

1 A computational framework to assess genome-wide distribution of polymorphic human
2 endogenous retrovirus-K in human populations

3

4 Weiling Li¹, Lin Lin², Raunaq Malhotra^{1,4}, Lei Yang³, Raj Acharya^{1,5}, and Mary Poss^{3*}

5 ¹ The School of Electrical Engineering and Computer Science,

6 ² Department of Statistics, ³ Department of Biology and Veterinary and Biomedical Sciences

7 The Pennsylvania State University, University Park, PA, 16802, USA

8 ⁴ Commense Health, Cambridge, MA 02142, USA

9 ⁵ School of Informatics, Computing and Engineering, Indiana University, Bloomington, IN 47405,

10 USA

11

12 Running Title: Global genomic distribution of polymorphic HERV-K

13 * Corresponding author:

14 Mary Poss

15 814 404-7061

16 maryposs@gmail.com

17

18 **Keywords** HERV-K, human endogenous retrovirus, insertional polymorphism, 1000 Genomes
19 Project, mixture model, visualization

20

21 **Abstract**

22 Human Endogenous Retrovirus type K (HERV-K) is the only HERV known to be
23 insertionally polymorphic. It is possible that HERV-Ks contribute to human disease because
24 people differ in both number and genomic location of these retroviruses. Indeed viral
25 transcripts, proteins, and antibody against HERV-K are detected in cancers, auto-immune, and
26 neurodegenerative diseases. However, attempts to link a polymorphic HERV-K with any disease
27 have been frustrated in part because population frequency of HERV-K provirus at each site is
28 lacking and it is challenging to identify closely related elements such as HERV-K from short read
29 sequence data. We present an integrated and computationally robust approach that uses
30 whole genome short read data to determine the occupation status at all sites reported to
31 contain a HERV-K provirus. Our method estimates the proportion of fixed length genomic
32 sequence (*k-mers*) from whole genome sequence data matching a reference set of *k-mers*
33 unique to each HERV-K loci and applies mixture model-based clustering to account for low
34 depth sequence data. Our analysis of 1000 Genomes Project Data (KGP) reveals numerous
35 differences among the five KGP super-populations in the frequency of individual and co-
36 occurring HERV-K proviruses; we provide a visualization tool to easily depict the prevalence of
37 any combination of HERV-K among KGP populations. Further, the genome burden of
38 polymorphic HERV-K is variable in humans, with East Asian (EAS) individuals having the fewest
39 integration sites. Our study identifies population-specific sequence variation for several HERV-K
40 proviruses. We expect these resources will advance research on HERV-K contributions to
41 human diseases.

42

43 **Author summary**

44 Human Endogenous Retrovirus type K (HERV-K) is the youngest of retrovirus families in the
45 human genome and is the only group that is polymorphic; a HERV-K can be present in one
46 individual but absent from others. HERV-Ks could contribute to disease risk but establishing a
47 link of a polymorphic HERV-K to a specific disease has been difficult. We develop an easy to use
48 method that reveals the considerable variation existing among global populations in the
49 frequency of individual and co-occurring polymorphic HERV-K, and in the total number of HERV-
50 K that any individual has in their genome. Our study provides a global reference set of HERV-K
51 genomic diversity and tools needed to determine the genomic landscape of HERV-K in any
52 patient population.

53

54 **Introduction**

55 Endogenous retroviruses (ERVs) are derived from infectious retroviruses that integrated
56 into a host germ cell at some time in the evolutionary history of a species [1–5]. ERVs in
57 humans (HERVs) comprise up to 8% of the genome and have contributed important functions
58 to their host [6–8]. The infection events that resulted in the contemporary profile of HERVs
59 occurred prior to emergence of modern humans so most HERVs are fixed in human populations
60 and those of closely related primates. However some HERVs retain the ability to replicate and
61 reintegrate into germline so that individuals differ in the number and genomic location
62 occupied by an ERV, a situation termed insertional polymorphism [9–11]. Among all families of
63 HERVs, HERV-K is the only one known to be insertional polymorphic in humans.

64

65 A full-length retroviral sequence is called a provirus and encodes several viral structural
66 or regulatory proteins that are flanked by two long terminal repeats (5' or 3' LTR). While there
67 are several HERV-K that are full length, none are infectious and most contain mutations or
68 deletions that affect the open reading frames or truncate the virus. Further, the identical LTRs
69 are substrates for homologous recombination, which deletes virus genes while retaining a
70 single, or solo, LTR at the integration site [12–14]. Thus, in a population, a site could be
71 unoccupied, occupied by a HERV-K provirus, or contain a solo LTR. Insertional polymorphism
72 typically refers to the occupancy at a loci [15,16]. However the occupied site can contain a
73 provirus or solo LTR and a provirus sequence can vary among individuals. Thus HERV-K and
74 other HERVs can contribute to genomic diversity in the global human population in several
75 ways [17].

76
77 HERVs from multiple families have been linked with both proliferative and degenerative
78 diseases in humans [18–24] . Although there are known mechanisms by which a HERV can
79 cause disease; for example, by inducing genome structural variation through recombination
80 [25–29], affecting host gene expression [30], and inappropriate activation of an immune
81 response by viral RNA or proteins [21], it has been difficult to establish an etiological role of a
82 HERV in any disease. HERV-K specifically has been associated with breast and other cancers
83 [3,31–35], and autoimmune diseases, such as rheumatoid arthritis [36,37], multiple sclerosis
84 [20,38] and systemic lupus erythematosus [8,20,39] without definitive evidence of causality or
85 of the specific loci involved. Recently, a HERV-K envelope protein was shown to recapitulate the

86 clinical and histological lesions characterizing Amyotrophic Lateral Sclerosis [40,41], providing an
87 important mechanistic advance of a role for a HERV-K protein in a disease.

88 In this paper, we focus on characterizing the genome landscape of known insertionally
89 polymorphic HERV-K proviruses in the 1000 Genomes Project (KGP) data. We present a data-
90 mining tool and a statistical framework that accommodates low depth data characteristic of the
91 KGP - and often patient - data to estimate the presence or absence of a provirus at known
92 HERV-K loci. Because combinations of HERV-K may act synergistically in the pathogenesis of a
93 disease [42], we estimate the co-occurrence of polymorphic HERV-K proviruses in different
94 populations and provide a tool to visualize HERV-K co-occurrence in global populations. Our
95 results provide a reference of global population diversity in HERV-K proviruses at all currently
96 known loci in the human genome and demonstrate that there are notable differences among
97 population frequencies of HERV-Ks and the total number of HERV-Ks found in a person's
98 genome.

99

100 **Results**

101

102 **A model to estimate polymorphic HERV-K from whole genome sequence data.**

103 The goal of this research was to develop a computationally efficient and easy to use tool
104 that could accurately report the status of all HERV-Ks with coding potential (provirus) from
105 whole genome sequence (WGS) data. We use the KGP database to establish the global
106 population diversity of each polymorphic HERV-K and the burden of HERV-K in individual

107 genomes to provide a foundation to study the role of HERV-K in human disease. Our method
108 takes as input all reads that map to identified HERV-K elements in hg19. The rationale here is
109 that polymorphic HERV-K are very similar to those in the reference genome and will map on
110 existing elements. The recovered reads are reduced to *k-mers* and mapped to a reference set
111 of *k-mers* representing all unique sites in every HERV-K in the database. The output is a ratio of
112 subject *k-mers* (*n*) that are 100% match to the reference *k-mers* (*T*) (see methods for full
113 details).

114 Our preliminary analysis of the KGP data demonstrated that our *k-mer*-based approach
115 is sensitive to sequence depth; some HERV-K are represented by an almost continuous range of
116 *n/T* from 0-1 (Fig 1A), making presence/absence classification difficult. A comparison of a
117 subset of the 28 individuals in the KGP data that have both low and high sequence depth data
118 shows how depth affects *n/T* (Fig 1B, see S2 Fig for data of all 28 persons). If read depth is
119 greater than 20, there is less dispersion of *n/T* values, most likely because more reads are
120 recovered from the mapped intervals. However, the majority of the KGP data is approximately
121 6x depth and thus to make use of this important resource, we developed a mixture model to
122 cluster the *n/T* values from genomes sequenced at low depth. *K* was optimized to 50 because
123 this value improved our model computational efficiency and output (Fig 1B, S1-S3 Methods, S1
124 Fig). The states, 'provirus', 'solo LTR', and 'absent' are preliminarily assigned to each cluster
125 based on the high depth data (Fig 1B). Individuals with *n/T*=1 have the reference allele and
126 *n/T*=0 indicates that the HERV-K is absent (no *k-mers* to unique sites in the HERV-K were
127 recovered from mapped sequence reads). The *k-mers* derived from persons with low and

128 intermediate n/T values were mapped to each HERV-K to determine whether they localized
129 only in the LTR (assign 'solo LTR') or in the coding region (assign 'provirus') (S3 Fig).

130

131 **Fig 1. A mixture model to account for low depth WGS data**

132 A) The plot displays the n/T value for 2535 individuals from KGP with low depth sequence data for
133 chr12:55727215-55728183 when K=70. There is little resolution of values to enable assignment to
134 provirus, solo LTR, or absent states.

135 B) The result of the mixture model on the same data. The individual clusters generated by the model are
136 indicated by a unique color; in the example shown, there are four clusters. K has been optimized to 50
137 to enable clear clustering in the 'absent', 'solo LTR' and 'provirus' states. In this example, eight of the 28
138 individuals that have both low and high depth sequence data (see S1_Dataset:KGP) are shown to
139 demonstrate the effect of sequence depth. The n/T ratio is 1 for persons with high depth data [red
140 numbers, #6 and 12] who have the reference allele, while the corresponding low depth data [black
141 numbers, yellow cluster] from the same individuals have n/T ranging from 0.7 to 0.9. There is less of an
142 effect of depth for individuals who do not have the HERV-K (n/T=0). However, optimizing K facilitates
143 separation of clusters for absent [red cluster, #23 and 28], and solo LTR [green cluster, #4 and 16]. States
144 are confirmed by mapping the *k-mers* from individuals in a cluster to the reference HERV-K (S3 Fig).

145

146

147 **Prevalence of polymorphic HERV-K in each KGP super-population**

148 The WGS data of each individual in the KGP dataset were evaluated using our analysis
149 workflow. HERV-Ks on chrY were not considered. Twenty sites, omitting one at chr1:73594980
150 [see methods] were identified that were polymorphic for containing a HERV-K provirus. A

151 phylogenetic analysis of all HERV-Ks greater than 6 kbp shows that polymorphic HERV-Ks are
152 closely related (S4 Fig). The prevalence of the 20 polymorphic HERV-Ks varied from 0.9% to
153 99.5% when averaged across the entire KGP dataset (Table 1). However, there were notable
154 differences in prevalence at each site among the five super-populations (AFR, EAS, AMR, EUR,
155 SAS). Of the 20, the prevalence of seven polymorphic HERV-Ks was greater than 90% and the
156 difference between populations with the lowest and highest prevalence was less than 6.5%
157 (Table 1). There was 100% occupancy for six of the seven high prevalence polymorphic HERV-Ks
158 (98.8% for the seventh), indicating that the rate of conversion to solo LTR is low for viruses at
159 these sites (S1 Table). Two polymorphic HERV-Ks had an overall prevalence of less than 10% in
160 any population (Table 1) and we found no evidence of a solo LTR at these sites; both are found
161 in individuals from AFR. Nine of the remaining 11 HERV-Ks are of interest because the
162 difference between super-populations with the highest and lowest prevalence is between 28
163 and 80 percentage points (Table 1). Of note, the prevalence is lowest in EAS populations for the
164 three HERV-Ks with the largest difference among super-populations.

165
166 Individuals from African populations differ significantly from the other four super-
167 populations in the prevalence of ten of the polymorphic HERV-K, three of which occur in close
168 proximity on chr19. (Table 1, S2_Dataset:compare_prevalence). EUR and AFR super-populations
169 are significantly different at all but one of the 20 polymorphic HERV-K based on adjusted p-
170 values (S2_Dataset:compare_prevalence).

171

172 **Table 1. Provirus frequencies of polymorphic HERV-K.**

173

	KGP	AFR	AMR	EAS	EUR	SAS	max- min
	pro- virus (%)						
<u>chr1:75842771</u>^c	42.88	26.76	56.53	6.02	68.91	66.80	62.89
chr3:112743479^a	98.46	96.71	99.72	99.81	99.60	97.37	3.09
chr3:148281477	41.89	38.86	42.61	45.05	46.53	37.45	9.09
<u>chr3:185280336</u>^a	99.49	98.06	100.00	100.00	100.00	100.00	1.94
<u>chr4:69463709</u>^c	72.50	93.87	88.92	31.07	85.35	61.94	62.80
chr5:156084717^a	99.41	98.36	99.72	100.00	99.80	99.60	1.64
chr6:57623896^a	93.65	90.73	97.16	90.87	97.23	94.33	6.50
chr6:78427019^a	97.71	95.52	97.16	99.61	97.23	99.60	4.10
chr7:4622057*^c	47.50	61.14	30.11	58.25	36.44	41.50	31.02
chr8:12316492^c	14.08	32.88	12.22	0	15.64	3.04	32.88
<u>chr8:7355397</u>^c	18.66	39.16	12.50	6.02	11.29	15.99	33.14
<u>chr10:27182399</u>^a	99.13	97.46	99.43	99.81	99.80	99.80	2.35
chr11:101565794^c	63.04	80.87	77.27	6.99	86.53	63.16	79.54

<u>chr12:55727215</u>	72.19	72.80	80.40	63.30	80.99	65.79	17.69
chr12:58721242^c	70.73	58.89	78.41	60.00	87.33	75.51	28.43
<u>chr19:21841536^c</u>	26.98	39.16	11.93	32.23	10.69	32.39	28.47
<u>chr19:22414379^c</u>	67.77	89.24	60.80	56.89	55.84	67.21	33.40
<u>chr19:22457244^b</u>	0.87	3.29	0.00	0.00	0.00	0.00	3.29
chr22:18926187^a	99.49	98.36	99.72	100.00	99.80	100.00	1.64
<u>chrX:93606603^b</u>	2.25	7.32	2.27	0.00	0.00	0.00	7.32

174 For simplicity, only the starting coordinate is listed.

175 * The value given represents those with the tandem repeat

176 ^a: prevalence > 90%

177 ^b: low prevalence and no solo LTR

178 ^c: max-min difference is > 28%

179 underline: AFR significantly different from other 4 super populations.

180 See S2_Dataset:compare_prevalence for full data set.

181

182 **The number of polymorphic HERV-Ks per individual**

183

184 The HERV-K genome is close to 10 kbp. As there are 20 HERV-Ks that are polymorphic in
 185 human populations, we asked if some individuals carry a different burden of these repetitive,
 186 and potentially functional, viral elements than others. This was indeed the case. The number of

187 polymorphic HERV-K proviruses per person ranges from 7-18 (Fig 2, S2_Dataset:HERV-K per
188 person). More than 63% of individuals from all super-populations except EAS carry 12 to 14
189 proviruses in their genome. Individuals from EAS have a lower burden with 69% of individuals
190 carrying 9-11 HERV-K proviruses. 7% of AFR individuals have 16 or 17 proviruses compared to a
191 maximum of 2% in other groups (S2_Dataset:HERV-K per person). These data highlight the
192 importance of using a comprehensive approach to study the potential disease impact of
193 polymorphic HERV-Ks, because total HERV-K burden in individuals might influence a disease
194 phenotype without significant variation in prevalence of each provirus in a patient pool.

195

196 **Fig 2. Histogram of the number of proviruses per individual among super-populations.**

197

198

199 **Co-occurrence of polymorphic HERV-Ks**

200 Our data provide a comprehensive picture of sites occupied by HERV-K provirus in each
201 genome, which enables analysis of polymorphic HERV-K co-occurrence in populations. We
202 assessed combinations of three, four and five polymorphic HERV-Ks and found that there are
203 many combinations of co-occurring viruses that are population-specific (S3_Dataset). To
204 facilitate exploration of HERV-K combinations among KGP populations, we developed a D3.js
205 visualization tool that allows a user to choose any combination of the 20 polymorphic HERV-K
206 proviruses and display the co-occurrence prevalence among the 26 populations represented in
207 the KGP data. As an example, we present a combination of four HERV-Ks to represent the
208 variation that occurs in KGP individuals, which in this case ranges from 3% in EAS to 59% in EUR

209 (Fig 3A). We also determine that the three polymorphic HERV-Ks found on chr19 co-occur only
210 from three AFR populations and in less than 2% of individuals (Fig 3B).

211

212 **Fig 3. A visualization tool to examine co-occurrence of polymorphic HERV-Ks.**

213 A) The co-occurrence of HERV-Ks at chr1:75842771-75849143, chr3:112743479-112752282,
214 chr6:57623896-57628704, and chr12:58721242-58730698 in the 26 populations are represented based
215 on their geographic location. The relative frequency for these four co-occurring HERV-Ks in each
216 population bubble is displayed based on the color gradient shown in the scale at the top. The actual
217 prevalence of the given combination of HERV-K provirus for each population and the cumulative
218 prevalence for the super-population are shown in text on the right. Note that AFR and EAS have the
219 lowest prevalence of these four polymorphic HERV-Ks.

220 B) As in (A) showing the co-occurrence of the three polymorphic HERV-Ks that are present on chr19 by
221 population. This is a rare combination only found in two AFR populations and individuals in the
222 Caribbean of African ancestry.

223

224

225 **HERV-K status informs KGP super-populations**

226 Because there are clearly population-specific differences in both individual HERV-K
227 frequency and in the frequency of HERV-K co-occurrence, we explored whether the presence or
228 absence of polymorphic HERV-Ks is sufficient to distinguish populations using Fisher's linear
229 discriminant analysis (LDA) [43]. Based on the status 'provirus', 'solo LTR', or 'absence', there is
230 little resolution of AFR, EUR, and EAS super-populations (Fig 4A). However, there is sufficient
231 signature to separate AFR, EUR, and EAS if we utilize the n/T ratio on the 20 polymorphic HERV-

232 Ks (S5 Fig) and we further improve population separation if we use the n/T ratio for all 96 HERV-
233 Ks (Fig 4B). This indicates that we are losing information by reducing the data to three states
234 and that fixed HERV-K also contain signal for population of origin.

235

236 **Fig 4. Linear discriminant analysis of HERV-K status among three super-populations.**

237 A) LDA based on the states 'provirus', 'solo LTR' and 'absence' of the 20 polymorphic HERV-K for AFR,
238 EAS, and EUR. AMR and SAS overlap these three populations and are removed for clarity B) LDA plot on
239 n/T ratio of all 96 HERV-K discriminates AFR, EAS, and EUR super-populations. See S6 Fig for plots with
240 all five super-populations.

241

242 An $n/T = 1$ indicates that we recovered all reads that map to the reference set T for a
243 specific HERV-K. If there is a HERV-K allele that has not been reported in any database but that
244 is common in a population, we expect $n/T < 1$ because we require 100% match to reference set
245 T and *k-mers* covering allelic sites will not be included. We assessed the density distributions of
246 n/T plots for each of the 96 HERV-Ks for evidence of population-specific alleles (S4 Method, S7
247 Fig). Five HERV-Ks have some indication of population specific distributions (S1_Dataset:virus).
248 The HERV-K at chr1:155596457-155605636, which we report as fixed, is notable because the
249 reference allele ($n/T=1$) is only found in AFR (Fig 5A, S7 Fig). Individuals from most other
250 populations have n/T near 0.5. We mapped *k-mers* from individuals with n/T near 0.5 to the
251 reference HERV-K sequence and confirmed that there is a loss of *k-mers* at several sites covered
252 by the unique reference *k-mers* for this virus (S8 Fig). There are also cases where the reference
253 allele is found in all populations except AFR (Fig 5B and see S7 Fig for additional examples).

254

255 **Fig 5. Population specificity of HERV-K alleles.** A) n/T plot for chr1:155596457-155605636 colored by
256 each of the 5 super-populations. Only individuals from AFR and a few from AMR have the HERV-K
257 reference sequence. B) Plot of chr5:156084717-156093896 colored by each of the 5 super-populations.
258 In this case, all populations except AFR have the reference allele and all super-populations have an
259 alternative allele that is not present in the databases.

260

261 **Discussion**

262 Our research provides a tool to mine whole genome sequence data to collectively
263 evaluate the status of HERV-K provirus at known polymorphic and fixed sites in the human
264 genome. The tool incorporates a statistical clustering algorithm to accommodate low depth
265 sequence data and a visualization tool to explore the co-occurrence of polymorphic HERV-K in
266 the global populations represented in the KGP data. There are numerous significant differences
267 in the prevalence of individual and co-occurring polymorphic HERV-K among the five KGP super-
268 populations. It is notable that individuals from EAS carry a lower total burden of HERV-K than
269 other represented populations. These data provide a comprehensive framework of HERV-K
270 genomic diversity to advance studies on potential roles for HERV-K in human disease, which
271 have been alluring yet difficult to establish [19,20,22].

272

273 Tools developed to interrogate ERV insertional polymorphism typically exploit the
274 unique signature created by the host-virus junction [11,44,45]. These approaches indicate that
275 a site is occupied by an ERV but not whether there is a provirus associated with the site, which
276 is more difficult to accomplish with short read sequence data. Our analysis tool provides an

277 efficient means to detect occupancy and provirus status in one step. We decrease
278 computational time by analyzing only the set of reads that map to existing HERV-K in the
279 reference genome. This approach is justified because the polymorphic HERV-K that are missing
280 from the human reference are closely related to those in the reference genome assembly and
281 hence reads derived from them map to a related HERV-K in the reference. We employ *k-mer*
282 counting methods, which also increase computational efficiency. A reference set of *k-mers* that
283 is unique to a HERV-K is generated for each location in the genome and the proportion of reads
284 from the query set that maps to the *k-mer* reference set is reported as a continuous variable;
285 there is no threshold of read count or depth imposed for classification. Instead we utilize a
286 mixture model to cluster values and assign the same HERV-K status to the entire cluster.
287 Clusters with n/T of 1 have all the unique *k-mers* identified in the HERV-K reference set. We
288 classify other clusters by determining if *k-mers* mapped on the reference allele are distributed
289 at sites in the coding portion of the genome or only in the LTR. This approach led to the
290 interesting finding that several HERV-K could have population specific alleles.

291
292 Wildschutte *et al* [11] have conducted the most comprehensive study of HERV-K
293 prevalence in the KGP data to date. While the goal of that paper was to identify new
294 polymorphic insertions of either provirus or solo LTR in the KGP data, their analyses provide the
295 prevalence of some polymorphic HERV-K provirus for comparison with our results. There are
296 five HERV-K previously reported in Subramanian *et al* 2011 [10] that were not included in the
297 Wildschutte paper [11]; all are polymorphic in our analysis (range 43-99%, see Table 1 and
298 S1_Dataset:virus-column N). Seven polymorphic HERV-K, which Wildschutte *et al* [11] indicate

299 occur in greater than 98% of KGP individuals, are fixed in our study. Our estimated prevalence
300 for 14 HERV-K differs from that reported in Wildschutte et al [11] by 5% or more. Of these 14,
301 the prevalence estimates at chr1:155596457-155605636 are most divergent. Our data show
302 this site is fixed for provirus and Wildschutte *et al* [11] report that only 14% of the KGP data, all
303 from AFR, have a HERV-K provirus integration. Our plots for chr1:155596457-155605636 show
304 that AFR individuals carry the reference allele at this site (n/T near 1, Fig 5A) and all other
305 individuals have n/T near 0.5. The *k-mers* from individuals with low n/T values for
306 chr1:155596457-155605636 map to only a subset of sites marked by unique *k-mers* in the
307 coding region (S8 Fig), which is consistent with sequence polymorphism or a deletion at these
308 positions. The reference set T is small for this HERV-K and therefore overall coverage of the
309 genome is low. Because Wildschutte *et al* [11] used a minimum coverage threshold for their *k-*
310 *mer* mapping method, it is possible that alleles present in non-AFR populations would be
311 outside their inclusion criteria. There is a similar signal for alleles, represented by lower n/T
312 values, at the other 13 HERV-K sites although the differences between our prevalence
313 estimates and those of Wildschutte et al [11] are small (S1_Dataset:virus). In most cases these
314 putative alleles are found in all populations at different frequencies but in five there is some
315 degree of population specificity (Fig 5, S7 Fig, S1_Dataset:virus). Our results indicate that there
316 could be considerably more sequence variation in HERV-K among human populations than
317 previously appreciated. These data also suggest that using HERV-K consensus sequences to
318 study pathogenic potential could miss important features of HERV-K polymorphism, which can
319 be characterized by both the site occupancy status (presence/absence) and, when present, by
320 sequence differences in among individuals.

321

322 HERV-Ks are the youngest family of endogenous retroviruses in humans and
323 consequently they share considerable sequence identity. This has the effect of limiting the
324 number of unique sites associated with some HERV-K, which decreases the size of the reference
325 set T (S1_Dataset:virus). Another example of a polymorphic HERV-K with a small set T is
326 chr8:12316492, reported to be human specific (9), which shares a recent common ancestor
327 with two older HERV-K (chr8:12073970 and chr8:8054700) all located at 8p23.1. Our data
328 indicate that 14% of KGP individuals have the reference allele and most n/T values are less than
329 0.4 and fall into two non-zero clusters (S9 Fig). These appear to represent various structural
330 variations (truncation or deletion) because there are several peaks in both the LTR and in the 5'
331 coding region. Thus although an n/T ratio of 0 or 1 reliably indicates absence and presence of
332 the reference HERV-K, respectively, when T is small, sequence polymorphism and a deletion
333 event can be difficult to distinguish from a solo LTR. However, because our mixture model
334 statistically clusters similar n/T values, all individuals in a cluster have the same status (e.g allele
335 or solo LTR) even if we do not know what that state is.

336

337 Our approach provides human disease researchers with a rapid means to determine if
338 the frequency, and overall burden of the 96 HERV-K proviruses evaluated differ between a
339 patient data set and populations represented in KGP. The visualization tool allows investigators
340 to determine if HERV-Ks co-occur in certain clinical settings. The potential that HERV-K has
341 multiple allelic forms in different populations is worthy of further analysis because a sequence
342 allele could also contribute to a disease condition.

343

344 **Materials and methods**

345

346 **HERV-K proviruses**

347 The 96 HERV-K proviruses previously reported [10,11,32,46] were supplemented with
348 HERV-K alleles present in the NCBI nt database (November 2016 release). We required that any
349 allele of a HERV-K from the nt database have at least 2kb of reference-matching host flanking
350 sequence to confirm genome location. In total, 234 alleles were collected at the 96 known
351 HERV-K loci (92 in hg19, and 4 from the nt database). The location information and virus
352 features are summarized in S1_Dataset: virus.

353

354 **Developing a *k-mer* based detection model**

355 We identified the *k-mers* that correspond to unique sequence characterizing each HERV-
356 K. *K-mers* are substrings (subsequences) of length *k* that exist in a string (DNA sequence). The
357 length *k* is determined empirically (S1 Method). Each *k-mer* is labeled with the corresponding
358 viruses in which it is observed.

359 Only those *k-mers* referring to a single virus, unique *k-mers*, are selected for the set T. Where
360 multiple alleles of a HERV-K are available, *k-mers* unique to all alleles at that location comprise
361 T. Multiple 2bps different *k-mers* (such as SNPs) corresponding to the same location on the
362 virus, are merged into a single entry for the purposes of computing T. We map unique *k-mers*
363 back to the corresponding alleles to determine depth of the HERV-K (S3 Fig) and whether *k-*
364 *mers* are located in LTRs. (S1_Dataset: virus)

365

366

367 **Analysis of 1000 Genome Project (KGP) Data**

368 To develop a method to recover sequences containing information on HERV-K we
369 leverage the fact that HERV-Ks are closely related. Thus, most sequence reads obtained from an
370 individual with a polymorphic HERV-K that is absent in the human reference, hg19, will map to
371 the location of one of the closely related HERV-K that is present in the human genome reference.
372 A file with the coordinates for all reported HERV-K insertions is used to extract mapped reads
373 from a genome sequence file (S1_Dataset:bed, which provides the coordinates for both hg19
374 and hg38). Note that the KGP data were mapped to GRCh37, which includes the decoy
375 sequence hs37d5. This decoy contains the HERV-K at chr1:73594980_73595948 and is not
376 present in hg19. Thus, we did not recover any reads for this HERV-K, which is polymorphic but
377 reportedly at high prevalence in most populations [11].

378

379 The KGP data were downloaded in aligned Binary Alignment/Map (BAM) format
380 (<ftp://ftp.ncbi.nlm.nih.gov/1000genomes/ftp/data/>). It contains data for 2,535 individuals
381 (S1_Dataset:KGP) sequenced via low-depth whole-genome sequencing (mean depth = 6.98X).
382 The individuals represent 26 populations, derived from 5 super-populations, including African
383 (AFR), Admixed America (AMR), East Asian (EAS), European (EUR), and South Asian (SAS)
384 [47,48]. Of 2,535 individuals, 28 also have high-depth DNA sequences (mean depth = 48.06X),
385 which we use as a pilot dataset for our clustering methods, described below and in
386 Supplementary Methods.

387

388 Our computational framework to indicate the status of each known HERV-K provirus is
389 based on the n/T ratio, which is the proportion of k -mers from each individual that are identical
390 to the reference set T for each HERV-K. Sequence reads are extracted from a mapped file of
391 whole human genome sequence data based on coordinates corresponding to each annotated
392 HERV-K. The reads are k -merized and mapped to the set T, which represents all unique k -mers
393 assigned to each HERV-K in the reference set. We use exact match to map k -mers data set to
394 the unique k -mers references. The n/T ratio is used as an indicator of the presence of each
395 HERV-K. Using a hash table (S5 Method), it takes 15 minutes to generate the n/T matrix for 100
396 files. The source code for the entire process is at <https://github.com/lwl1112/polymorphicHERV>

397

398

399 **Dirichlet process Gaussian mixture model (DPGMM)**

400 Because n (the number of k -mers obtained from a persons' sequence data, that map to
401 a specific HERV-K) is affected by sequencing depth, we utilized a statistical model to cluster n/T
402 for each HERV-K for each individual and assigned HERV-K status to a cluster. Mixture modeling
403 is arguably the most widely used statistical method for clustering. In this analysis, we follow the
404 work proposed by Lin et al. 2015 [49], which employs a Gaussian Mixture Model (GMM) with
405 density function given by

$$406 \quad f(x | \theta) = \sum_{j=1}^M \pi_j N(\mu_j, \Sigma_j), \quad (1)$$

407 and with a mixture of relatively large number M Gaussian components (denoted by $N(\mu_j, \Sigma_j)$,
408 for $j = 1:M$) to represent the data density. To allow a flexible modeling approach, we employ
409 the standard Bayesian (truncated) Dirichlet Process prior for the parameters $\theta = (\pi_j, \mu_j, \Sigma_j$
410 , $j = 1:M$ [50,51]. The idea is that some of the mixture probabilities (π_j) can be zero, hence the
411 actual number of mixture components needed may be smaller than the upper bound M . This
412 mechanism allows automatic determination of the number of mixture components needed by
413 the data set at hand. Given a fitted model via the Bayesian expectation–maximization
414 algorithm, in terms of estimates of all parameters θ , we identify clusters by aggregating
415 Gaussian components. Merging components into clusters can be done by associating each of
416 the Gaussian components to the closest mode of $f(x|\theta)$. Hence, the number of modes identified
417 is the realized number of clusters. [See S2 Method for full detail]

418

419 **Co-occurrence of polymorphic HERV-K**

420 We consider that both the individual frequency of a HERV-K and the co-occurrence of
421 multiple HERV-K could differ among populations.

422 The time of a brute-force approach for finding all combinations C_m of size m from p

423 polymorphic HERV-K is $(\sum_{m=1}^p \binom{p}{m} = 2^p - 1)$, which is not efficient and redundant. We

424 employed the Apriori algorithm [52], which is commonly used for finding frequent pattern sets;

425 in our case indicating which polymorphic HERV-K frequently appear together. It first generates

426 combinations C_m (initialized to 1). In the optimization, frequent combinations F_m are returned

427 from candidates C_m when frequency exceeds the minimum threshold of co-occurrence. F_m are

428 then self-joined to generate combinations C_{m+1} of size $m+1$ and out of which F_{m+1} satisfy the
429 minimum co-occurrence. In each pass, candidate combinations are pruned so as to avoid
430 generating all combinations, which reduces running time significantly.

431

432

433 **Statistical analysis of HERV-K frequencies across populations**

434 We make statistical comparisons across 5 super-populations for the following three
435 problems. For each problem, there are $\binom{5}{2} = 10$ families of 1-to-1 comparisons conducted. The
436 ‘prop-test’ function in R is used to test whether the proportions for two super-populations are
437 the same.

- 438 1) individual prevalence of polymorphic HERV-K. (20 comparisons for each polymorphic HERV-K
439 in a family)
- 440 2) the number of polymorphic HERV-K present per individual. (21 comparisons as the number
441 of co-occurring polymorphic HERV-K is from 0 to 20)
- 442 3) the co-occurrence for combinations of polymorphic HERV-K.

443 Therefore, multiple hypotheses would be conducted on frequencies F across super-populations
444 $P_{1...5}$ as follows:

445 Null hypothesis, $H_0: F_{P_i} = F_{P_j}$, where $i \neq j$;

446 Alternative hypothesis, $H_A: F_{P_i} \neq F_{P_j}$, where $i \neq j$.

447 A separate P-value is computed for each test and the Benjamini-Hochberg procedure [53] is
448 used to account for multiple comparisons.

449

450 **Visualization in D3.js**

451 We utilized D3.js (Data Driven Documents) [54], an open-source java script library to
452 create an interactive visualization to display co-occurrence of polymorphic HERV-Ks in human
453 populations. Our visualization system includes two modules, a welcome page and a result page.
454 Input JSON data include locations of polymorphic HERV-K, population information, and the 0/1
455 (absence / presence) matrix. (See S3 Method). Source code is available at:

456 <https://github.com/lwl1112/polymorphicHERV/tree/master/visualization>

457

458

459 **Acknowledgements**

460 WL was supported in part by the Louis S. and Sara S. Michael Endowed Graduate
461 Fellowship in Engineering and the Fred A. and Susan Breidenbach Graduate Fellowship in
462 Engineering.

463

464

465 **References**

466

- 467 1. Hayward A, Grabherr M, Jern P. Broad-scale phylogenomics provides insights into retrovirus-host
468 evolution. *Proc Natl Acad Sci.* 2013;110: 20146–20151. doi:10.1073/pnas.1315419110
- 469 2. Feschotte C, Gilbert C. Endogenous viruses: insights into viral evolution and impact on host
470 biology. *Nat Rev Genet.* 2012;13: 283–296. doi:10.1038/nrg3199
- 471 3. Stoye JP. Studies of endogenous retroviruses reveal a continuing evolutionary saga. *Nat Rev*

- 472 Microbiol. 2012;10: 395–406. doi:10.1038/nrmicro2783
- 473 4. Gifford R, Tristem M. The evolution, distribution and diversity of endogenous retroviruses. *Virus*
474 *Genes*. 2003;26: 291–315. Available: <http://www.ncbi.nlm.nih.gov/pubmed/12876457>
- 475 5. Weiss RA. The discovery of endogenous retroviruses. *Retrovirology*. 2006;3: 67.
476 doi:10.1186/1742-4690-3-67
- 477 6. Jern P, Coffin JM. Effects of Retroviruses on Host Genome Function. *Annu Rev Genet*. 2008;42:
478 709–732. doi:10.1146/annurev.genet.42.110807.091501
- 479 7. Löwer R, Löwer J, Kurth R. The viruses in all of us: characteristics and biological significance of
480 human endogenous retrovirus sequences. *Proc Natl Acad Sci U S A*. 1996;93: 5177–84.
481 Available: <http://www.ncbi.nlm.nih.gov/pubmed/8643549>
- 482 8. Bannert N, Kurth R. Retroelements and the human genome: New perspectives on an old relation.
483 *Proc Natl Acad Sci*. 2004;101: 14572–14579. doi:10.1073/pnas.0404838101
- 484 9. Moyes D, Griffiths DJ, Venables PJ. Insertional polymorphisms: a new lease of life for endogenous
485 retroviruses in human disease. *Trends Genet*. 2007;23: 326–333. doi:10.1016/j.tig.2007.05.004
- 486 10. Subramanian RP, Wildschutte JH, Russo C, Coffin JM. Identification, characterization, and
487 comparative genomic distribution of the HERV-K (HML-2) group of human endogenous
488 retroviruses. *Retrovirology*. 2011;8: 90. doi:10.1186/1742-4690-8-90
- 489 11. Wildschutte JH, Williams ZH, Montesion M, Subramanian RP, Kidd JM, Coffin JM. Discovery of
490 unfixed endogenous retrovirus insertions in diverse human populations. *Proc Natl Acad Sci*.
491 2016;113: E2326–E2334. doi:10.1073/pnas.1602336113
- 492 12. Belshaw R, Watson J, Katzourakis A, Howe A, Woolven-Allen J, Burt A, et al. Rate of
493 Recombinational Deletion among Human Endogenous Retroviruses. *J Virol*. 2007;81: 9437–9442.
494 doi:10.1128/JVI.02216-06
- 495 13. Medstrand P, Mager DL. Human-specific integrations of the HERV-K endogenous retrovirus

- 496 family. *J Virol.* 1998;72: 9782–7. Available: <http://www.ncbi.nlm.nih.gov/pubmed/9811713>
- 497 14. Hughes JF, Coffin JM. Human endogenous retrovirus K solo-LTR formation and insertional
498 polymorphisms: Implications for human and viral evolution. *Proc Natl Acad Sci.* 2004;101: 1668–
499 1672. doi:10.1073/pnas.0307885100
- 500 15. Belshaw R, Dawson ALA, Woolven-Allen J, Redding J, Burt A, Tristem M. Genomewide Screening
501 Reveals High Levels of Insertional Polymorphism in the Human Endogenous Retrovirus Family
502 HERV-K(HML2): Implications for Present-Day Activity. *J Virol.* 2005;79: 12507–12514.
503 doi:10.1128/JVI.79.19.12507-12514.2005
- 504 16. Marchi E, Kanapin A, Magiorkinis G, Belshaw R. Unfixed Endogenous Retroviral Insertions in the
505 Human Population. *J Virol.* 2014;88: 9529–9537. doi:10.1128/JVI.00919-14
- 506 17. Shin W, Lee J, Son S-Y, Ahn K, Kim H-S, Han K. Human-Specific HERV-K Insertion Causes
507 Genomic Variations in the Human Genome. Cordaux R, editor. *PLoS One.* 2013;8: e60605.
508 doi:10.1371/journal.pone.0060605
- 509 18. Gröger V, Cynis H. Human Endogenous Retroviruses and Their Putative Role in the Development
510 of Autoimmune Disorders Such as Multiple Sclerosis. *Front Microbiol.* 2018;9.
511 doi:10.3389/fmicb.2018.00265
- 512 19. Young GR, Stoye JP, Kassiotis G. Are human endogenous retroviruses pathogenic? An approach
513 to testing the hypothesis. *BioEssays.* 2013;35: 794–803. doi:10.1002/bies.201300049
- 514 20. Ryan FP. Human endogenous retroviruses in health and disease: a symbiotic perspective. *JRSM.*
515 2004;97: 560–565. doi:10.1258/jrsm.97.12.560
- 516 21. Volkman HE, Stetson DB. The enemy within: endogenous retroelements and autoimmune
517 disease. *Nat Immunol.* 2014;15: 415–422. doi:10.1038/ni.2872
- 518 22. Magiorkinis G, Belshaw R, Katzourakis A. “There and back again”: revisiting the
519 pathophysiological roles of human endogenous retroviruses in the post-genomic era. *Philos Trans*

- 520 R Soc B Biol Sci. 2013;368: 20120504–20120504. doi:10.1098/rstb.2012.0504
- 521 23. Löwer R. The pathogenic potential of endogenous retroviruses: facts and fantasies. Trends
522 Microbiol. 1999;7: 350–356. doi:10.1016/S0966-842X(99)01565-6
- 523 24. Hohn O, Hanke K, Bannert N. HERV-K(HML-2), the Best Preserved Family of HERVs:
524 Endogenization, Expression, and Implications in Health and Disease. Front Oncol. 2013;3.
525 doi:10.3389/fonc.2013.00246
- 526 25. Hughes JF. Human Endogenous Retroviral Elements as Indicators of Ectopic Recombination
527 Events in the Primate Genome. Genetics. 2005;171: 1183–1194.
528 doi:10.1534/genetics.105.043976
- 529 26. Hughes JF, Coffin JM. Evidence for genomic rearrangements mediated by human endogenous
530 retroviruses during primate evolution. Nat Genet. 2001;29: 487–489. doi:10.1038/ng775
- 531 27. Romanish MT, Cohen CJ, Mager DL. Potential mechanisms of endogenous retroviral-mediated
532 genomic instability in human cancer. Semin Cancer Biol. 2010;20: 246–253.
533 doi:10.1016/j.semcancer.2010.05.005
- 534 28. Kamp C, Hirschmann P, Voss H, Huellen K, Vogt PH. Two long homologous retroviral sequence
535 blocks in proximal Yq11 cause AZFa microdeletions as a result of intrachromosomal
536 recombination events. Hum Mol Genet. 2000;9: 2563–72. Available:
537 <http://www.ncbi.nlm.nih.gov/pubmed/11030762>
- 538 29. Kidd JM, Graves T, Newman TL, Fulton R, Hayden HS, Malig M, et al. A Human Genome
539 Structural Variation Sequencing Resource Reveals Insights into Mutational Mechanisms. Cell.
540 2010;143: 837–847. doi:10.1016/j.cell.2010.10.027
- 541 30. Cohen CJ, Lock WM, Mager DL. Endogenous retroviral LTRs as promoters for human genes: A
542 critical assessment. Gene. 2009;448: 105–114. doi:10.1016/j.gene.2009.06.020
- 543 31. Simmons W. The Role of Human Endogenous Retroviruses (HERV-K) in the Pathogenesis of

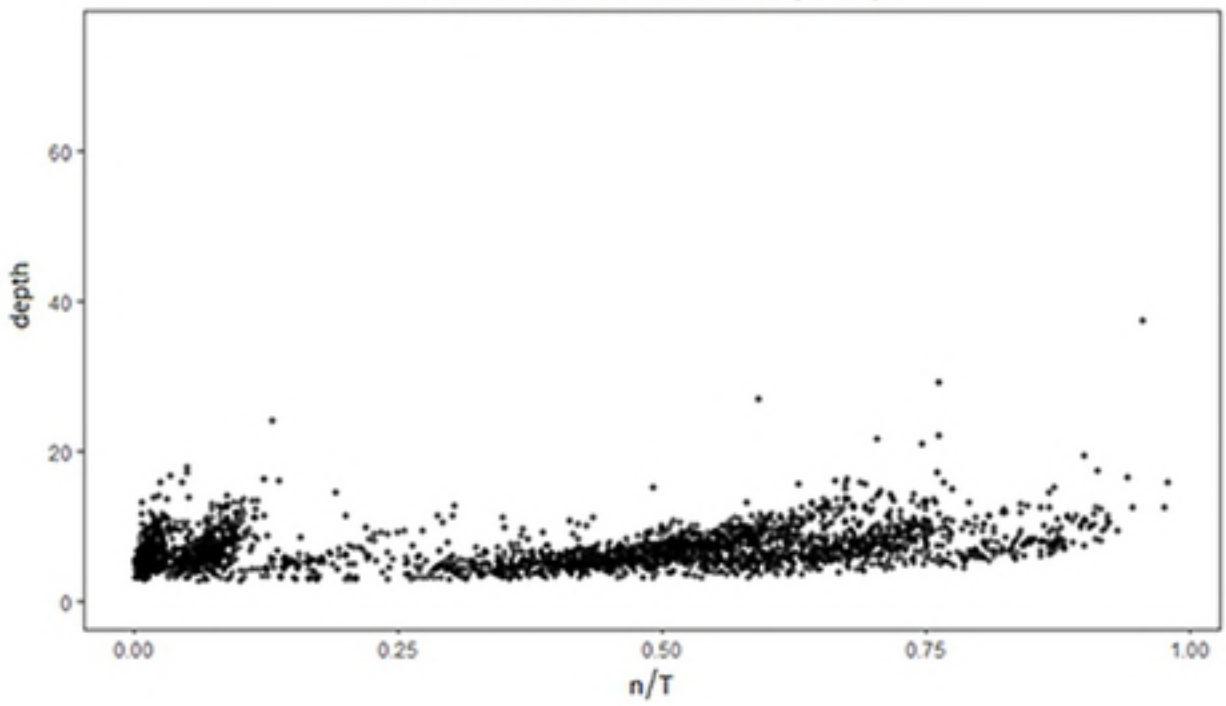
- 544 Human Cancers. *Mol Biol.* 2016;05. doi:10.4172/2168-9547.1000169
- 545 32. Wildschutte JH, Ram D, Subramanian R, Stevens VL, Coffin JM. The distribution of insertionally
546 polymorphic endogenous retroviruses in breast cancer patients and cancer-free controls.
547 *Retrovirology.* 2014;11: 62. doi:10.1186/s12977-014-0062-3
- 548 33. Kassiotis G, Stoye JP. Making a virtue of necessity: the pleiotropic role of human endogenous
549 retroviruses in cancer. *Philos Trans R Soc B Biol Sci.* 2017;372: 20160277.
550 doi:10.1098/rstb.2016.0277
- 551 34. Johanning GL, Malouf GG, Zheng X, Esteva FJ, Weinstein JN, Wang-Johanning F, et al.
552 Expression of human endogenous retrovirus-K is strongly associated with the basal-like breast
553 cancer phenotype. *Sci Rep.* 2017;7: 41960. doi:10.1038/srep41960
- 554 35. Bhardwaj N, Coffin JM. Endogenous Retroviruses and Human Cancer: Is There Anything to the
555 Rumors? *Cell Host Microbe.* 2014;15: 255–259. doi:10.1016/j.chom.2014.02.013
- 556 36. Hanke K, Hohn O, Bannert N. HERV-K(HML-2), a seemingly silent subtenant - but still waters run
557 deep. *APMIS.* 2016;124: 67–87. doi:10.1111/apm.12475
- 558 37. Trela M, Nelson PN, Rylance PB. The role of molecular mimicry and other factors in the
559 association of Human Endogenous Retroviruses and autoimmunity. *APMIS.* 2016;124: 88–104.
560 doi:10.1111/apm.12487
- 561 38. Antony JM, DesLauriers AM, Bhat RK, Ellestad KK, Power C. Human endogenous retroviruses
562 and multiple sclerosis: Innocent bystanders or disease determinants? *Biochim Biophys Acta - Mol*
563 *Basis Dis.* 2011;1812: 162–176. doi:10.1016/j.bbadis.2010.07.016
- 564 39. Tugnet N. Human Endogenous Retroviruses (HERVs) and Autoimmune Rheumatic Disease: Is
565 There a Link? *Open For Sci J.* 2013;7: 13–21. doi:10.2174/1874312901307010013
- 566 40. Li W, Lee M-H, Henderson L, Tyagi R, Bachani M, Steiner J, et al. Human endogenous retrovirus-
567 K contributes to motor neuron disease. *Sci Transl Med.* 2015;7: 307ra153-307ra153.

- 568 doi:10.1126/scitranslmed.aac8201
- 569 41. Douville RN, Nath A. Human Endogenous Retrovirus-K and TDP-43 Expression Bridges ALS and
570 HIV Neuropathology. *Front Microbiol.* 2017;8. doi:10.3389/fmicb.2017.01986
- 571 42. Nexø BA, Villesen P, Nissen KK, Lindegaard HM, Rossing P, Petersen T, et al. Are human
572 endogenous retroviruses triggers of autoimmune diseases? Unveiling associations of three
573 diseases and viral loci. *Immunol Res.* 2016;64: 55–63. doi:10.1007/s12026-015-8671-z
- 574 43. Fukunaga K. Introduction to statistical pattern recognition. Academic press; 1990.
- 575 44. Ciuffi A, Ronen K, Brady T, Malani N, Wang G, Berry CC, et al. Methods for integration site
576 distribution analyses in animal cell genomes. *Methods.* 2009;47: 261–268.
577 doi:10.1016/j.ymeth.2008.10.028
- 578 45. Witherspoon DJ, Xing J, Zhang Y, Watkins WS, Batzer MA, Jorde LB. Mobile element scanning
579 (ME-Scan) by targeted high-throughput sequencing. *BMC Genomics.* 2010;11: 410.
580 doi:10.1186/1471-2164-11-410
- 581 46. Bhardwaj N, Montesion M, Roy F, Coffin J. Differential Expression of HERV-K (HML-2) Proviruses
582 in Cells and Virions of the Teratocarcinoma Cell Line Tera-1. *Viruses.* 2015;7: 939–968.
583 doi:10.3390/v7030939
- 584 47. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated
585 map of structural variation in 2,504 human genomes. *Nature.* 2015;526: 75–81.
586 doi:10.1038/nature15394
- 587 48. Gibbs RA, Boerwinkle E, Doddapaneni H, Han Y, Korchina V, Kovar C, et al. A global reference
588 for human genetic variation. *Nature.* 2015;526: 68–74. doi:10.1038/nature15393
- 589 49. Lin L, Chan C, West M. Discriminative variable subsets in Bayesian classification with mixture
590 models, with application in flow cytometry studies. *Biostatistics.* 2015; kxv021.
591 doi:10.1093/biostatistics/kxv021

- 592 50. Escobar MD, West M. Bayesian Density Estimation and Inference Using Mixtures. *J Am Stat*
593 *Assoc.* 1995;90: 577. doi:10.2307/2291069
- 594 51. Ishwaran H, James LF. Gibbs Sampling Methods for Stick-Breaking Priors. *J Am Stat Assoc.*
595 2001;96: 161–173. doi:10.1198/016214501750332758
- 596 52. Huang L, Chen H, Wang X, Chen G. A fast algorithm for mining association rules. *J Comput Sci*
597 *Technol.* 2000;15: 619–624. doi:10.1007/BF02948845
- 598 53. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach
599 to multiple testing. *J R Stat Soc Ser B.* 1995; 289–300. doi:10.2307/2346101
- 600 54. Bostock M, Ogievetsky V, Heer J. D³ Data-Driven Documents. *IEEE Trans Vis Comput Graph.*
601 2011;17: 2301–2309. doi:10.1109/TVCG.2011.185
- 602
- 603

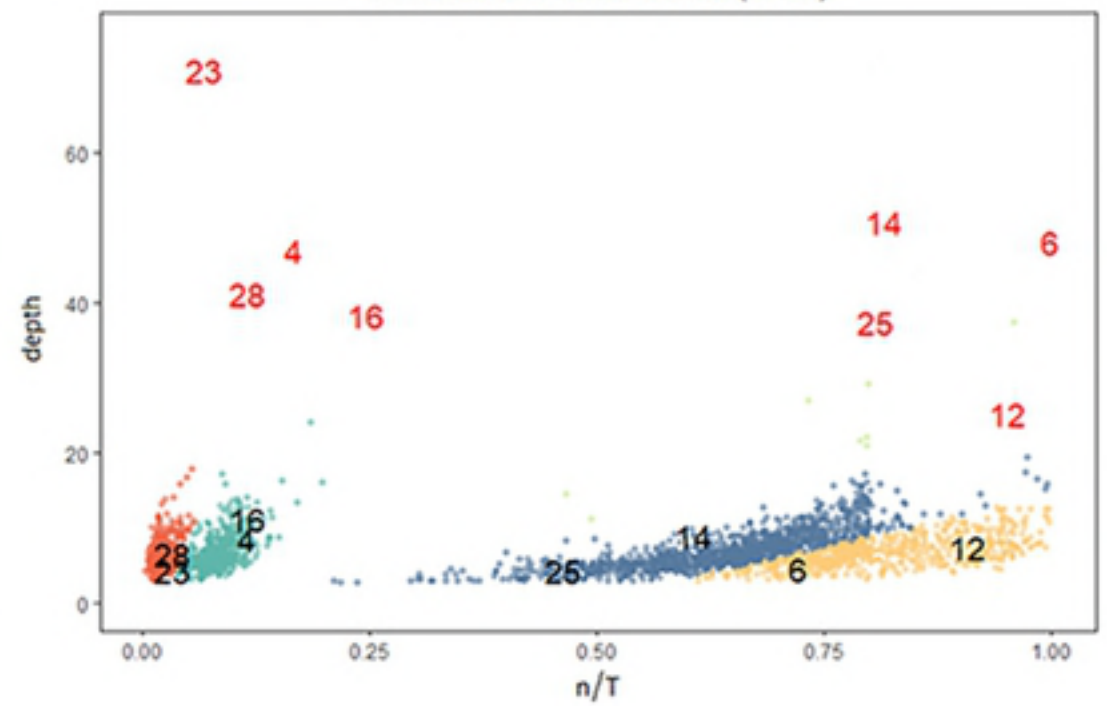
(A)

chr12:55727215-55728183 (K=70)

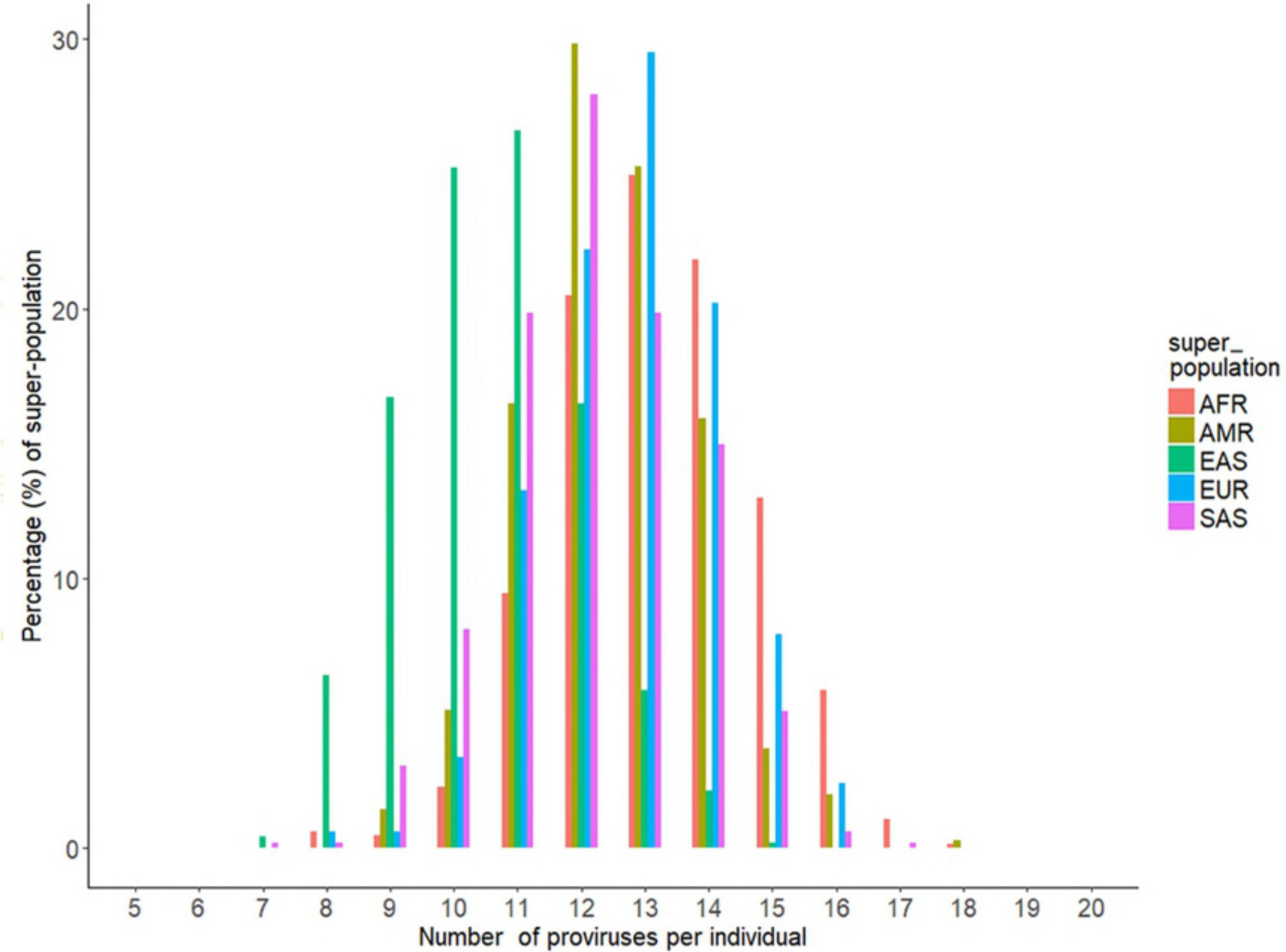


(B)

chr12:55727215-55728183 (K=50)

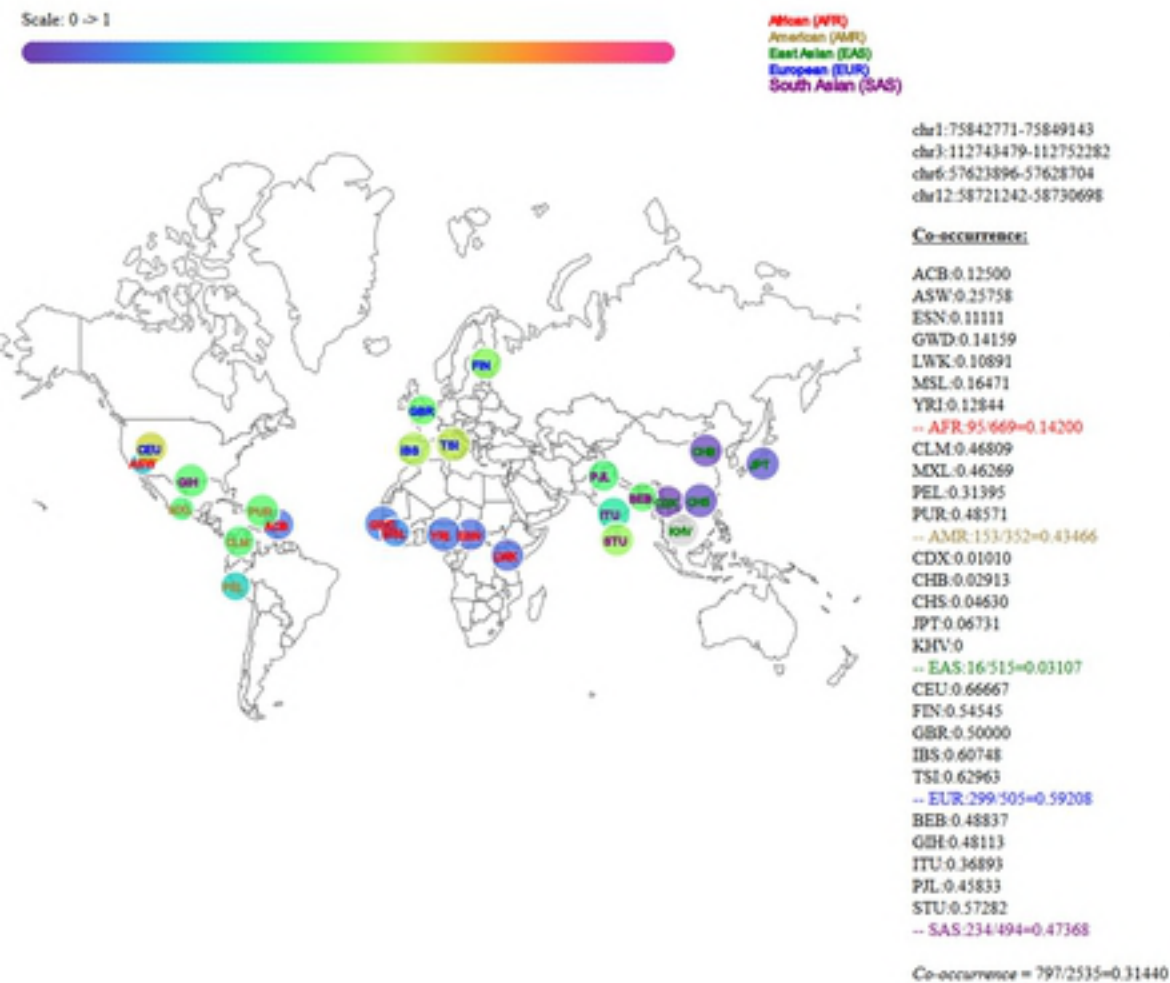


Figure

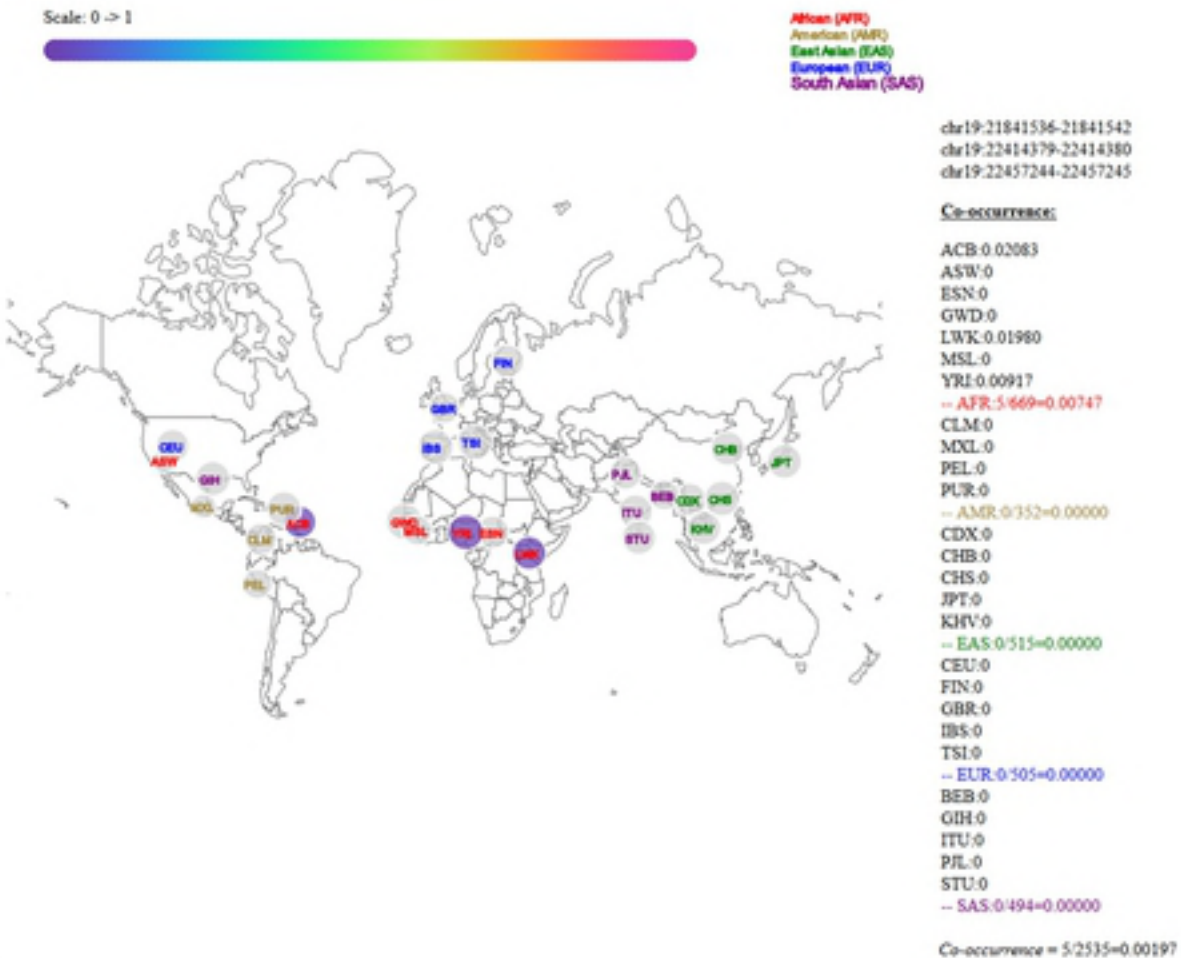


Figure

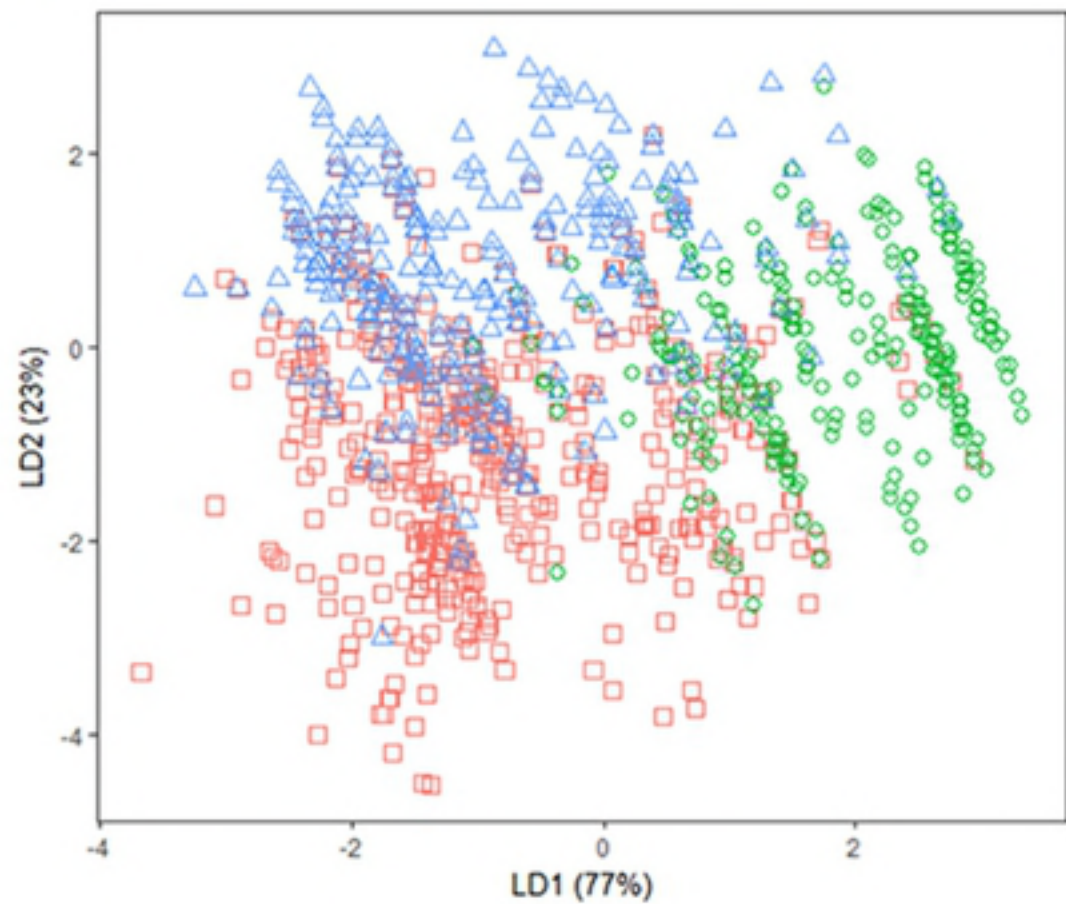
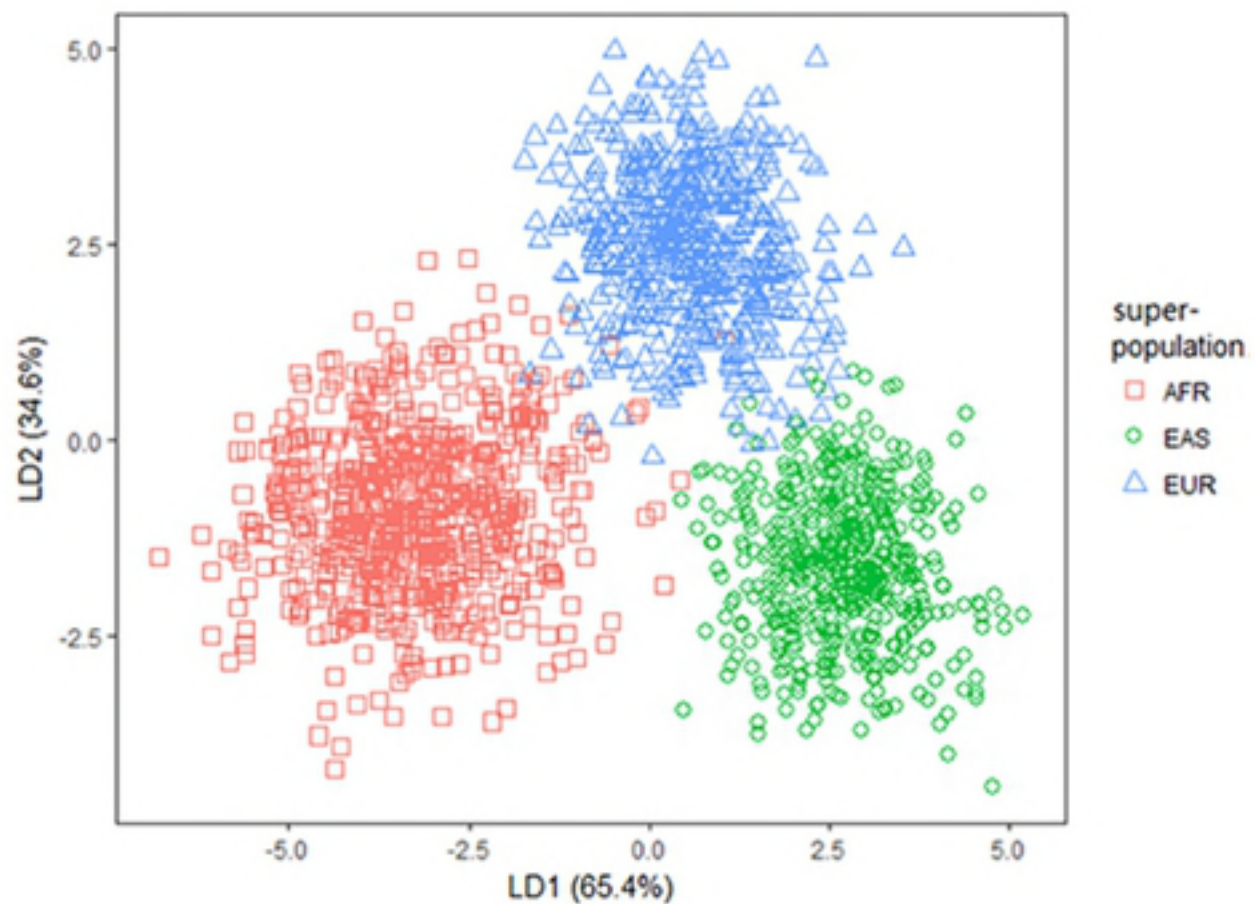
A)

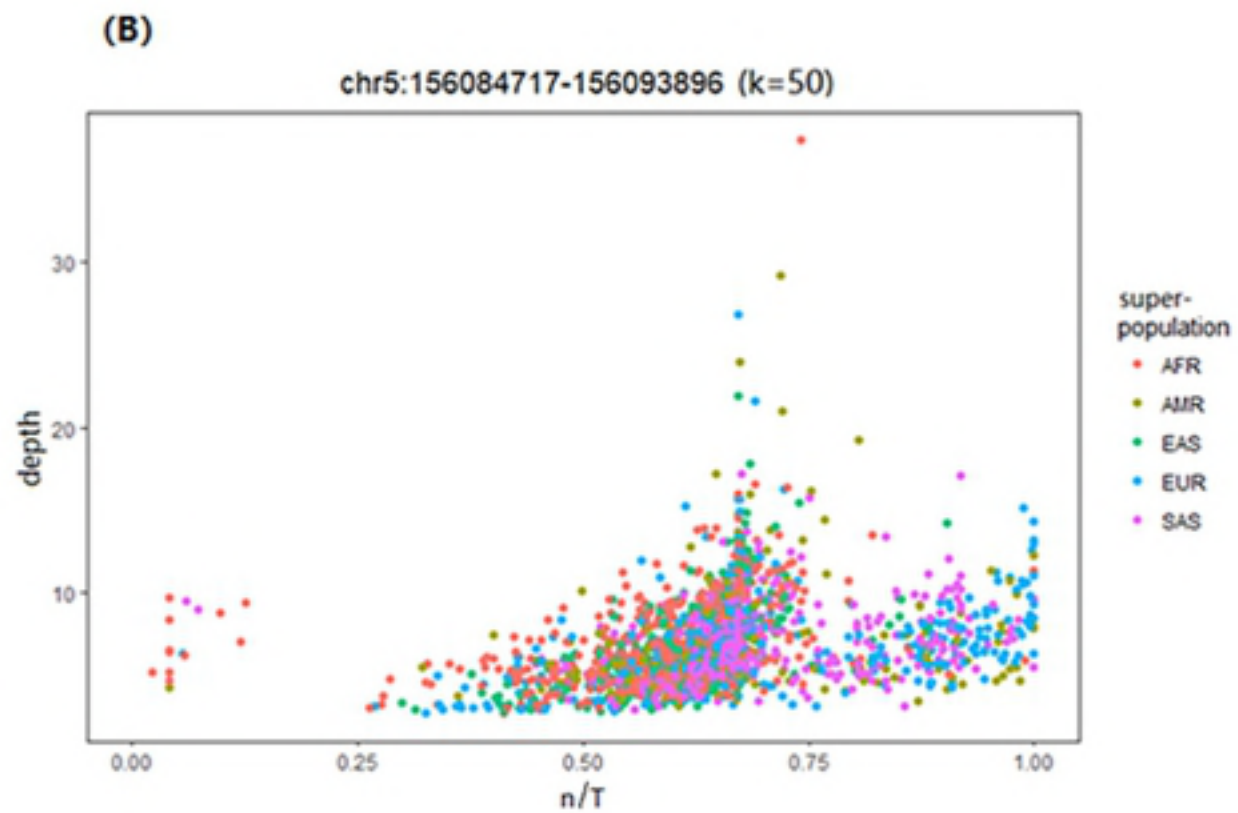
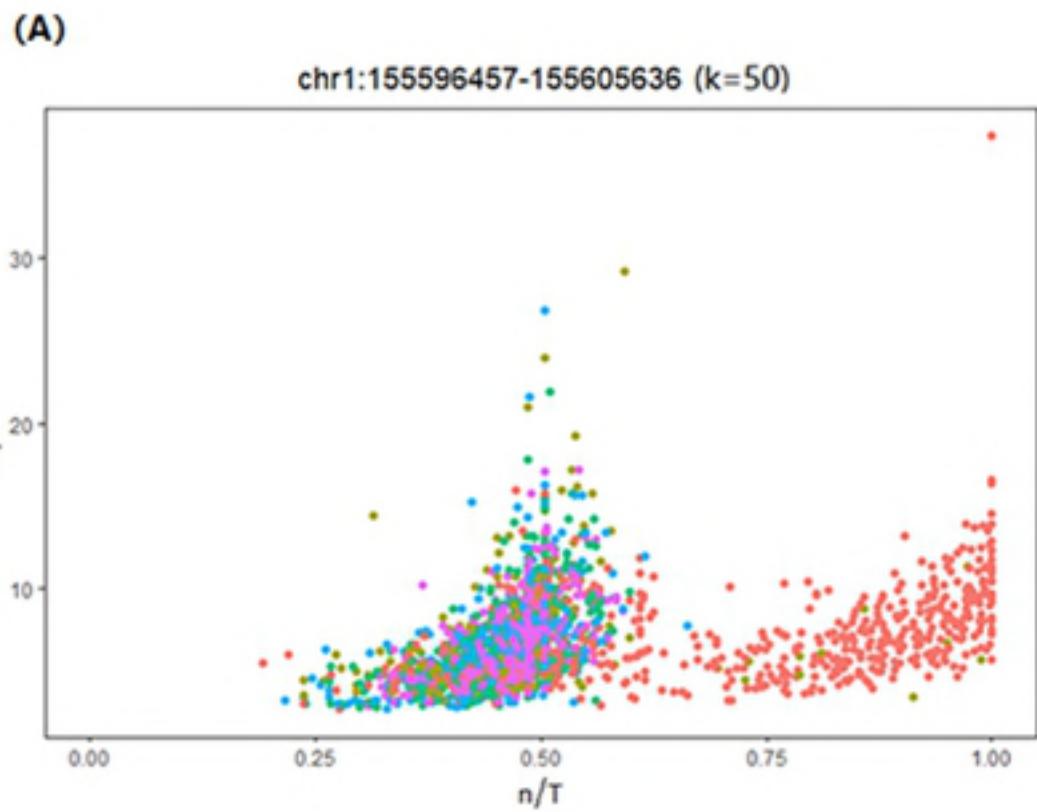


B)



Figure

(A)**(B)****Figure**



Figure