltransferase activity (GO:0042054), and a list of genes with RNA modification
738 activity (GO:0009451). The two GO gene sets were obtained on 2024-07-29 from the
739 Molecular Signatures Database (MSigDB) v2023.2.Hs [101, 102], corresponding to GO
740 release 2023-07-27.

741

742 Given the limited number of significant SVs and the unavailability of SV
743 sequences, we largely focused on blacklist-filtered, significant SNVs in downstream
744 enrichment analyses. For each of the above annotations, we assessed whether
745 greenlisted SNVs were significantly enriched or depleted for that annotation compared

746 to background SNVs—all SNVs that were tested in the GWAS. The numbers of
747 background and greenlisted SNVs overlapping or not overlapping a set of annotations
748 were placed into a contingency table, and statistical significance was assessed using
749 Fisher's exact test (with the `fisher.test` function in R v4.3.3). After all tests were carried
750 out, p-values were FDR corrected using the `p.adjust` function in R. All
751 enrichments/depletions with an FDR < 0.05 were considered significant. For the repeat
752 enrichment analyses, we also analyzed our greenlisted SNVs using the TEENA web
753 server [66] (on August 8, 2024), specifying the hg38 assembly and using all other
754 default options.

755

756 **4.5 RNA-seq and gene co-expression network analyses**

757 For lymphoblastoid cell line transcriptional analyses, mRNA-sequencing was
758 initially carried out by the GEUVADIS consortium [74] on LCLs from a small subset of
759 European and African (Yoruban, specifically) samples from the 1000 Genomes Project.
760 Recently, we re-processed this data to quantify gene and transposon expression levels
761 [42]. Briefly, reads were trimmed using `fastp v0.20.1` [103], trimmed reads were aligned
762 to the GRCh38 human genome assembly using `STAR v2.7.3a` [104], and the
763 `TEtranscripts v2.1.4` [105] package was used to obtain gene and TE counts, using the
764 GENCODE release 33 [106] annotations and a repeat GTF file provided on the
765 Hammell lab website. To note, the EBV genome (GenBank ID V01555.2) was included
766 as an additional contig in our reference genome, since LCLs are generated by infecting
767 B-cells with Epstein-Barr virus (EBV).

768

769 Using these gene/TE count matrices, lowly expressed genes were filtered out if
770 50% of European or Yoruban samples did not have over 0.44 counts per million (cpm)
771 or 0.43 cpm, respectively, which correspond to 10 reads in each cohort's median-length
772 library. Since we were interested in building consensus co-expression networks
773 between the European and Yoruban samples, we also removed genes that were not
774 expressed in both groups. After, the filtered counts underwent a variance stabilizing
775 transformation (`vst`) using `DESeq2 v1.42.1` [107] and the following covariates were
776 regressed out using the `'removeBatchEffect'` function in `Limma v3.58.1` [108]: biological

777 sex, sequencing lab, population category, principal components 1-2 of the pruned
778 genotype matrices containing both SNVs and SVs, and EBV expression levels. The
779 population category variable was omitted in the Yoruban batch correction since that did
780 not vary.

781

782 The batch-corrected VST matrices were then used to perform weighted gene co-
783 expression network analysis (WGCNA) [75] using the WGCNA v1.72-5 R package. We
784 used the ‘blockwiseConsensusModules’ function to automate consensus network
785 construction for both the European and African expression data, specifying these
786 parameters: corType = “bicor”, power = 12, networkType = "signed", maxPOutliers =
787 0.05, mergeCutHeight = 0.25, deepSplit = 2, minKMEtoStay = 0, pamRespectsDendro =
788 FALSE, minModuleSize = 30, and randomSeed = 90280. Phenotype-module
789 correlations, and the corresponding p-values, were calculated using WGCNA’s ‘cor’ and
790 ‘corPvalueFisher’ functions, respectively. The p-values for the European and Yoruban
791 correlations were meta-analyzed using Fisher’s method. For visualization purposes
792 only, to show a correlation direction in the meta-analysis, we took the average of the
793 European and Yoruban correlations. Modules showing opposite correlations across the
794 two consensus networks were disregarded in the meta-analysis. Correlations with a p-
795 value < 0.01 were considered significant.

796

797 **4.6 Functional enrichment analyses**

798 We used the over-representation analysis (ORA) paradigm as implemented in
799 the R package clusterProfiler v4.10.1 [49]. Gene Ontology Biological Process gene sets
800 were obtained from the R package msgdbr v7.5.1, an Ensembl ID-mapped collection of
801 gene sets from the Molecular Signatures Database [101, 102]. For ORA with genes
802 linked to greenlisted SNVs and SVs, we used the genes linked to the background SNVs
803 and SVs, respectively, for the universe background to compute enrichment significance.
804 For ORA analysis of co-expression network modules, we used all genes in the network
805 for the universe background. All gene sets with an FDR < 0.05 were considered
806 significant, and the top 10 significant gene sets, at most, were plotted. All enrichments
807 results are included in **Table S1C**, **S1D**, and **S1I**.

808

809

810 **DECLARATIONS**

811 **Ethics approval and consent to participate**

812 Not applicable.

813

814 **Consent for publication**

815 Not applicable.

816

817 **Availability of data and materials**

818 All code is available on the Benayoun lab GitHub
819 (https://github.com/BenayounLaboratory/TE_GWAS). Analyses were conducted using R
820 version 4.3.3. Code was re-run independently on R version 4.3.0 to check for
821 reproducibility.

822

823 **Competing interests**

824 The authors declare that they have no competing interests.

825

826 **Funding**

827 This work was supported by the National Science Foundation
828 [<https://www.nsf.gov/>] Graduate Research Fellowship Program (NSF GRFP) DGE-
829 1842487 (J.I.B.), the National Institute on Aging [<https://www.nia.nih.gov/>] T32
830 AG052374 (J.I.B.), the University of Southern California with a Provost Fellowship
831 (J.I.B.), and the National Institute of General Medical Sciences
832 [<https://www.nigms.nih.gov/>] R35 GM142395 (to B.A.B).

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

835

836 **Authors' contributions**

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

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

842

843 **Acknowledgements**

844 We would like to thank Prof. Rachel Brem at the University of California Berkeley
845 for her feedback and insights on the GWAS analyses. We would also like to thank Dr.
846 Heather C. Mefford at St. Jude Children's Research Hospital for referring us to
847 publications that helped shape the analyses linking our variants to regions potentially
848 prone to genomic instability and involved in disease. Finally, we would like to thank Dr.
849 Minhoo Kim and Mr. Aaron Lemus for their comments on the manuscript.

850

851 REFERENCES

- 852 1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K,
853 Doyle M, FitzHugh W, et al: **Initial sequencing and analysis of the human genome.**
854 *Nature* 2001, **409**:860-921.
- 855 2. Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, Smith HO, Yandell M,
856 Evans CA, Holt RA, et al: **The Sequence of the Human Genome.** *Science* 2001, **291**:1304-
857 1351.
- 858 3. Ostertag EM, Kazazian HH, Jr.: **Biology of mammalian L1 retrotransposons.** *Annu Rev*
859 *Genet* 2001, **35**:501-538.
- 860 4. Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran JV, Kazazian HH: **Hot**
861 **L1s account for the bulk of retrotransposition in the human population.** *Proceedings of*
862 *the National Academy of Sciences* 2003, **100**:5280-5285.
- 863 5. Rubin CM, Houck CM, Deininger PL, Friedmann T, Schmid CW: **Partial nucleotide**
864 **sequence of the 300-nucleotide interspersed repeated human DNA sequences.** *Nature*
865 1980, **284**:372-374.
- 866 6. Dewannieux M, Esnault C, Heidmann T: **LINE-mediated retrotransposition of marked**
867 **Alu sequences.** *Nature Genetics* 2003, **35**:41-48.
- 868 7. Wallace N, Wagstaff BJ, Deininger PL, Roy-Engel AM: **LINE-1 ORF1 protein enhances Alu**
869 **SINE retrotransposition.** *Gene* 2008, **419**:1-6.
- 870 8. Bennett EA, Keller H, Mills RE, Schmidt S, Moran JV, Weichenrieder O, Devine SE: **Active**
871 **Alu retrotransposons in the human genome.** *Genome Research* 2008, **18**:1875-1883.
- 872 9. Wei W, Gilbert N, Ooi SL, Lawler JF, Ostertag EM, Kazazian HH, Boeke JD, Moran JV:
873 **Human L1 Retrotransposition: cisPreference versus trans Complementation.** *Molecular*
874 *and Cellular Biology* 2001, **21**:1429-1439.
- 875 10. Mendez-Dorantes C, Bhargava R, Stark JM: **Repeat-mediated deletions can be induced**
876 **by a chromosomal break far from a repeat, but multiple pathways suppress such**
877 **rearrangements.** *Genes Dev* 2018, **32**:524-536.
- 878 11. Mendez-Dorantes C, Tsai LJ, Jahanshir E, Lopezcolorado FW, Stark JM: **BLM has Contrary**
879 **Effects on Repeat-Mediated Deletions, based on the Distance of DNA DSBs to a Repeat**
880 **and Repeat Divergence.** *Cell Reports* 2020, **30**:1342-1357.e1344.
- 881 12. Boone Philip M, Yuan B, Campbell Ian M, Scull Jennifer C, Withers Marjorie A, Baggett
882 Brett C, Beck Christine R, Shaw Christine J, Stankiewicz P, Moretti P, et al: **The**
883 **Alu-Rich Genomic Architecture of SPAST Predisposes to**
884 **Diverse and Functionally Distinct Disease-Associated CNV Alleles.** *The American*
885 *Journal of Human Genetics* 2014, **95**:143-161.
- 886 13. Campbell IM, Gambin T, Dittwald P, Beck CR, Shuvarikov A, Hixson P, Patel A, Gambin A,
887 Shaw CA, Rosenfeld JA, Stankiewicz P: **Human endogenous retroviral elements**
888 **promote genome instability via non-allelic homologous recombination.** *BMC Biology*
889 2014, **12**:74.
- 890 14. Startek M, Szafranski P, Gambin T, Campbell IM, Hixson P, Shaw CA, Stankiewicz P,
891 Gambin A: **Genome-wide analyses of LINE–LINE-mediated nonallelic homologous**
892 **recombination.** *Nucleic Acids Research* 2015, **43**:2188-2198.

- 893 15. Robberecht C, Voet T, Esteki MZ, Nowakowska BA, Vermeesch JR: **Nonallelic**
894 **homologous recombination between retrotransposable elements is a driver of de**
895 **novo unbalanced translocations.** *Genome Research* 2013, **23**:411-418.
- 896 16. Belancio VP, Deininger PL, Roy-Engel AM: **LINE dancing in the human genome:**
897 **transposable elements and disease.** *Genome Medicine* 2009, **1**:97.
- 898 17. Cordaux R, Batzer MA: **The impact of retrotransposons on human genome evolution.**
899 *Nature Reviews Genetics* 2009, **10**:691-703.
- 900 18. Bailey JA, Liu G, Eichler EE: **An *Alu* Transposition Model for the Origin and**
901 **Expansion of Human Segmental Duplications.** *The American Journal of Human Genetics*
902 2003, **73**:823-834.
- 903 19. Kolomietz E, Meyn MS, Pandita A, Squire JA: **The role of Alu repeat clusters as**
904 **mediators of recurrent chromosomal aberrations in tumors.** *Genes, Chromosomes and*
905 *Cancer* 2002, **35**:97-112.
- 906 20. Bravo JI, Nozownik S, Danthi PS, Benayoun BA: **Transposable elements, circular RNAs**
907 **and mitochondrial transcription in age-related genomic regulation.** *Development* 2020,
908 **147**.
- 909 21. Zhang X, Zhang R, Yu J: **New Understanding of the Relevant Role of LINE-1**
910 **Retrotransposition in Human Disease and Immune Modulation.** *Frontiers in Cell and*
911 *Developmental Biology* 2020, **8**.
- 912 22. López-Otín C, Blasco MA, Partridge L, Serrano M, Kroemer G: **Hallmarks of aging: An**
913 **expanding universe.** *Cell* 2023, **186**:243-278.
- 914 23. De Cecco M, Criscione SW, Peterson AL, Neretti N, Sedivy JM, Kreiling JA: **Transposable**
915 **elements become active and mobile in the genomes of aging mammalian somatic**
916 **tissues.** *Aging (Albany NY)* 2013, **5**:867-883.
- 917 24. De Cecco M, Criscione SW, Peckham EJ, Hillenmeyer S, Hamm EA, Manivannan J,
918 Peterson AL, Kreiling JA, Neretti N, Sedivy JM: **Genomes of replicatively senescent cells**
919 **undergo global epigenetic changes leading to gene silencing and activation of**
920 **transposable elements.** *Aging Cell* 2013, **12**:247-256.
- 921 25. Campisi J: **Aging, Cellular Senescence, and Cancer.** *Annual Review of Physiology* 2013,
922 **75**:685-705.
- 923 26. Franceschi C, Garagnani P, Parini P, Giuliani C, Santoro A: **Inflammaging: a new**
924 **immune–metabolic viewpoint for age-related diseases.** *Nature Reviews Endocrinology*
925 2018, **14**:576-590.
- 926 27. Wallace NA, Belancio VP, Deininger PL: **L1 mobile element expression causes multiple**
927 **types of toxicity.** *Gene* 2008, **419**:75-81.
- 928 28. Belancio VP, Roy-Engel AM, Pochampally RR, Deininger P: **Somatic expression of LINE-1**
929 **elements in human tissues.** *Nucleic acids research* 2010, **38**:3909-3922.
- 930 29. De Cecco M, Ito T, Petrashen AP, Elias AE, Skvir NJ, Criscione SW, Caligiana A, Broccoli G,
931 Adney EM, Boeke JD, et al: **L1 drives IFN in senescent cells and promotes age-**
932 **associated inflammation.** *Nature* 2019, **566**:73-78.
- 933 30. Yamada K, Kaneko H, Shimizu H, Suzumura A, Namba R, Takayama K, Ito S, Sugimoto M,
934 Terasaki H: **Lamivudine Inhibits Alu RNA-induced Retinal Pigment Epithelium**
935 **Degeneration via Anti-inflammatory and Anti-senescence Activities.** *Translational*
936 *Vision Science & Technology* 2020, **9**:1-1.

- 937 31. Wang J, Geesman GJ, Hostikka SL, Atallah M, Blackwell B, Lee E, Cook PJ, Pasaniuc B,
938 Shariat G, Halperin E, et al: **Inhibition of activated pericentromeric SINE/Alu repeat**
939 **transcription in senescent human adult stem cells reinstates self-renewal.** *Cell Cycle*
940 2011, **10**:3016-3030.
- 941 32. Levin HL, Moran JV: **Dynamic interactions between transposable elements and their**
942 **hosts.** *Nature Reviews Genetics* 2011, **12**:615-627.
- 943 33. Rebollo R, Romanish MT, Mager DL: **Transposable Elements: An Abundant and Natural**
944 **Source of Regulatory Sequences for Host Genes.** *Annual Review of Genetics* 2012,
945 **46**:21-42.
- 946 34. Liu N, Lee CH, Swigut T, Grow E, Gu B, Bassik MC, Wysocka J: **Selective silencing of**
947 **euchromatic L1s revealed by genome-wide screens for L1 regulators.** *Nature* 2018,
948 **553**:228-232.
- 949 35. Tristan-Ramos P, Morell S, Sanchez L, Toledo B, Garcia-Perez JL, Heras SR: **sRNA/L1**
950 **retrotransposition: using siRNAs and miRNAs to expand the applications of the cell**
951 **culture-based LINE-1 retrotransposition assay.** *Philosophical Transactions of the Royal*
952 *Society B: Biological Sciences* 2020, **375**:20190346.
- 953 36. Mita P, Sun X, Fenyö D, Kahler DJ, Li D, Agmon N, Wudzinska A, Keegan S, Bader JS, Yun
954 C, Boeke JD: **BRCA1 and S phase DNA repair pathways restrict LINE-1**
955 **retrotransposition in human cells.** *Nat Struct Mol Biol* 2020, **27**:179-191.
- 956 37. Briggs EM, Mita P, Sun X, Ha S, Vasilyev N, Leopold ZR, Nudler E, Boeke JD, Logan SK:
957 **Unbiased proteomic mapping of the LINE-1 promoter using CRISPR Cas9.** *Mobile DNA*
958 2021, **12**:21.
- 959 38. Li X, Bie L, Wang Y, Hong Y, Zhou Z, Fan Y, Yan X, Tao Y, Huang C, Zhang Y, et al: **LINE-1**
960 **transcription activates long-range gene expression.** *Nature Genetics* 2024, **56**:1494-
961 1502.
- 962 39. Sun X, Wang X, Tang Z, Grivainis M, Kahler D, Yun C, Mita P, Fenyö D, Boeke JD:
963 **Transcription factor profiling reveals molecular choreography and key regulators of**
964 **human retrotransposon expression.** *Proc Natl Acad Sci U S A* 2018, **115**:E5526-e5535.
- 965 40. Chung N, Jonaid GM, Quinton S, Ross A, Sexton CE, Alberto A, Clymer C, Churchill D,
966 Navarro Leija O, Han MV: **Transcriptome analyses of tumor-adjacent somatic tissues**
967 **reveal genes co-expressed with transposable elements.** *Mobile DNA* 2019, **10**:39.
- 968 41. Tristán-Ramos P, Rubio-Roldan A, Peris G, Sánchez L, Amador-Cubero S, Viollet S,
969 Cristofari G, Heras SR: **The tumor suppressor microRNA let-7 inhibits human LINE-1**
970 **retrotransposition.** *Nature Communications* 2020, **11**:5712.
- 971 42. Bravo JI, Mizrahi CR, Kim S, Zhang L, Suh Y, Benayoun BA: **An eQTL-based approach**
972 **reveals candidate regulators of LINE-1 RNA levels in lymphoblastoid cells.** *PLOS*
973 *Genetics* 2024, **20**:e1011311.
- 974 43. Johnston HR, Hu Y, Cutler DJ: **Population Genetics Identifies Challenges in Analyzing**
975 **Rare Variants.** *Genetic Epidemiology* 2015, **39**:145-148.
- 976 44. Evangelou E, Ioannidis JPA: **Meta-analysis methods for genome-wide association**
977 **studies and beyond.** *Nature Reviews Genetics* 2013, **14**:379-389.
- 978 45. Peterson RE, Kuchenbaecker K, Walters RK, Chen CY, Popejoy AB, Periyasamy S, Lam M,
979 Iyegbe C, Strawbridge RJ, Brick L, et al: **Genome-wide Association Studies in Ancestrally**

- 980 **Diverse Populations: Opportunities, Methods, Pitfalls, and Recommendations.** *Cell*
981 2019, **179**:589-603.
- 982 46. Amemiya HM, Kundaje A, Boyle AP: **The ENCODE Blacklist: Identification of Problematic**
983 **Regions of the Genome.** *Scientific Reports* 2019, **9**:9354.
- 984 47. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G:
985 **GREAT improves functional interpretation of cis-regulatory regions.** *Nature*
986 *Biotechnology* 2010, **28**:495-501.
- 987 48. Tanigawa Y, Dyer ES, Bejerano G: **WhichTF is functionally important in your open**
988 **chromatin data?** *PLOS Computational Biology* 2022, **18**:e1010378.
- 989 49. Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, Feng T, Zhou L, Tang W, Zhan L, et al:
990 **clusterProfiler 4.0: A universal enrichment tool for interpreting omics data.** *The*
991 *Innovation* 2021, **2**:100141.
- 992 50. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM: **A**
993 **program for annotating and predicting the effects of single nucleotide polymorphisms,**
994 **SnEff.** *Fly* 2012, **6**:80-92.
- 995 51. Liu Y, Yang Q, Zhao F: **Synonymous but Not Silent: The Codon Usage Code for Gene**
996 **Expression and Protein Folding.** *Annual Review of Biochemistry* 2021, **90**:375-401.
- 997 52. Hampel H, Hardy J, Blennow K, Chen C, Perry G, Kim SH, Villemagne VL, Aisen P,
998 Vendruscolo M, Iwatsubo T, et al: **The Amyloid- β Pathway in Alzheimer's Disease.**
999 *Molecular Psychiatry* 2021, **26**:5481-5503.
- 1000 53. Abascal F, Acosta R, Addleman NJ, Adrian J, Afzal V, Ai R, Aken B, Akiyama JA, Jammal
1001 OA, Amrhein H, et al: **Expanded encyclopaedias of DNA elements in the human and**
1002 **mouse genomes.** *Nature* 2020, **583**:699-710.
- 1003 54. Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, Hasz R, Walters G, Garcia F,
1004 Young N, et al: **The Genotype-Tissue Expression (GTEx) project.** *Nature Genetics* 2013,
1005 **45**:580-585.
- 1006 55. Carithers LJ, Ardlie K, Barcus M, Branton PA, Britton A, Buia SA, Compton CC, DeLuca DS,
1007 Peter-Demchok J, Gelfand ET, et al: **A Novel Approach to High-Quality Postmortem**
1008 **Tissue Procurement: The GTEx Project.** *Biopreservation and Biobanking* 2015, **13**:311-
1009 319.
- 1010 56. Consortium TG, Aguet F, Anand S, Ardlie KG, Gabriel S, Getz GA, Graubert A, Hadley K,
1011 Handsaker RE, Huang KH, et al: **The GTEx Consortium atlas of genetic regulatory effects**
1012 **across human tissues.** *Science* 2020, **369**:1318-1330.
- 1013 57. Faulkner GJ, Billon V: **L1 retrotransposition in the soma: a field jumping ahead.** *Mobile*
1014 *DNA* 2018, **9**:22.
- 1015 58. Dai L, Taylor MS, O'Donnell KA, Boeke JD: **Poly(A) binding protein C1 is essential for**
1016 **efficient L1 retrotransposition and affects L1 RNP formation.** *Mol Cell Biol* 2012,
1017 **32**:4323-4336.
- 1018 59. Taylor Martin S, LaCava J, Mita P, Molloy Kelly R, Huang Cheng Ran L, Li D, Adney
1019 Emily M, Jiang H, Burns Kathleen H, Chait Brian T, et al: **Affinity Proteomics Reveals**
1020 **Human Host Factors Implicated in Discrete Stages of LINE-1 Retrotransposition.** *Cell*
1021 2013, **155**:1034-1048.

- 1022 60. Goodier JL, Cheung LE, Kazazian HH, Jr: **Mapping the LINE1 ORF1 protein interactome**
1023 **reveals associated inhibitors of human retrotransposition.** *Nucleic Acids Research*
1024 2013, **41**:7401-7419.
- 1025 61. Gu Z, Liu Y, Zhang Y, Cao H, Lyu J, Wang X, Wylie A, Newkirk SJ, Jones AE, Lee M, et al:
1026 **Silencing of LINE-1 retrotransposons is a selective dependency of myeloid leukemia.**
1027 *Nature Genetics* 2021, **53**:672-682.
- 1028 62. Müller I, Moroni AS, Shlyueva D, Sahadevan S, Schoof EM, Radzsheuskaya A, Højfeldt
1029 JW, Tatar T, Koche RP, Huang C, Helin K: **MPP8 is essential for sustaining self-renewal**
1030 **of ground-state pluripotent stem cells.** *Nature Communications* 2021, **12**:3034.
- 1031 63. Danac JMC, Matthews RE, Gungi A, Qin C, Parsons H, Antrobus R, Timms RT,
1032 Tchasovnikarova IA: **Competition between two HUSH complexes orchestrates the**
1033 **immune response to retroelement invasion.** *Molecular Cell* 2024, **84**:2870-2881.e2875.
- 1034 64. Raghava Kurup R, Oakes EK, Manning AC, Mukherjee P, Vadlamani P, Hundley HA: **RNA**
1035 **binding by ADAR3 inhibits adenosine-to-inosine editing and promotes expression of**
1036 **immune response protein MAVS.** *Journal of Biological Chemistry* 2022, **298**.
- 1037 65. Penzkofer T, Jäger M, Figlerowicz M, Badge R, Mundlos S, Robinson PN, Zemojtel T:
1038 **L1Base 2: more retrotransposition-active LINE-1s, more mammalian genomes.** *Nucleic*
1039 *Acids Research* 2016, **45**:D68-D73.
- 1040 66. Li Y, Lyu R, Chen S, Wang Y, Sun M-a: **TEENA: an integrated web server for transposable**
1041 **element enrichment analysis in various model and non-model organisms.** *Nucleic Acids*
1042 *Research* 2024, **52**:W126-W131.
- 1043 67. Sexton CE, Tillett RL, Han MV: **The essential but enigmatic regulatory role of HERVH in**
1044 **pluripotency.** *Trends in Genetics* 2022, **38**:12-21.
- 1045 68. Sharp AJ, Locke DP, McGrath SD, Cheng Z, Bailey JA, Vallente RU, Pertz LM, Clark RA,
1046 Schwartz S, Seagraves R, et al: **Segmental Duplications and Copy-Number Variation in**
1047 **the Human Genome.** *The American Journal of Human Genetics* 2005, **77**:78-88.
- 1048 69. Perry GH, Tchinda J, McGrath SD, Zhang J, Picker SR, Cáceres AM, Iafrate AJ, Tyler-Smith
1049 C, Scherer SW, Eichler EE, et al: **Hotspots for copy number variation in chimpanzees**
1050 **and humans.** *Proc Natl Acad Sci U S A* 2006, **103**:8006-8011.
- 1051 70. Lin Y-L, Gokcumen O: **Fine-Scale Characterization of Genomic Structural Variation in**
1052 **the Human Genome Reveals Adaptive and Biomedically Relevant Hotspots.** *Genome*
1053 *Biology and Evolution* 2019, **11**:1136-1151.
- 1054 71. Bailey JA, Yavor AM, Massa HF, Trask BJ, Eichler EE: **Segmental Duplications:**
1055 **Organization and Impact Within the Current Human Genome Project Assembly.**
1056 *Genome Research* 2001, **11**:1005-1017.
- 1057 72. Bailey JA, Gu Z, Clark RA, Reinert K, Samonte RV, Schwartz S, Adams MD, Myers EW, Li
1058 PW, Eichler EE: **Recent Segmental Duplications in the Human Genome.** *Science* 2002,
1059 **297**:1003-1007.
- 1060 73. Porubsky D, Höps W, Ashraf H, Hsieh P, Rodriguez-Martin B, Yilmaz F, Ebler J, Hallast P,
1061 Maria Maggiolini FA, Harvey WT, et al: **Recurrent inversion polymorphisms in humans**
1062 **associate with genetic instability and genomic disorders.** *Cell* 2022, **185**:1986-
1063 2005.e1926.
- 1064 74. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas MA,
1065 González-Porta M, Kurbatova N, Griebel T, Ferreira PG, et al: **Transcriptome and**

- 1066 **genome sequencing uncovers functional variation in humans.** *Nature* 2013, **501**:506-
1067 511.
- 1068 75. Langfelder P, Horvath S: **WGCNA: an R package for weighted correlation network**
1069 **analysis.** *BMC Bioinformatics* 2008, **9**:559.
- 1070 76. Kubo S, Seleme MdC, Soifer HS, Perez JLG, Moran JV, Kazazian HH, Kasahara N: **L1**
1071 **retrotransposition in nondividing and primary human somatic cells.** *Proceedings of the*
1072 *National Academy of Sciences* 2006, **103**:8036-8041.
- 1073 77. Macia A, Widmann TJ, Heras SR, Ayllon V, Sanchez L, Benkaddour-Boumzaouad M,
1074 Muñoz-Lopez M, Rubio A, Amador-Cubero S, Blanco-Jimenez E, et al: **Engineered LINE-1**
1075 **retrotransposition in nondividing human neurons.** *Genome Research* 2017, **27**:335-348.
- 1076 78. Mita P, Wudzinska A, Sun X, Andrade J, Nayak S, Kahler DJ, Badri S, LaCava J, Ueberheide
1077 B, Yun CY, et al: **LINE-1 protein localization and functional dynamics during the cell**
1078 **cycle.** *eLife* 2018, **7**:e30058.
- 1079 79. Potapova T, Gorbsky GJ: **The Consequences of Chromosome Segregation Errors in**
1080 **Mitosis and Meiosis.** *Biology* 2017, **6**:12.
- 1081 80. Mazoyer S: **Genomic rearrangements in the BRCA1 and BRCA2 genes.** *Human Mutation*
1082 2005, **25**:415-422.
- 1083 81. Peixoto A, Pinheiro M, Massena L, Santos C, Pinto P, Rocha P, Pinto C, Teixeira MR:
1084 **Genomic characterization of two large Alu-mediated rearrangements of the BRCA1**
1085 **gene.** *Journal of Human Genetics* 2013, **58**:78-83.
- 1086 82. Wang Y, Bernhardt AJ, Nacson J, Kraiss JJ, Tan Y-F, Nicolas E, Radke MR, Handorf E, Llop-
1087 Guevara A, Balmaña J, et al: **BRCA1 intronic Alu elements drive gene rearrangements**
1088 **and PARP inhibitor resistance.** *Nature Communications* 2019, **10**:5661.
- 1089 83. Hanahan D, Weinberg Robert A: **Hallmarks of Cancer: The Next Generation.** *Cell* 2011,
1090 **144**:646-674.
- 1091 84. Rodić N, Sharma R, Sharma R, Zampella J, Dai L, Taylor MS, Hruban RH, Iacobuzio-
1092 Donahue CA, Maitra A, Torbenson MS, et al: **Long Interspersed Element-1 Protein**
1093 **Expression Is a Hallmark of Many Human Cancers.** *The American Journal of Pathology*
1094 2014, **184**:1280-1286.
- 1095 85. Evering TH, Marston JL, Gan L, Nixon DF: **Transposable elements and Alzheimer's**
1096 **disease pathogenesis.** *Trends Neurosci* 2023, **46**:170-172.
- 1097 86. Buiting K, Nazlican H, Galetzka D, Wawrzik M, Groß S, Horsthemke B: **C15orf2 and a**
1098 **novel noncoding transcript from the Prader-Willi/Angelman syndrome region show**
1099 **monoallelic expression in fetal brain.** *Genomics* 2007, **89**:588-595.
- 1100 87. Angulo MA, Butler MG, Cataletto ME: **Prader-Willi syndrome: a review of clinical,**
1101 **genetic, and endocrine findings.** *J Endocrinol Invest* 2015, **38**:1249-1263.
- 1102 88. Mefford HC, Eichler EE: **Duplication hotspots, rare genomic disorders, and common**
1103 **disease.** *Current Opinion in Genetics & Development* 2009, **19**:196-204.
- 1104 89. Makoff AJ, Flomen RH: **Detailed analysis of 15q11-q14 sequence corrects errors and**
1105 **gaps in the public access sequence to fully reveal large segmental duplications at**
1106 **breakpoints for Prader-Willi, Angelman, and inv dup(15) syndromes.** *Genome Biology*
1107 2007, **8**:R114.

- 1108 90. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A,
1109 Clark AG, Donnelly P, Eichler EE, et al: **A global reference for human genetic variation.**
1110 *Nature* 2015, **526**:68-74.
- 1111 91. Lowy-Gallego E, Fairley S, Zheng-Bradley X, Ruffier M, Clarke L, Flicek P, null n: **Variant**
1112 **calling on the GRCh38 assembly with the data from phase three of the 1000 Genomes**
1113 **Project [version 2; peer review: 2 approved].** *Wellcome Open Research* 2019, **4**.
- 1114 92. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye
1115 K, Jun G, Hsi-Yang Fritz M, et al: **An integrated map of structural variation in 2,504**
1116 **human genomes.** *Nature* 2015, **526**:75-81.
- 1117 93. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter
1118 G, Marth GT, Sherry ST, et al: **The variant call format and VCFtools.** *Bioinformatics* 2011,
1119 **27**:2156-2158.
- 1120 94. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
1121 McCarthy SA, Davies RM, Li H: **Twelve years of SAMtools and BCFtools.** *GigaScience*
1122 2021, **10**.
- 1123 95. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar P,
1124 de Bakker PIW, Daly MJ, Sham PC: **PLINK: A Tool Set for Whole-Genome Association**
1125 **and Population-Based Linkage Analyses.** *The American Journal of Human Genetics*
1126 2007, **81**:559-575.
- 1127 96. Gao F, Chang D, Biddanda A, Ma L, Guo Y, Zhou Z, Keinan A: **XWAS: A Software Toolset**
1128 **for Genetic Data Analysis and Association Studies of the X Chromosome.** *J Hered* 2015,
1129 **106**:666-671.
- 1130 97. Keur N, Ricaño-Ponce I, Kumar V, Matzaraki V: **A systematic review of analytical**
1131 **methods used in genetic association analysis of the X-chromosome.** *Briefings in*
1132 *Bioinformatics* 2022, **23**.
- 1133 98. Jurka J: **Rebase Update: a database and an electronic journal of repetitive elements.**
1134 *Trends in Genetics* 2000, **16**:418-420.
- 1135 99. Nassar LR, Barber GP, Benet-Pagès A, Casper J, Clawson H, Diekhans M, Fischer C,
1136 Gonzalez JN, Hinrichs Angie S, Lee Brian T, et al: **The UCSC Genome Browser database:**
1137 **2023 update.** *Nucleic Acids Research* 2022, **51**:D1188-D1195.
- 1138 100. Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic**
1139 **features.** *Bioinformatics* 2010, **26**:841-842.
- 1140 101. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A,
1141 Pomeroy SL, Golub TR, Lander ES, Mesirov JP: **Gene set enrichment analysis: A**
1142 **knowledge-based approach for interpreting genome-wide expression profiles.**
1143 *Proceedings of the National Academy of Sciences* 2005, **102**:15545-15550.
- 1144 102. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP:
1145 **Molecular signatures database (MSigDB) 3.0.** *Bioinformatics* 2011, **27**:1739-1740.
- 1146 103. Chen S, Zhou Y, Chen Y, Gu J: **fastp: an ultra-fast all-in-one FASTQ preprocessor.**
1147 *Bioinformatics* 2018, **34**:i884-i890.
- 1148 104. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,
1149 Gingeras TR: **STAR: ultrafast universal RNA-seq aligner.** *Bioinformatics* 2012, **29**:15-21.

- 1150 105. Jin Y, Tam OH, Paniagua E, Hammell M: **TEtranscripts: a package for including**
1151 **transposable elements in differential expression analysis of RNA-seq datasets.**
1152 *Bioinformatics* 2015, **31**:3593-3599.
- 1153 106. Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, Mudge JM, Sisu
1154 C, Wright J, Armstrong J, et al: **GENCODE reference annotation for the human and**
1155 **mouse genomes.** *Nucleic Acids Research* 2018, **47**:D766-D773.
- 1156 107. Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for**
1157 **RNA-seq data with DESeq2.** *Genome Biology* 2014, **15**:550.
- 1158 108. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK: **limma powers**
1159 **differential expression analyses for RNA-sequencing and microarray studies.** *Nucleic*
1160 *Acids Research* 2015, **43**:e47-e47.
1161
1162

1163 **Legends to Figures**

1164

1165 **Fig 1. Overview of the pipeline to scan for genetic variants associated with L1/Alu**
1166 **global singletons.**

1167 **(A)** An illustration of the samples used in this study. SNV and SV genetic data was
1168 available for 2503 individuals from 5 super-populations, including 660 Africans (AFR),
1169 504 East Asians (EAS), 503 Europeans (EUR), 489 South Asians (SAS), and 347
1170 Admixed Americans (AMR). Males and females were approximately equally
1171 represented, with male-to-female ratios (M/F ratios) ranging from 0.91 to 1.14. **(B)** A
1172 schematic illustrating the trans-ethnic integration of available SNV and SV data to
1173 identify variants associated with L1/Alu insertion global singletons. Within each super-
1174 population, samples were segregated into cases and controls depending on whether or
1175 not they harbored a global Alu or L1 insertion singleton. GWAS was carried out within
1176 each super-population to identify polymorphic SNVs and SVs associated with case-
1177 control status. Finally, GWAS results from all 5 super-populations were meta-analyzed
1178 using a random effects statistical model, yielding a summary meta-analysis odds ratio
1179 and p-value for each variant. **(C)** The frequency of Alu and L1 insertion singletons in
1180 each super-population (*left panel*) or among cases within each super-population (*middle*
1181 *panel*). The distribution of insertion singletons across autosomes is also shown (*right*
1182 *panel*). **(D)** A Manhattan plot for the trans-ethnic GWAS meta-analysis. The dashed line
1183 at $p = 1.40E-5$ corresponds to an average empirical FDR < 0.05 , based on 20 random
1184 permutations. One such permutation is shown in the bottom panel for illustrative
1185 purposes. The solid line at $p = 6.00E-6$ corresponds to a Benjamini-Hochberg FDR $<$
1186 0.05 . The stricter of the two thresholds, $p = 6.00E-6$, was used to define significant
1187 SNVs and SVs. Significant variants overlapping regions in the ENCODE blacklist v2 are
1188 shown in blue and were omitted from downstream analyses. FDR: False Discovery
1189 Rate.

1190

1191 **Fig 2. Significant SNVs lie in genomic regions containing genes involved in**
1192 **transposon control.**

1193 **(A)** Scheme for assessing whether greenlisted SNVs were enriched in regions
1194 containing genes with TE regulatory potential. For a given gene set with regulatory
1195 potential (regulatory set A or B), the proportion of SNVs near genes in that gene set
1196 were calculated for the background and significant SNV lists, and statistical significance
1197 was assessed using Fisher's exact test. **(B)** Enrichment analysis of greenlisted SNVs
1198 near genes previously implicated in L1 expression control [38] or L1 transposition
1199 control [34] by CRISPR screening in cancer cell lines. Three specific examples of
1200 greenlisted SNVs near **(C)** genes controlling L1 expression and **(D)** genes controlling L1
1201 transposition are shown. **(E)** A summary of the associations we identified with various
1202 TE regulatory gene sets, highlighting the number of associated SNVs and the regulatory
1203 genes those SNVs were proximal to. Though not exclusive regulators of TE activity *per*
1204 *se*, we included in our analysis gene sets involved in "histone methyltransferase activity"
1205 and "RNA modification" functions, since those processes have been implicated in
1206 transposon control. FDR: False Discovery Rate.

1207

1208 **Fig 3. Significant SNVs are enriched in genomic regions containing features**
1209 **associated with genome instability.**

1210 **(A)** Scheme for assessing whether greenlisted SNVs are enriched in regions containing
1211 elements known for promoting genome instability. For a given set of potentially
1212 genetically unstable regions (unstable element set A or B), the proportion of SNVs
1213 overlapping regions in each set are calculated for the background and significant SNV
1214 lists, and statistical significance is assessed using Fisher's exact test. **(B)** Enrichment
1215 analysis of greenlisted SNVs overlapping evolutionary age-stratified Alu (*left column*)
1216 and L1 (*right column*) copies. **(C)** Enrichment analysis of greenlisted SNVs overlapping
1217 curated L1 loci in L1Base v2 [65]. This database contains putatively active L1 copies
1218 (with either full-length, fully intact L1 copies or L1 copies with only ORF2 intact), as well
1219 as non-autonomous, full-length, non-intact L1 copies with regulatory potential. **(D)**
1220 Enrichment analysis of greenlisted SNVs overlapping genomic regions containing
1221 segmental duplications [71, 72] or structural variation hotspots [70]. FDR: False
1222 Discovery Rate.

1223

1224 **Fig 4. Alterations in the cell cycle are positively correlated with case status.**

1225 **(A)** Scheme for characterizing transcriptomic differences between case and control
1226 samples. Gene expression profiles were quantified using mRNA-sequencing data from
1227 lymphoblastoid cells belonging to 358 European and 86 African individuals. To note, all
1228 African individuals here were from the Yoruban population. These gene expression
1229 profiles were used to construct consensus gene co-expression networks with WGCNA.
1230 We then quantified the correlations between each module in the network and the case-
1231 control status of all samples (encoded as 0 for controls and 1 for cases). Finally, over-
1232 representation analysis (ORA) using the Gene Ontology (GO) Biological Process gene
1233 set collection was used to assign functions to significantly correlated modules. **(B)** The
1234 correlations and correlation p-values between consensus network modules and case-
1235 control status in the European and African cohorts. Boxes were color-coded according
1236 to the strength of the correlation. A meta-analysis was also carried out to summarize
1237 statistical results by combining European and African correlation p-values using Fisher's
1238 method. For visualization purposes only, the average of the European and African
1239 correlations was assigned to the meta-analysis. Correlations with opposite trends in the
1240 two cohorts were disregarded in the meta-analysis and colored grey. Correlations with p
1241 < 0.01 were considered statistically significant and highlighted in bold. **(C)** The top 10
1242 ORA results for the MEroyalblue module using the GO Biological Process gene set
1243 collection. The colors represent the gene ratio (i.e. the fraction of module genes from
1244 the listed gene set) and the sizes of the dots represent the Benjamini-Hochberg FDR.
1245 NA: Not Applicable. FDR: False Discovery Rate.

1246

1247

1248 **Legends to Supplementary Figures**

1249

1250 **S1 Fig. Quality control of combined SNV and SV 1000 Genomes Project data.**

1251 PCA plots for pruned SNV and SV genotype data from **(A)** African, **(B)** East Asian, **(C)**
1252 European, **(D)** South Asian, and **(E)** Admixed American samples. Colors and shapes
1253 represent different populations within each super-population.

1254

1255 **S2 Fig. GWAS in individual super-populations is underpowered.**

1256 Manhattan plots for GWAS results in individual super-populations, including for the **(A)**
1257 African, **(B)** East Asian, **(C)** European, **(D)** South Asian, and **(E)** Admixed American
1258 cohorts. Solid lines correspond to a Benjamini-Hochberg FDR < 0.05 and dashed lines
1259 correspond to an average empirical FDR < 0.05, based on 20 random permutations.
1260 The Benjamini-Hochberg FDR and average empirical FDR, respectively, corresponded
1261 to the following p-values in each super-population: $p = 1.18E-6$ and $p = 4.61E-6$ in the
1262 African cohort, $p = 9.06E-7$ and $p = 8.46E-7$ in the South Asian cohort, and $p = 3.53E-7$
1263 and $p = 1.07E-6$ in the Admixed American cohort. The stricter of the two thresholds in
1264 each super-population was used to define significant SNVs and SVs. No variant at an
1265 FDR < 0.05 was identified in the East Asian and European cohorts. Significant variants
1266 overlapping regions in the ENCODE blacklist v2 are shown in blue. FDR: False
1267 Discovery Rate.

1268

1269 **S3 Fig. Functional annotation of significant variants.**

1270 **(A)** Scheme for predicting functions of genomic regions containing significant variants.
1271 All SNVs and SVs used in this study were assigned genes using the GREAT [47] online
1272 platform. Significant SNV- and SV-associated genes were then tested for functional
1273 gene set enrichment by over-representation analysis (ORA) using clusterProfiler [49],
1274 specifying the corresponding background-associated genes as the universe. **(B)** The
1275 number of genes associated to greenlisted SNVs (*left*), the distance between
1276 greenlisted SNVs and the transcriptional start sites (TSS) of associated genes (*middle*),
1277 and the top 10 ORA results for associated genes using the GO Biological Process gene
1278 set collection (*right*). The colors represent the gene ratio (i.e. the fraction of significant

1279 SNV-associated genes from the listed gene set) and the sizes of the dots represent the
1280 Benjamini-Hochberg FDR. **(C)** The number of genes associated to greenlisted SVs (*left*)
1281 and the distance between greenlisted SVs and the transcriptional start sites (TSS) of
1282 associated genes (*middle*). **(D)** Enrichment analysis of greenlisted SNVs overlapping
1283 genomic regions with candidate cis-Regulatory Elements (cCREs) from the ENCODE
1284 Registry v4 [53]. **(E)** Heatmap comparing the median expression levels of significant
1285 SNV-associated genes in each tissue included in the GTEx Analysis v8. FDR: False
1286 Discovery Rate.

1287

1288 **S4 Fig. Polymorphic SVs of different classes are associated with L1/Alu insertion**
1289 **singletons.**

1290 **(A)** One example of each type of significant, polymorphic SV that was associated with
1291 L1/Alu singletons. These classes included inversions, Alu insertions, an L1 insertion,
1292 SINE-VNTR-Alu (SVA) insertions, and a multiallelic copy number variant. FDR: False
1293 Discovery Rate.

1294

1295

1296 **Inventory of Supplementary Tables**

1297

1298 **S1 Table. Sample singleton counts, significant variant annotations, and gene co-**
1299 **expression network results.**

1300 **(A)** Number of singletons for each sample. **(B)** All variants passing $FDR < 0.05$ in the
1301 GWAS meta-analysis. **(C)** Over-representation analysis of genes associated to
1302 greenlisted, significant SNVs using GO Biological Process gene sets. **(D)** Over-
1303 representation analysis of genes associated to greenlisted, significant SVs using GO
1304 Biological Process gene sets. **(E)** SnpEff annotations of significant variants. **(F)**
1305 Expression levels (median tissue-specific TPMs) of significant SNV-associated genes in
1306 the GTEx Analysis v8. The cluster of brain-associated genes is in blue, and the cluster
1307 of testes-associated genes is in orange. **(G)** TE enrichment analysis with TEENA. **(H)**
1308 Lymphoblastoid cell line WGCNA network gene-module assignments. **(I)** Over-
1309 representation analysis of MEroyalblue genes using GO Biological Process gene sets.

1310

1311

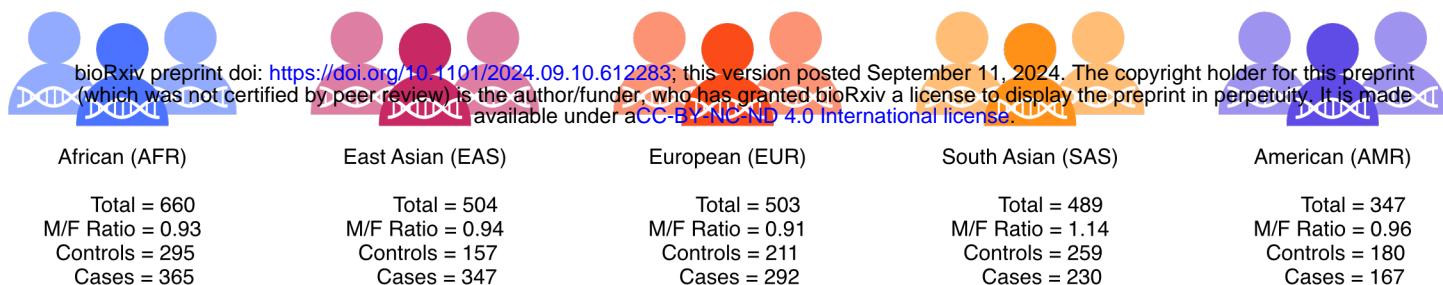
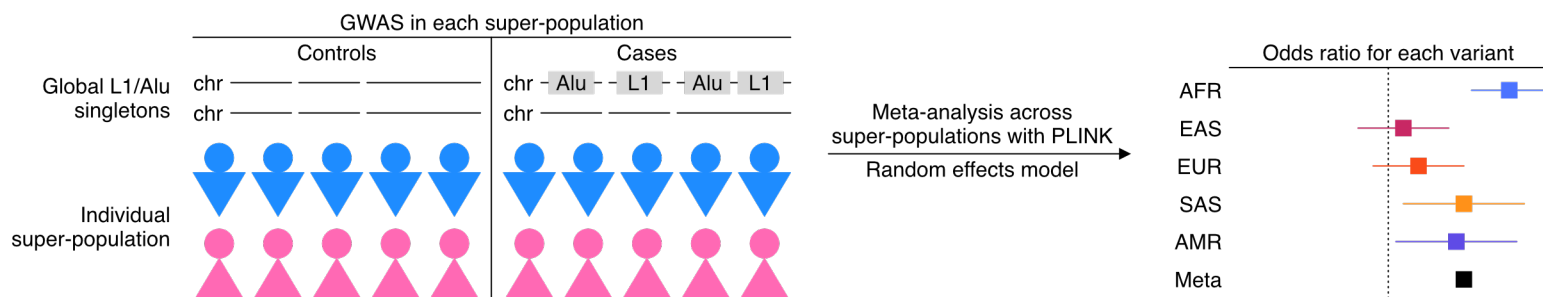
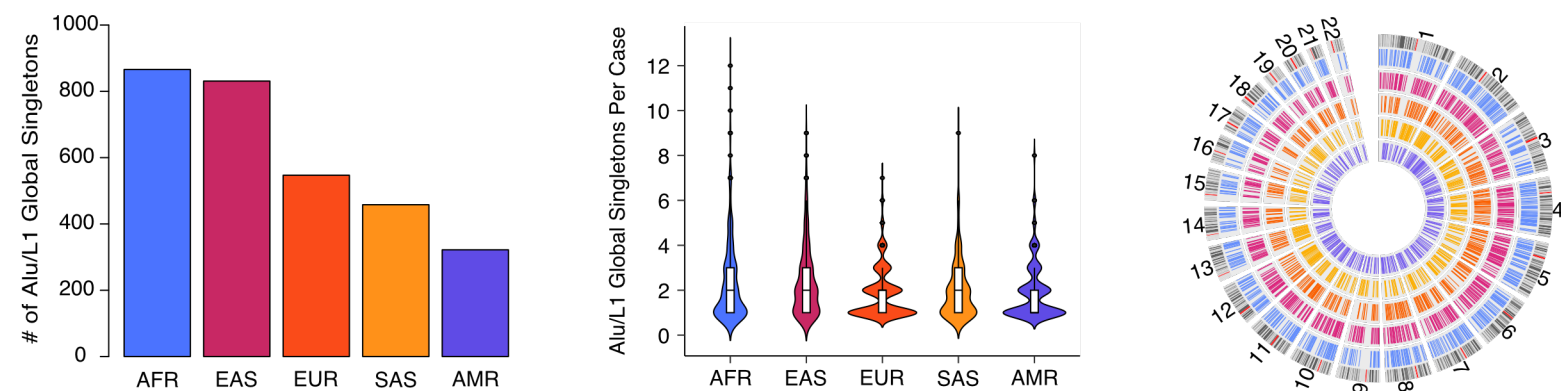
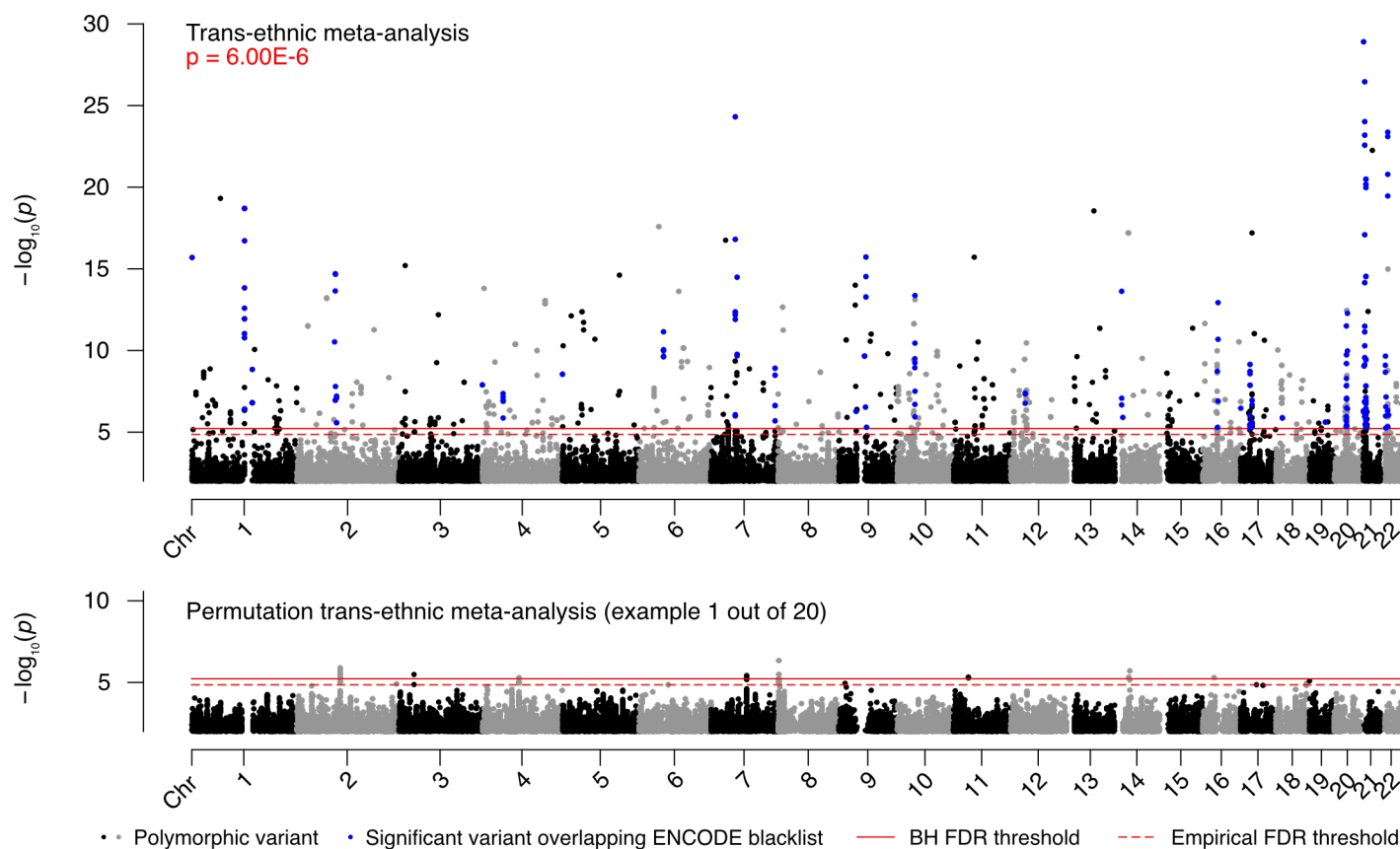
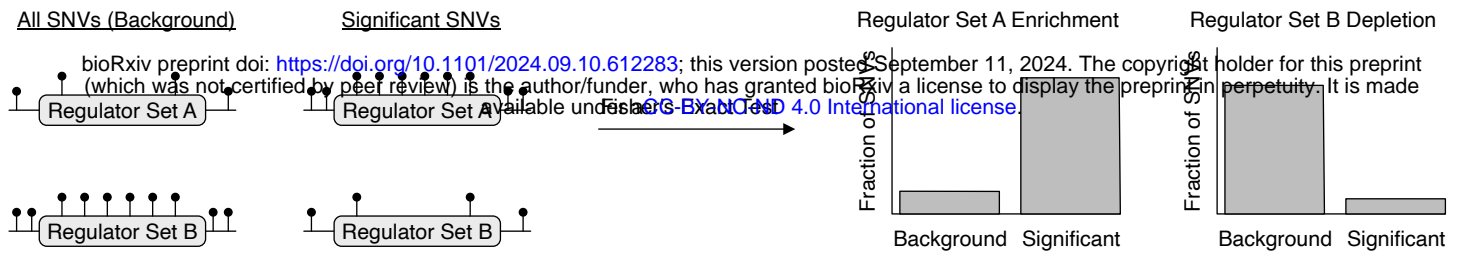
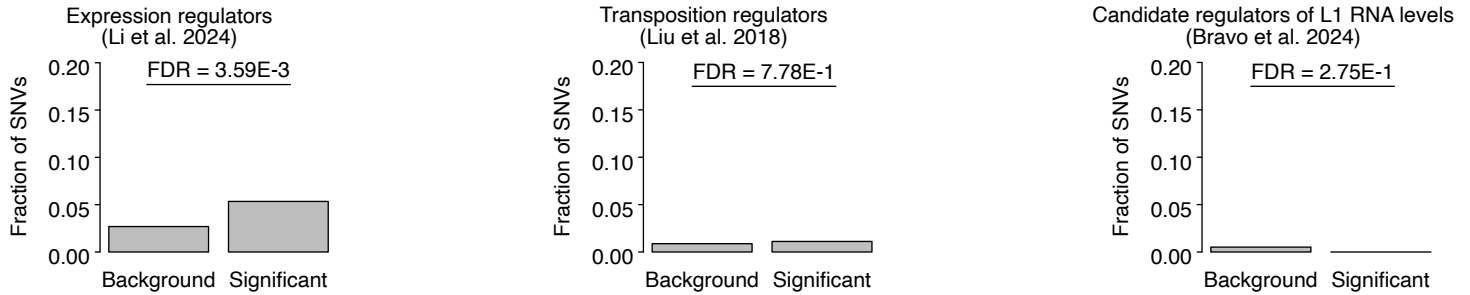
Figure 1**A** The 2503 samples from the 1000Genomes Project used in this study**B** Case-control study design for trans-ethnic GWAS**C** Frequency and genomic distribution of L1/Alu global singletons across superpopulations**D** Manhattan plot for the trans-ethnic GWAS meta-analysis associations

Figure 2

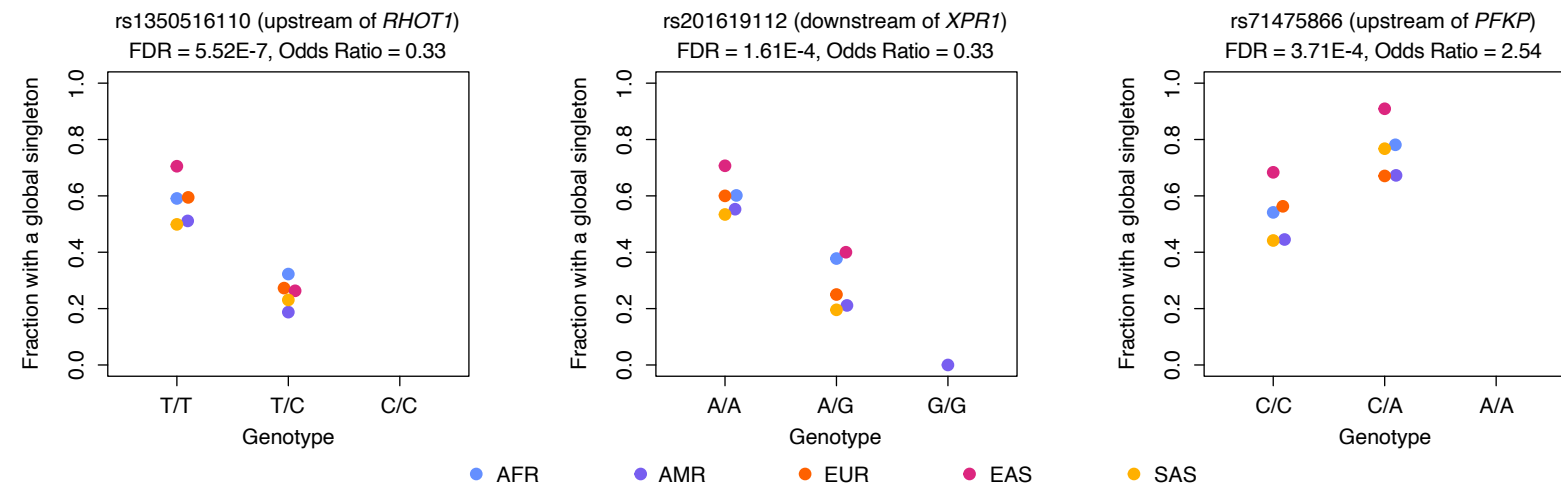
A Scheme for linking significant SNVs with potential transposon regulators



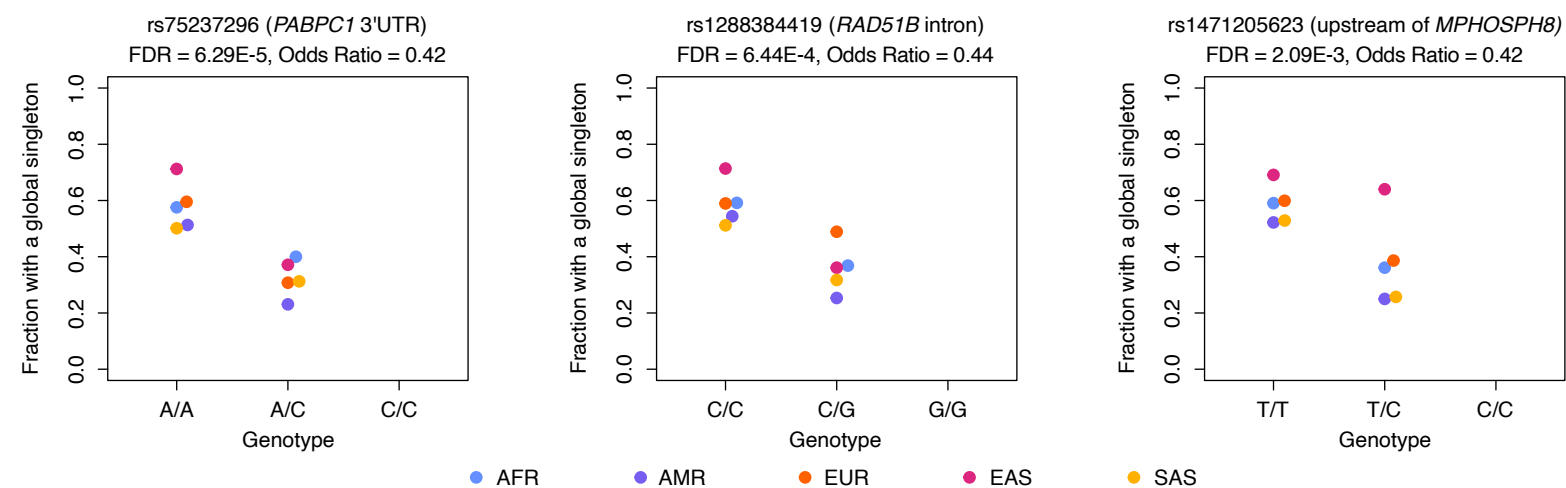
B SNVs near known L1 regulators



C Example SNVs near L1 expression regulators



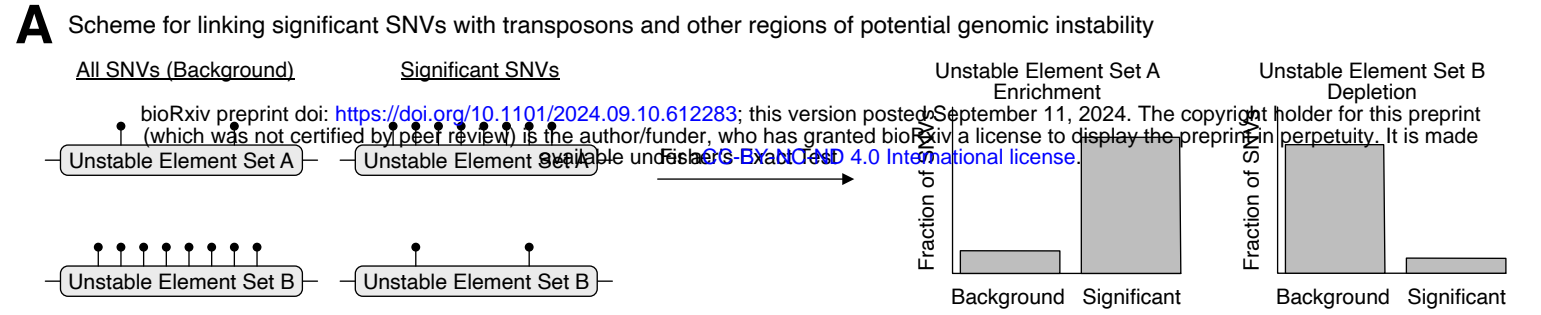
D Example SNVs near L1 transposition regulators



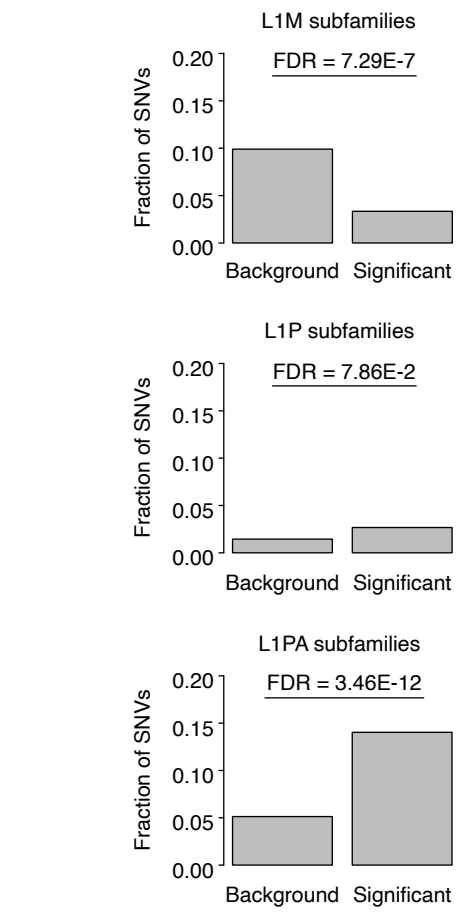
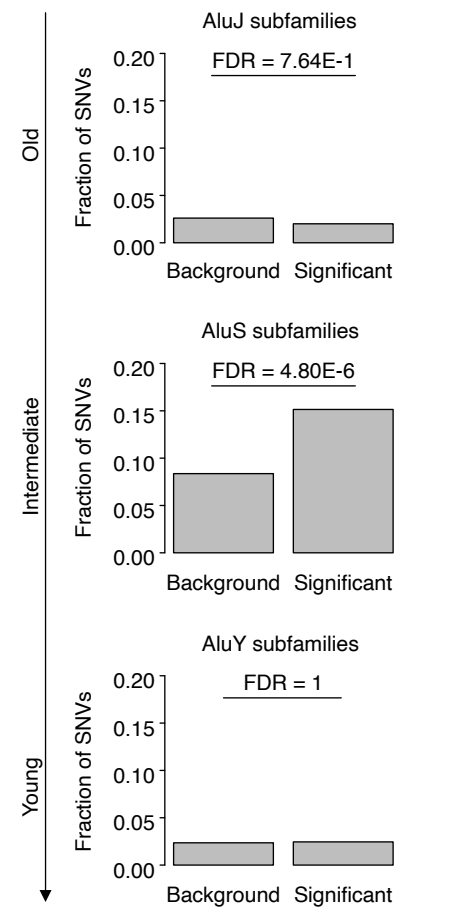
E Summary of SNVs near known TE regulators and genes with functions involved in TE regulation

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>

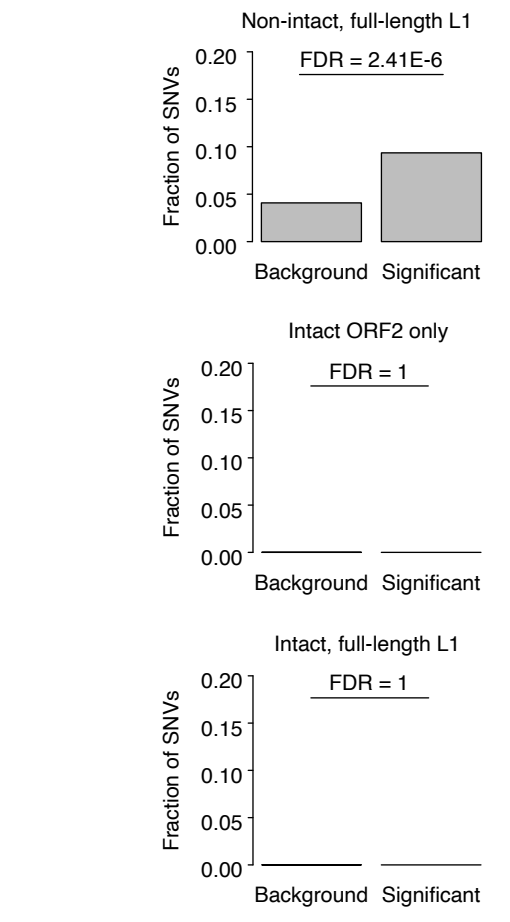
Figure 3



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

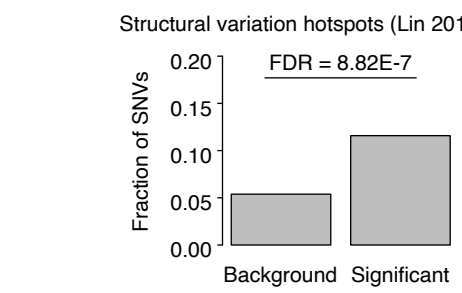
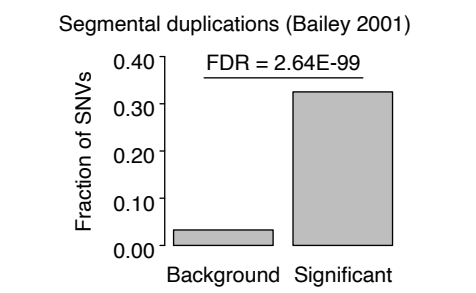
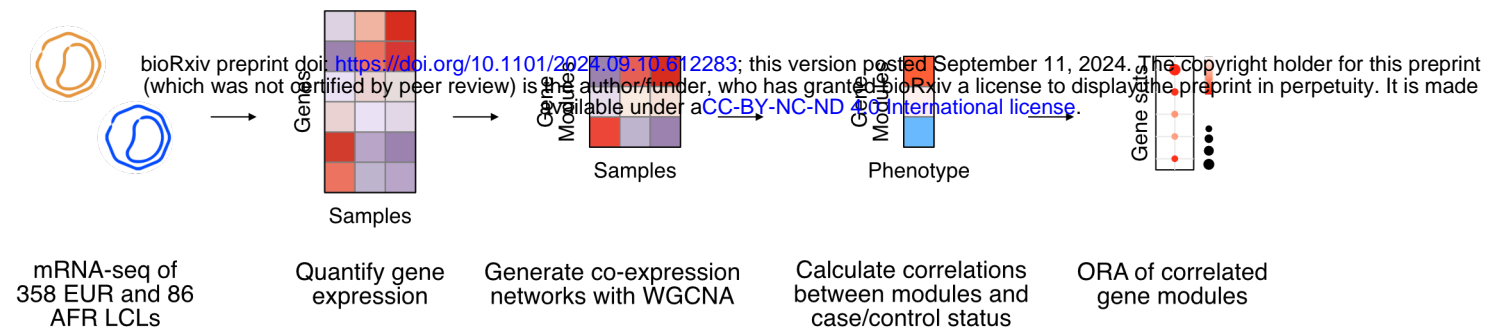
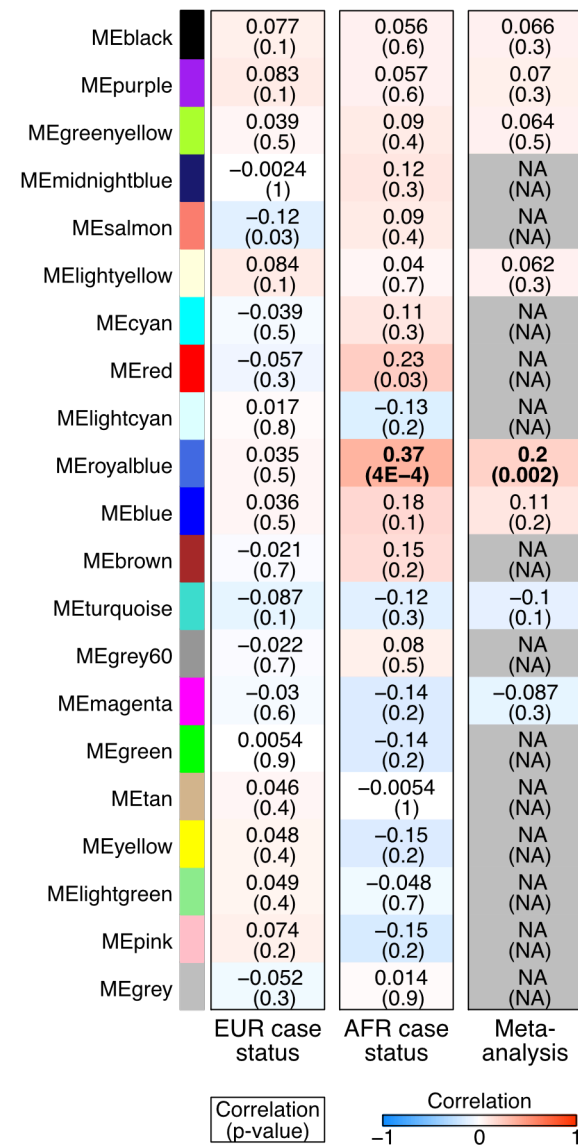


Figure 4

A Scheme for assessing transcriptomic differences between cases and controls



B Network module correlations with case/control status



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

