1    **TITLE: Virus-derived variation in diverse human genomes**

2

3    Shohei Kojima, Anselmo Jiro Kamada, Nicholas F. Parrish*

4

5    Genome Immunobiology RIKEN Hakubi Research Team, RIKEN Center for Integrative Medical

6    Sciences and RIKEN Cluster for Pioneering Research, Yokohama, Japan 230-0045

7

8    *corresponding author: nicholas.parrish@riken.jp

9

10

11

**Abstract**

Acquisition of genetic material from viruses by their hosts can generate inter-host structural genome variation. We developed computational tools enabling us to study virus-derived structural variants (SVs) in population-scale whole genome sequencing (WGS) datasets and applied them to 3,332 humans. Although SVs had already been cataloged in these subjects, we found previously-overlooked virus-derived SVs. We detected somatic SVs present in the sequenced lymphoblastoid cell lines (LCLs) derived from squirrel monkey retrovirus (SMRV), human immunodeficiency virus 1 (HIV-1), and human T lymphotropic virus (HTLV-1); these variants are attributable to infection of LCLs or their progenitor cells and may impact gene expression results and the biosafety of experiments using these cells. In addition, we detected new heritable SVs derived from human herpesvirus 6 (HHV-6) and human endogenous retrovirus-K (HERV-K). We report the first solo-DR HHV-6 that likely to reflects rearrangement of a known full-length endogenous HHV-6. We used linkage disequilibrium between single nucleotide variants (SNVs) and variants in reads that align to HERV-K, which often cannot be mapped uniquely using conventional short-read sequencing analysis methods, to locate previously-unknown polymorphic HERV-K loci. Some of these loci are tightly linked to trait-associated SNVs, some are in complex genome regions inaccessible to prior methods, and some contain novel HERV-K haplotypes likely derived from gene conversion from an unknown source or introgression. These tools and results broaden our perspective on the coevolution between viruses and humans, including ongoing virus-to-human gene transfer contributing to genetic variation between humans.

(241 words)

**Introduction**

Union of genomes from discrete biological entities is a major engine of genetic diversity. Fusion of gametes, each bearing a set of recombinant chromosomes, is the immediate source of the genetic material that uniquely identifies each human. Taking a wider viewpoint, much of human genome can be recognized to have been acquired from a non-human source. For example, about 2% of the genome of many living humans can be attributed to introgression from Neanderthals (1). Movement of genetic information between biological entities apart from sexual reproduction, known as horizontal gene transfers (HGT), has also occurred in the human lineage. Some HGT happened so long ago that it is difficult to accurately classify the entity contributing the horizontally-transferred sequences according to extant taxonomies. This case for the bacteria, acquired millennia ago, now represented as our mitochondrial genomes. Other HGT occurred more recently. For example, about 8% of human genetic material is derived from human endogenous retroviruses (HERV) that integrated into our ancestors' germline and then developed an intracellular replication cycle; some HERVs integrated recently enough that they can be classified based on homology to extant exogenous retroviruses.

The most recent among these retroviral integrations, of a lysine tRNA-primed HERV (i.e. HERV-K) subgroup called HML-2, occurred less than a million years ago (2). A single HERV-K element showing insertional polymorphism in different humans was known at the time of completion of the draft human genome (3), but during the past 20 years, over 40 insertionally-polymorphic elements have been described (2, 4–9). In addition to retroviruses, sequences from ancient relatives of Borna disease virus, an RNA-only virus, were horizontally acquired in the haplorrhine lineage (10). Human herpesviruses 6A and 6B, double stranded DNA viruses, have also been horizontally transferred to some human genomes during the holocene (11–13). These observations show that viruses acquired during the lifespan of an individual organism, including humans, have sometimes contributed to the genetic material passed on to their offspring, seemingly in violation of Weismann's proposed barrier between soma and germline (14). When these viral sequences are acquired, the resulting mutation would be classified as a structural variant, defined as a DNA rearrangement greater than 50 nucleotides in length. Structural variants in human genomes are increasingly characterized at population scale (15–17). In these studies, SVs caused by polymorphic insertion of mobile genetic elements classified as transposons (including Alu, LINE-1, and SVA) have been considered explicitly. On the other hand, structural variants derived from viruses, another potentially important class of mobile genetic elements, have yet to be analyzed comprehensively.

3

68    Here we designed and applied new tools to comprehensively assess virus-derived

69    structural variants in short-read genome sequencing data at population scale. We are not the

70    first to consider viral sequences present in shotgun WGS datasets, however others have done

71    so under the assumption that viral reads reflect a somatically-acquired "virome," similar to the

72    bacterial microbiome (18, 19). To distinguish exogenous virus contamination from germline

73    integration, we applied several criteria, including read depth relative to autosomal genes and

74    patterns of linkage disequilibrium with SNVs. Although we used human WGS datasets which

75    have already been deeply analyzed to establish global SV references (15, 17), we discovered

76    previously-undescribed heritable SVs derived from virus-origin genetic material. We detect

77    squirrel monkey retrovirus (SMRV), human immunodeficiency virus 1 (HIV-1), and human T

78    lymphotropic virus 1 (HTLV-1) in LCLs widely distributed as reference materials for

79    characterizing human genetic and phenotypic variation, raising both biosafety and

80    reproducibility concerns. We developed a new approach to detect and map polymorphisms in

81    HERV-K that allows us to infer polymorphisms at over 60 loci previously unknown to be

82    polymorphic, including new loci associated with human phenotypes. We show that viruses

83    contribute unexpectedly to human genome structural variation and describe new tools for
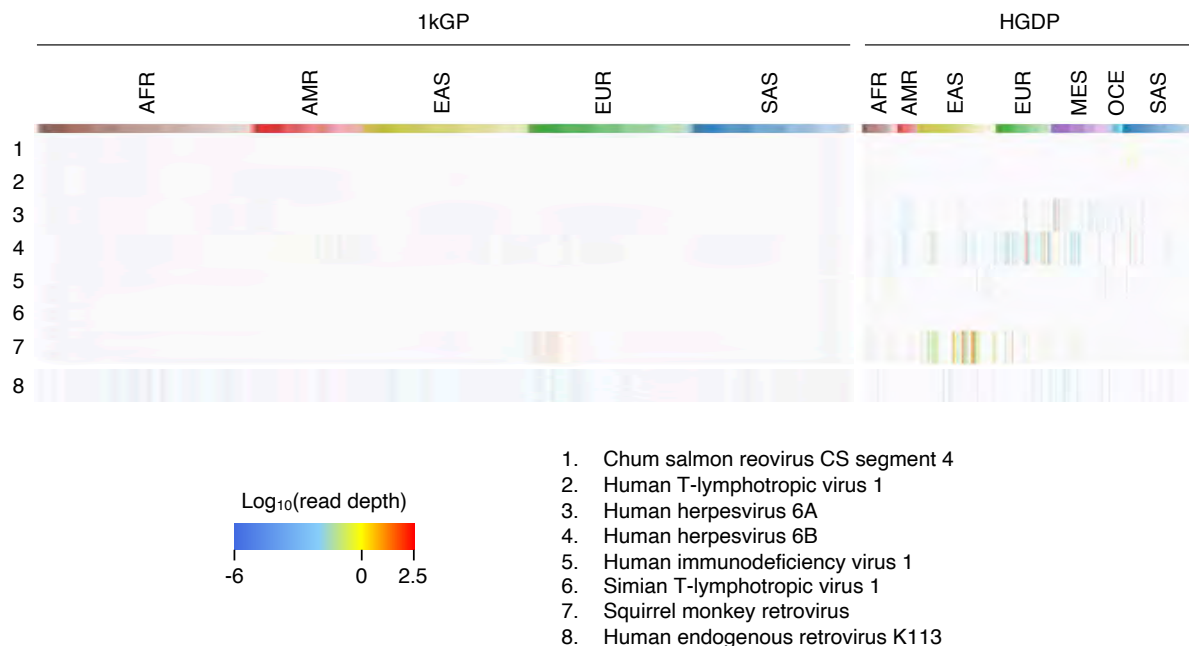
84    analyzing these variants at increasing scales.

85

86

**Results**

*Detection of virus-mapped reads from diverse human populations*

To discover human structural variation derived from viral sequences, we analyzed WGS reads that failed to map to the reference human genome (GRCh38DH). We used 3,332 high-coverage WGS datasets from the 1,000 Genomes Project (1kGP) and the Human Genome Diversity Project (HGDP) (20, 21), all derived from lymphoblastoid cell lines (LCLs). Unmapped reads were re-mapped to reference virus genomes from NCBI (see methods). We focused on viruses with abundantly-mapped reads, requiring that 5% of a viral genome be covered at more than 2x read depth. Applying this filter, we detected 7 high-coverage viruses (Figure 1, Supplementary Figure 1). Next we checked the patterns of viral genome coverage using the plots automatically generated as an output from our tool. We detected reads mapping to Chum salmon reovirus in 2 datasets from individuals of South Asian ancestry; only the first 200-bp of the viral genome was covered by reads, which were abundant in these two datasets but absent in others (Supplementary Figure 2). We detected reads mapping to simian T-lymphotropic virus 1 (STLV-1) in 4 samples. STLV-1-mapped reads were only detected in the datasets in which HTLV-1-mapped reads were also found (see below), and the same reads mapped to both HTLV-1 and STLV-1 (Supplementary Figure 2). In contrast, reads from SMRV, HIV-1, HTLV-1, human herpesvirus 6A (HHV-6A) and human herpesvirus 6B (HHV-6B) were abundantly detected in at least one subject, potentially consistent with presence in the germline, and reads covered the entire viral genome.



Log$_{10}$(read depth)

-6        0   2.5

1. Chum salmon reovirus CS segment 4
2. Human T-lymphotropic virus 1
3. Human herpesvirus 6A
4. Human herpesvirus 6B
5. Human immunodeficiency virus 1
6. Simian T-lymphotropic virus 1
7. Squirrel monkey retrovirus
8. Human endogenous retrovirus K113

5

109  **Figure 1** Virus search from 3,332 WGS.
110  Heatmap shows read depth of seven viruses with abundant reads in at least one dataset and HERV-K113. The
111  column colors show the human populations in the two databases. See Supplementary Figure 1 for the details of the
112  names of the indicated populations. (1kGP: 1,000 Genomes Project; HGDP: Human Genome Diversity Project; AFR:
113  African; AMR: American; EAS: East Asian; EUR: European; SAS: South Asian).
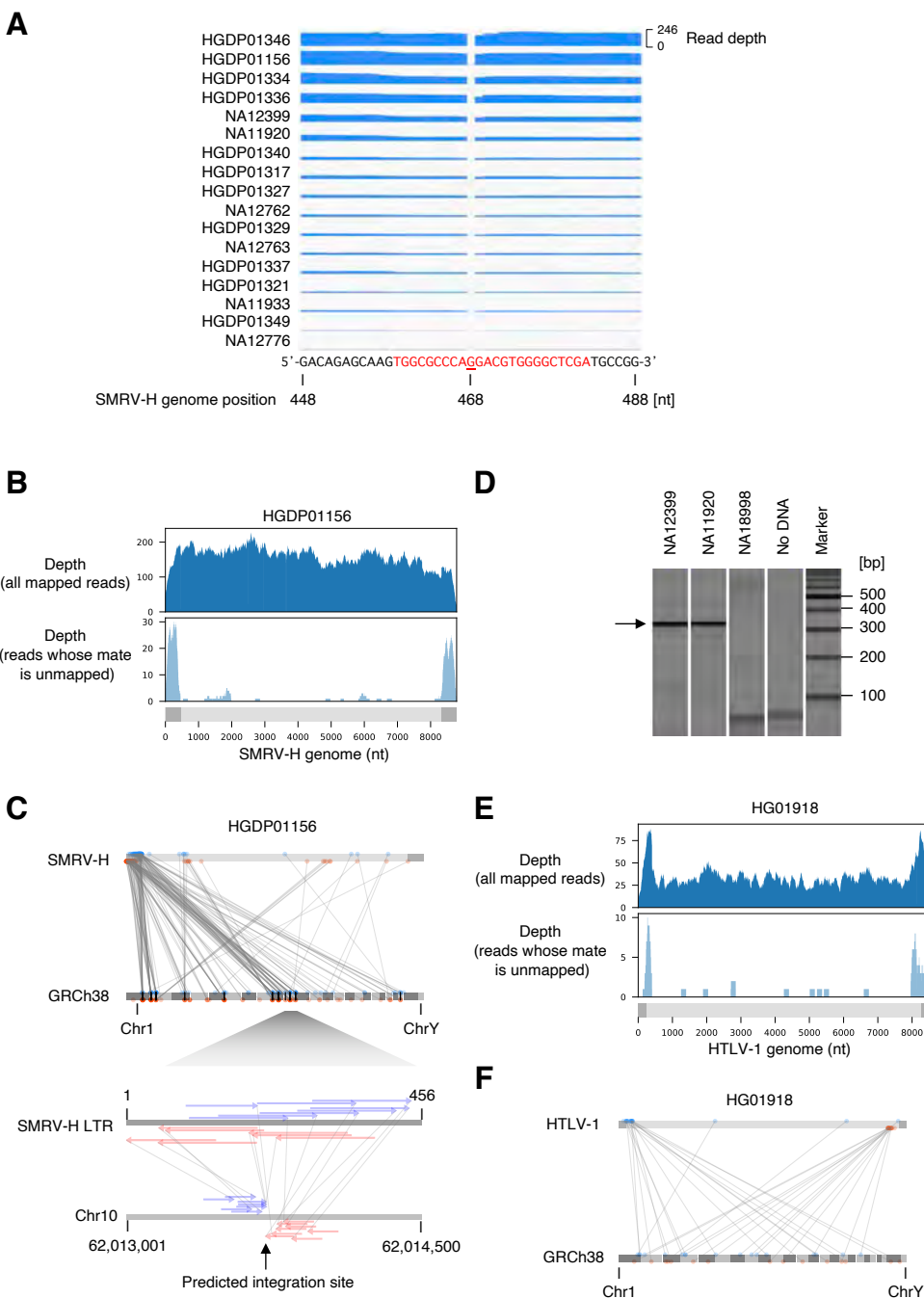114
115
116  *Squirrel monkey retrovirus*

117  SMRV-mapped reads were abundantly detected in 18 datasets, with a wide range of

118  SMRV-mapped read depths in datasets from different subjects (Figure 2A, Supplementary

119  Figure 3). From the 1kGP, all 12 datasets with SMRV were from subjects from Utah, including 2

120  subjects with read depths greater than 1x autosomal depth. All reads lack a guanosine in the

121  tRNA primer binding site (PBS) relative to SMRV-H, which was isolated from macaque cell-

122  derived preparations of Epstein Barr virus (EBV) (22); the tRNA PBS of the SMRV detected

123  here is identical to SMRV sequences recently obtained from Vero cells (23), a cell line used for

124  biologicals production (Figure 2A) (24). To determine if SMRV was integrated into the genomes

125  of the sequenced LCLs, we searched for paired reads with one read mapped to the virus

126  genome and one pair to the human genome (i.e. hybrid reads). Virus-human hybrid reads were

127  observed, often mapped to the SMRV long terminal repeat (LTR), consistent with the virus

128  being integrated into human chromosomes (Figure 2B); the human-mapped reads of these

129  hybrid pairs mapped to multiple chromosomal loci (Figure 2C, Supplementary Figure 4). We

130  found no enriched integration sites consistent with a clonal integration in the germline, nor

131  shared integration sites across different subjects. These observations suggest these SMRV

132  integrations are somatic rather than present in the germline. To assess whether the SMRV

133  integrations occurred before or after the peripheral blood mononuclear cells (PBMCs) used to

134  produce LCLs were removed from these subjects, we analyzed sequencing results obtained

135  directly from these subject's nucleated blood cells (25, 26). No reads mapped to SMRV. The

136  most plausible source of SMRV in the sequenced cell lines is thus laboratory contamination (22,

137  27). This observation notwithstanding, SMRV requires biosafety level 2 (BSL2) containment,

138  and the adventitious presence of SMRV in these samples could influence other results using

139  these reference materials. We obtained LCLs from 2 subjects from whom high SMRV reads

140  were found and confirmed the presence of SMRV DNA by PCR (Fig 2D). We analyzed existing

141  RNAseq data (28) and confirmed that SMRV is transcribed and is associated with differential

142  gene expression relative to LCLs in which SMRV RNA is not detected (Supplementary Figure

143  5).

144

6

**Figure 2** Chromosomal integrations of SMRV and HTLV-1

A.  Depth of WGS reads mapping to the primer binding site (PBS) of the SMRV-H genome. Seventeen datasets with at least one read mapping to PBS are shown. One dataset did not have any read mapping to PBS. The PBS of SMRV-H is shown with red characters. In all WGS datasets, the SMRV reads lack the guanosine present at the 468th nucleotide of the SMRV-H genome.

B.  Depth of HGDP01156 reads mapping to SMRV-H. Upper panel shows the depth of all reads in the dataset mapping to the SMRV-H genome. Lower panel shows the depth of reads mapping to the SMRV-H genome whose mate is not mapped to the SMRV-H genome. Virus genome structure is shown as gray bars. LTR are shown as dark gray rectangles.

C.  Mapping positions of SMRV-chromosome hybrid reads. Read-1 and Read-2 of a read-pair are connected with a line. All LTR-mapped reads are shown on the left LTR. The lower panel shows the predicted SMRV

7

157        integration site on chromosome 10. Gray bar in the top of the upper panel represents the virus genome
158             structure. Dark gray rectangles represent LTR. Reads mapping to the forward and reverse directions are
159             shown as blue and red arrows, respectively.

160    D.   Detection of SMRV DNA from 1kGP LCLs by PCR. Genomic DNA extracted from the indicated LCLs were
161             used as templates for PCR. WGS datasets from NA12399 and NA11920 are positive for SMRV, while that of
162             NA18998 is negative.

163    E.   Depth of HG01918 reads mapping to HTLV-1. Upper panel shows the depth of all reads in the dataset
164             mapping to HTLV-1. Lower panel shows the depth of reads whose pair is not mapped to the HTLV-1
165             genome.

166    F.   Mapping positions of HTLV-1-chromosome hybrid reads. Read-1 and Read-2 of a read-pair are connected
167             with a line. The reads mapping to left LTR was kept when a read was multi-mapped to both left and right
168             LTR. The genome position of reads mapping only to right LTR were replaced to the left LTR.

169

170

171   *Human immunodeficiency virus 1 and human T-lymphotropic virus 1*

172       We detected reads mapping to HIV-1, whose primary targets in the peripheral blood are

173   T lineage cells, in 8 datasets with a maximum coverage of the viral genome 8.5% and depth

174   0.29x, inconsistent with germline integration (Supplementary Figure 2). LCLs are generated by

175   infecting PBMCs with EBV, which infects mature B cells. Accordingly, most LCL B cell receptors

176   (BCR) have undergone V(D)J recombination, the signature of mature B cells. Moreover, the

177   mode of BCR clonality in a subset of the LCLs analyzed here is one; i.e. they are monoclonal

178   (29). Expression of a rearranged T cell receptor, consistent with presence of T lineage cells,

179   was observed in only one LCL among over 450 screened. HIV-1-mapped reads thus likely

180   either result from infection of hematopoietic progenitor cells (30), ongoing infection of LCLs (31),

181   or from contamination. We did not attempt to confirm the presence of infectious HIV-1 from

182   these cell lines.

183       Like HIV-1, HTLV-1 is a known human pathogen endemic in populations studied here.

184   HTLV-1 is often transmitted perinatally; analyzing WGS is thus an opportunity to distinguish

185   somatic vertical transmission from potential occult germline horizontal transfer of HTLV-1. Five

186   datasets showed HTLV-1 reads, with read depths ranging from 0.03x (a single paired-end read)

187   to 1.1x relative to autosomes. Two datasets contained HTLV-1-mapping reads at a depth

188   potentially consistent with heritable integrations, 0.55x and 1.1x respectively. Using the hybrid

189   reads approach described above, we demonstrated multiple integrations (Figure 2E, F,

190   Supplementary Figure 4), arguing against germline-inherited integration as the cause of the high

191   abundance of HTLV-1-mapped reads in these datasets. Like HIV-1, HTLV-1 is not well known to

192   infect and integrate into B cells, the source of most LCLs. Thus integration of HTLV-1 into

193   hematopoietic progenitor cells and maintenance of integration site diversity through the LCL

194   generation process, or ongoing replication of HTLV-1 in LCLs, may explain these findings.

195    Subjects whose LCLs sequencing datasets suggest presence of SMRV, HIV-1, and HTLV-1

196    are listed in Supplementary Table 1.

197

198    *Human herpesvirus 6*

199         Germline-integrated HHV-6 has been reported in some of the same datasets analyzed

200    here (32), however we recently described another form of integrated HHV-6 in which a single

201    HHV-6 direct repeat (DR) is present (termed "solo-DR"). The solo-DR form presumably reflects

202    recombination between the two DR regions present in the full-length integrated HHV-6 genome

203    leading to excision of the unique portion of the viral genome (12). Abundant HHV-6 reads were

204    present in 18 datasets (Table 1, Figure 3A, Supplementary Figure 6), suggesting that these

205    subjects likely have chromosomally-integrated copies of HHV-6. One of these samples

206    contained reads mapped only to the DR region of HHV-6B, characteristic of the solo-DR form of

207    integrated HHV-6.

208

209    **Table 1** Summary of integrated HHV-6 identified from 1kGP and HGDP

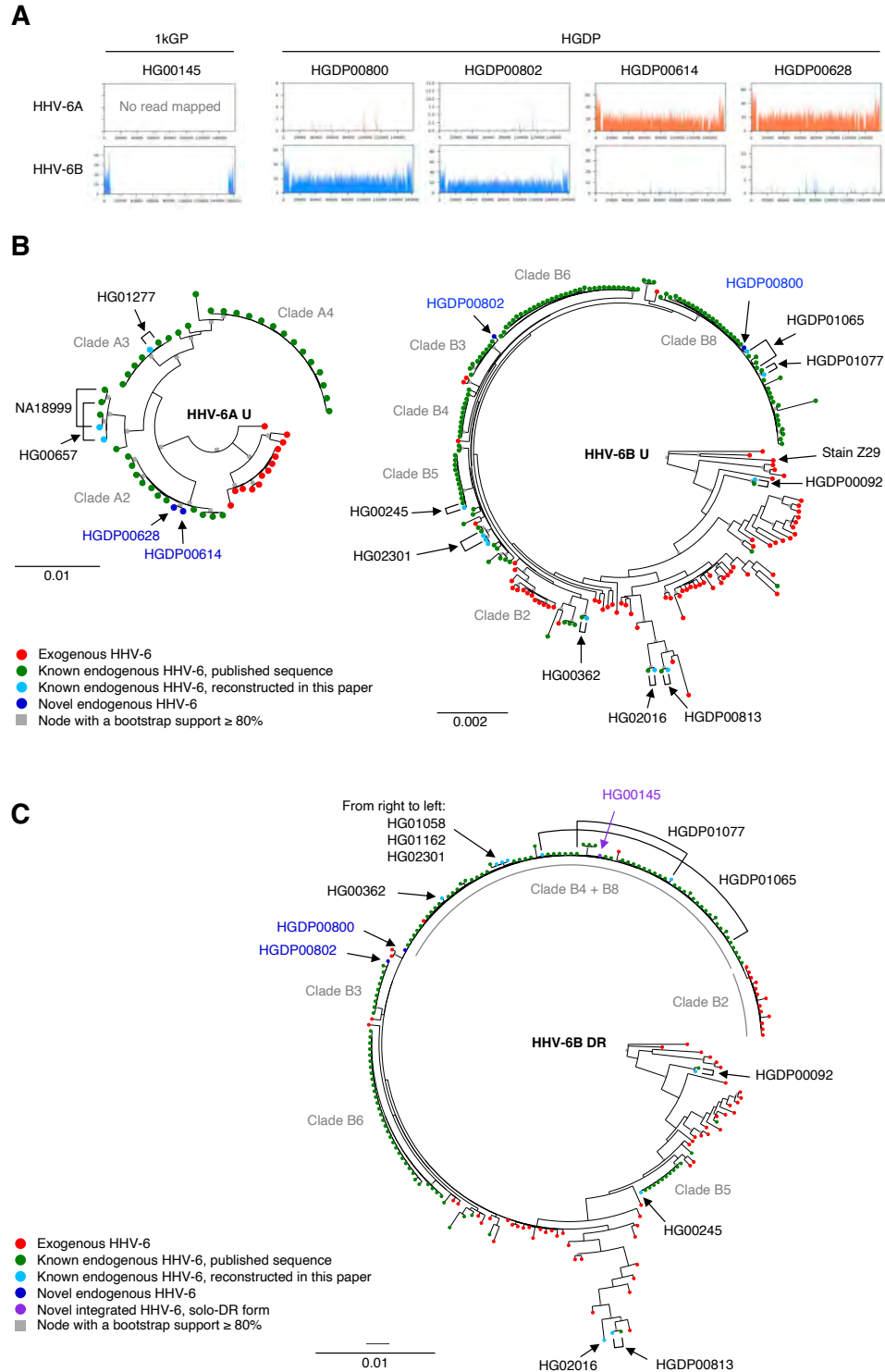| Database | Sample | HHV-6 | Structure | Population | Reference |
|---|---|---|---|---|---|
| 1kGP phase3 | HG00245 | HHV-6B | Full | GBR/EUR | Telford *et al*. |
| 1kGP phase3 | HG00362 | HHV-6B | Full | FIN/EUR | Telford *et al*. |
| 1kGP phase3 | HG01058 | HHV-6B | Full | PUR/AMR | Telford *et al*. |
| 1kGP phase3 | HG01162 | HHV-6B | Full | PUR/AMR | Telford *et al*. |
| 1kGP phase3 | HG02016 | HHV-6B | Full | KHV/EAS | Telford *et al*. |
| 1kGP phase3 | HG02301 | HHV-6B | Full | PEL/AMR | Telford *et al*. |
| 1kGP phase3 | HG00145 | HHV-6B | Solo-DR | GBR/EUR | This study |
| 1kGP phase3 | HG00657 | HHV-6A | Full | CHS/EAS | Telford *et al*. |
| 1kGP phase3 | HG01277 | HHV-6A | Full | CLM/AMR | Telford *et al*. |
| 1kGP phase3 | NA18999 | HHV-6A | Full | JPT/EAS | Zhang *et al*. |
| 1kGP pilot | NA19381 | HHV-6B | Full | LWK/AFR | Telford *et al*. |
| 1kGP pilot | NA19382 | HHV-6B | Full | LWK/AFR | Telford *et al*. |
| | | | | | |
| HGDP | HGDP00092 | HHV-6B | Full | Balochi/SAS | Zhang *et al*. |
| HGDP | HGDP00614 | HHV-6A | Full | Bedouin/MES | This study |
| HGDP | HGDP00628 | HHV-6A | Full | Bedouin/MES | This study |
| HGDP | HGDP00800 | HHV-6B | Full | Orcadian/EUR | This study |
| HGDP | HGDP00802 | HHV-6B | Full | Orcadian/EUR | This study |
| HGDP | HGDP00813 | HHV-6B | Full | Han/EAS | Zhang *et al*. |
| HGDP | HGDP01065 | HHV-6B | Full | Sardinian/EUR | Zhang *et al*. |
| HGDP | HGDP01077 | HHV-6B | Full | Sardinian/EUR | Zhang *et al*. |

210

211

212

213         To understand the origin of these integrated HHV-6 variants, we analyzed their

214    relationship to previously-reported exogenous and integrated HHV-6 sequences. All

215    reconstructed sequences clustered with previously-reported sequences (Figure 3B,

216    Supplementary Figure 8, 9). Two endogenous HHV-6A found in Bedouin subjects clustered with

217    clade A2 sequences, which were previously found in subjects in the US and UK, suggesting

218    these subjects share endogenous HHV-6A derived from a single integration event. Shared

9

219    ancestry of this chromosomal fragment, previously shown to correspond to the telomere of

220    chromosome 17p, between subjects on three continents is consistent with the deep evolutionary

221    relationship of endogenous HHV-6A (13). Newly-reported endogenous HHV-6B in two subjects

222    in HGDP were grouped with clade B8 and Clade B3, respectively. Clade B8 integrations also

223    map to chromosome 17p, which bears a short telomere (33). The integration site of clade B3

224    endogenous HHV-6 has not yet been determined.

225         To clarify the origin of the newly-detected solo-DR variant, we generated a phylogenetic

226    tree using only DR sequences (Figure 3C, Supplementary Figure 10). The solo-DR form

227    reconstructed from the HG00145 genome was present in the same clade as DRs from clade B4

228    and B8 full-length endogenous HHV-6B. This suggests that the solo-DR form likely arose from

229    an HHV-6B source closely related to that of clade B4 and B8, but precludes confident inference

230    that the solo-DR represents a germline rearrangement of a full-length endogenous HHV-6.

231    Detection of solo-DR integrated HHV-6 in this already well-characterized dataset shows that

232    screening WGS databases may provide additional information regarding the excision and

233    potential for reactivation of endogenous HHV-6.

234

235

**Figure 3** Detection and phylogenetic analysis of endogenous HHV-6.
- A. Depth of reads mapping to HHV-6 in the 5 WGS datasets from 1kGP and HGDP.
- B. Phylogenetic trees inferred from U regions of HHV-6A and B. The publicly available sequences of endogenous and exogenous HHV-6 as well as ones reconstructed in the present study were used.
- C. Phylogenetic tree inferred from DR regions of HHV-6B. The publicly available sequences of endogenous and exogenous HHV-6B, as well as ones reconstructed in the present study, were used. B, C. Clade names defined in the phylogenetic analysis in Aswad *et al*. are shown.

11

244

245

246     *Human endogenous retrovirus K*

247         Previous studies of the DNA virome detectable from human genome sequencing

248     datasets have noted inter-individual variation in reads mapped to HERV-K (18, 19). This

249     variation (e.g. Figure 1) has been speculated to result from polymorphic HERV-K integrations

250     (19). While the viruses described above are absent from reference human genomes, HERV-K is

251     present in multiple nearly-identical copies in reference genomes. This makes detecting

252     additional non-reference integrations and mapping the chromosomal location of polymorphisms

253     challenging. Previous advances in mapping HERV-K polymorphisms have been made by local

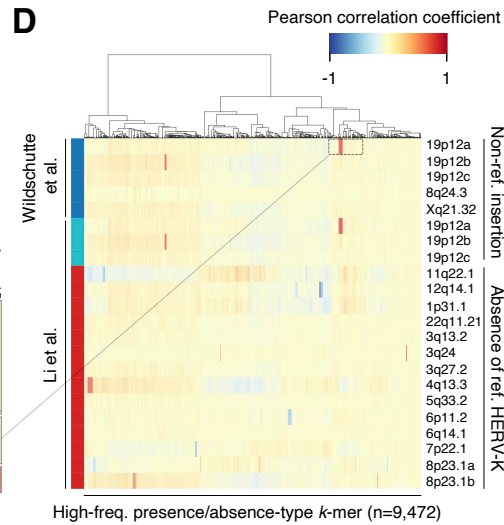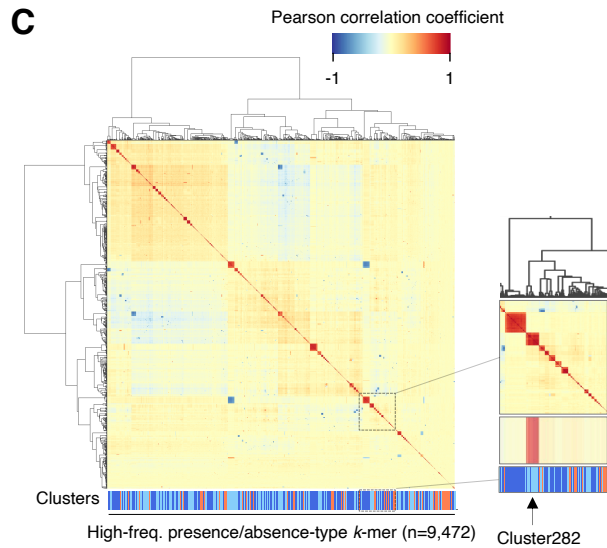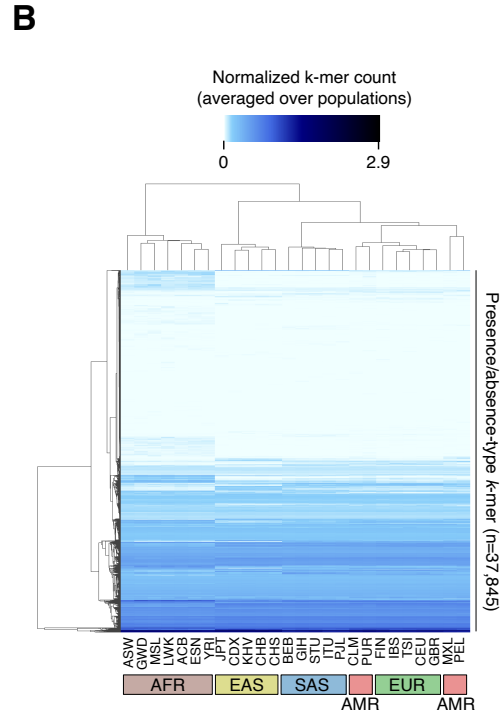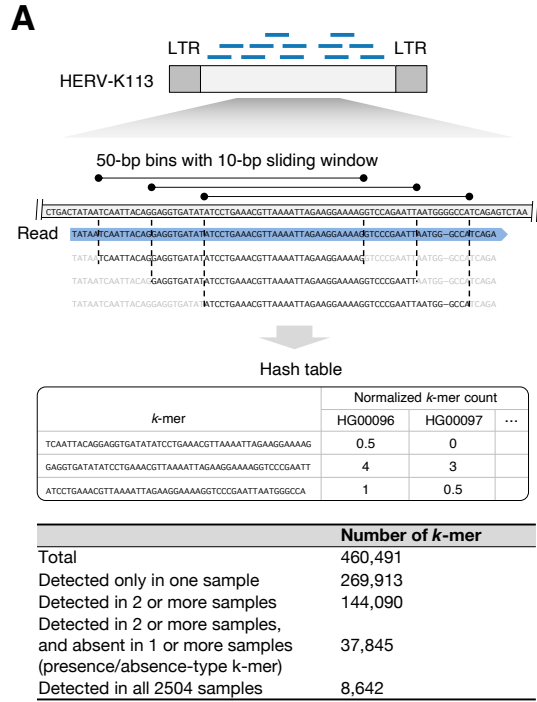254     breakpoint reconstruction using read-pairs that span insertion junctions (2). Taking a different

255     approach to this problem, we extracted all unique *k*-mers of length 50 from the aligned portion of

256     reads mapped to the HERV-K113 provirus (i.e. excluding the LTRs, Figure 4A). To filter only

257     those *k*-mers derived from HERV-K loci that are polymorphic between humans, we extracted *k*-

258     mers which were absent in at least one subject and present in at least two subjects (n=37,845

259     "presence-absence type" *k*-mers). Hierarchical clustering of presence/absence-type *k*-mer

260     occurrences recapitulates the continental human population supergroups (Figure 4B), as does

261     clustering based on the allele frequency of previously reported polymorphic HERV-K in human

262     subpopulations (2). This suggested that presence-absence *k*-mers may be a suitable proxy to

263     allow for discovery of additional polymorphic HERV-K alleles.

264         We hypothesized that structurally-polymorphic HERV-K alleles generate multiple unique

265     *k*-mers with the same pattern of presence or absence in multiple subjects. To test this, we

266     generated an all-by-all Pearson correlation coefficient matrix for presence/absence type *k*-mers

267     which were detected in more than 50 subjects (n=8,642) then performed hierarchical clustering

268     of *k*-mers. This revealed multiple groups of *k*-mers with presence/absence patterns that were

269     highly correlated with one other (Figure 4C), suggesting that a single polymorphic HERV-K

270     locus could generate multiple presence/absence-type *k*-mers. We formally defined clusters

271     using DBSCAN (see methods), resulting in 597 clusters of highly co-associated *k*-mers.

272         We next investigated how the observed clusters of presence/absence-type *k*-mers relate

273     to known HERV-K polymorphisms. Some *k*-mer clusters correspond very well with those of

274     known HERV-K polymorphisms previously described in these same subjects (Figure 4D). For

275     example, the presence/absence pattern of *k*-mers in cluster282 is highly correlated with that of a

276     non-reference HERV-K insertion on chr19. This suggests that in some cases, clusters of

277     presence/absence-type *k*-mers reflect HERV-K polymorphisms.

**A**

LTR        LTR

HERV-K113

50-bp bins with 10-bp sliding window

CTGACTATAATCAATTACAGGAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCAGAATTAATGGGGCCATCAGAGTCTAA

Read  TATAATCAATTACAGGAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCCGAATTRAATGG–GCCATCAGA

TATAATCAATTACAGGAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCCGAATT–AATGG–GCCATCAGA

TATAATCAATTACAGGAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCCGAATTAATGG–GCCATCAGA

TATAATCAATTACAGGAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCCGAATTAATGG–GCCATCAGA

Hash table

| k-mer | Normalized k-mer count | | |
| --- | --- | --- | --- |
| | HG00096 | HG00097 | ... |
| TCAATTACAGGAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAG | 0.5 | 0 | |
| GAGGTGATATATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCCGAATT | 4 | 3 | |
| ATCCTGAAACGTTAAAATTAGAAGGAAAAGGTCCCGAATTAATGGGCCA | 1 | 0.5 | |

| | Number of k-mer |
| --- | --- |
| Total | 460,491 |
| Detected only in one sample | 269,913 |
| Detected in 2 or more samples | 144,090 |
| Detected in 2 or more samples, and absent in 1 or more samples (presence/absence-type k-mer) | 37,845 |
| Detected in all 2504 samples | 8,642 |

**B**

Normalized k-mer count
(averaged over populations)

0        2.9

Presence/absence-type k-mer (n=37,845)

ASW GWD MSL LWK ACB ESN YRI JPT CDX KHV CHB CHS BEB GIH STU ITU PJL CLM PUR FIN IBS TSI CEU GBR MXL PEL

AFR    EAS    SAS    EUR
                AMR        AMR

**C**

Pearson correlation coefficient

−1        1

Clusters

High-freq. presence/absence-type k-mer (n=9,472)

Cluster282

**D**

Pearson correlation coefficient

−1        1

Wildschutte et al.

Li et al.

Non-ref. insertion
19p12a
19p12b
19p12c
8q24.3
Xq21.32
19p12a
19p12b
19p12c

Absence of ref. HERV-K
11q22.1
12q14.1
1p31.1
22q11.21
3q13.2
3q24
3q27.2
4q13.3
5q33.2
6p11.2
6q14.1
7p22.1
8p23.1a
8p23.1b

High-freq. presence/absence-type k-mer (n=9,472)

**E**

−log10(p-value)

Cluster282

count of a non-reference HERV-K insertion (19p12a)

1  2  3  4  5  6 7 8 9 10 11 12 13  15  17  19 22
Chromosome

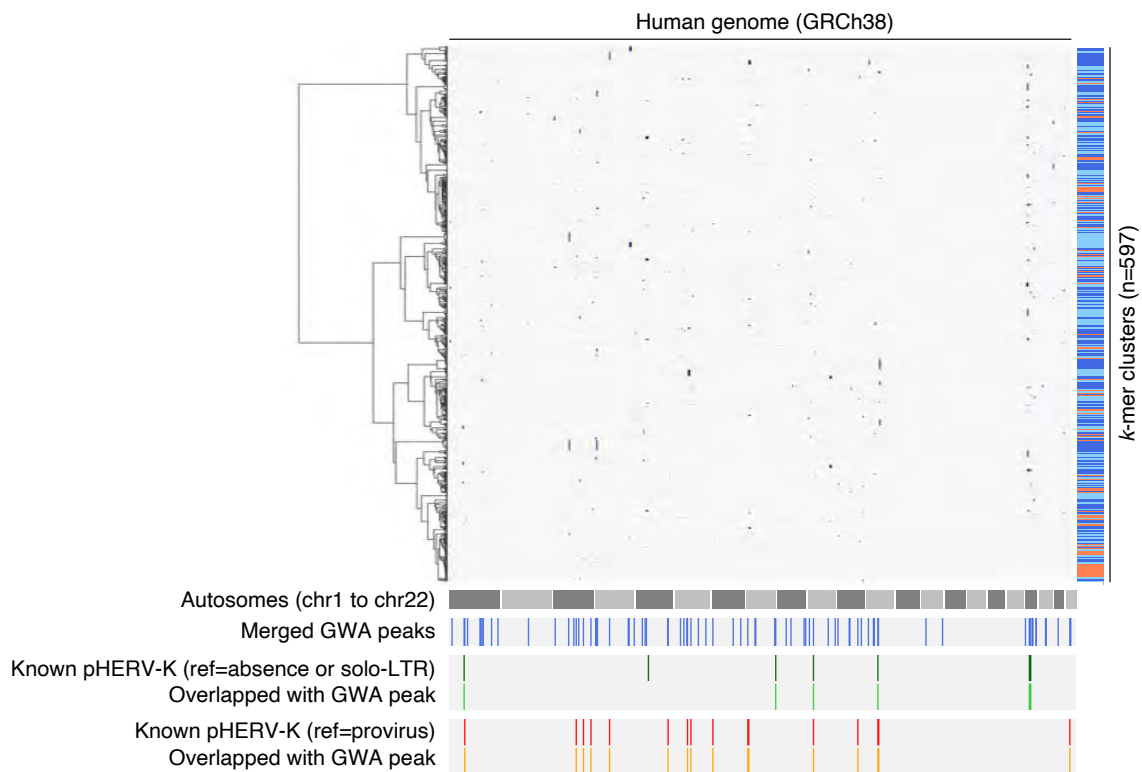2.14  2.15  2.16  2.17  2.18  2.19  2.20
Chromosome 19        (x10⁷)

279 **Figure 4** HERV-K *k*-mer detection from 1,000 Genomes Project WGS
280   A.  Schematic representation of *k*-mer counting from WGS reads mapping to HERV-K113. The HERV-K113
281       genome was split into 50-bp bins with a 10-bp sliding window, then, sequences of the mapped reads
282       corresponding to the HERV-K 50-bp bins were listed. The lower table shows the number of *k*-mers detected
283       from 2,504 WGS datasets from the 1,000 Genomes Project.
284   B.  Hierarchical clustering of *k*-mers based on their frequencies in 26 populations. Heatmap shows the normalized
285       *k*-mer count averaged over populations.
286   C.  Clustering of presence-absence type *k*-mers by Pearson correlation coefficient. Clustering were performed by
287       Ward's method (upper heatmap) and DBSCAN (lower color-bar). The heatmap shows the Pearson correlation
288       coefficient between *k*-mers, and the lower color-bar shows the clusters. Neighboring *k*-mer clusters are shown
289       as either dark or light blue. Orange represents the *k*-mers which were not clustered.
290   D.  Correlation between the occurrence of HERV-K *k*-mers and previously reported HERV-K polymorphisms.
291       Heatmap shows the Pearson correlation coefficient between the presence of *k*-mers and polymorphic HERV-K
292       reported in the two previous studies. (C, D) Insets in the panel C and D shows that the occurrence of *k*-mers in
293       cluster282 have high correlation to the presence of known polymorphic HERV-K in 19p12a.
294   E.  GWA using occurrence of *k*-mers detects known polymorphic HERV-K. Manhattan plots show SNVs with
295       association to the occurrence of *k*-mers in the cluster282. SNVs with p-value lower than 8.33e-11 are shown as
296       blue dots. Red solid line in the right panel shows the position of known non-reference HERV-K.

299         Determining where on the chromosome a particular polymorphic repetitive genetic

300   element is located can be challenging using short read sequencing data, because the read

301   evincing a polymorphism could potentially have arisen from a number of loci bearing nearly

302   identical elements. In addition, while some reads mapping to HERV-K LTRs are paired with

303   uniquely-mappable reads from flanking non-repeat DNA, variant reads mapping to the HERV-K

304   provirus are rarely paired with uniquely-mappable non-repetitive reads, as a consequence of the

305   size of most sequencing library inserts. To overcome these challenges, we took advantage of

306   the linkage disequilibrium (LD) structure of human chromosomes. We hypothesized that a *bona*

307   *fide* polymorphic HERV-K element, giving rise to a *k*-mer cluster, would be in linkage

308   disequilibrium with nearby SNVs. If so, analyzing genome-wide association (GWA) with SNVs

309   would allow us to locate the polymorphic HERV-K. To validate whether this approach is able to

310   accurately report the genome positions of polymorphic HERV-K, we examined a known non-

311   reference HERV-K insertion. GWA analysis using the presence/absence-pattern of cluster282 *k*-

312   mers, considered as a binary trait, detected a significant association with a single approximately

313   300-kb region on chromosome 19 known to contain the non-reference HERV-K insertion (Figure

314   4E). This approach of using l̲inkage d̲isequilibrium to f̲ind r̲epetitive e̲lement d̲ifferences, an

315   approach which may be applicable to other repetitive elements, is abbreviated here as "LDfred."

316         We performed GWA analyses using the presence/absence patterns of all 597 clusters.

317   As a result, 503 clusters detected at least one genome region with a Bonferroni-corrected

318   genome-wide significant association; clusters with associated regions tended to consist of more

319   *k*-mers and/or be present in more subjects (Supplementary Figure 11). We merged clusters that

14

320  were associated with overlapping regions (see methods), resulting in a total of 79 HERV-K *k*-

321  mer-associated loci spanning a total of 74.7 Mb. These loci most often include regions in which

322  the *P* values peak sharply, pinpointing the most tightly-linked LD block and narrowing the

323  presumed location of the HERV-K variant (e.g. Figure 4E). Consistent with previous work

324  showing that mobile elements are often linked to trait-associated SNVs (34, 35), the SNVs

325  comprising these HERV-K polymorphism-associated haplotypes are associated with numerous

326  human traits (Supplementaly Table 2), including 5 loci in which SNVs from the GWAS catalog

327  (36) overlap precisely with the genomic regions evincing polymorphic HERV-K (Supplemental

328  Figure 12). To check whether these new *k*-mer-associated loci indeed contain polymorphic

329  HERV-K, we evaluated overlap with known polymorphic HERV-K loci (2, 4–6). Seven out of 10

330  known non-reference proviruses are present within the observed loci (2 known non-reference

331  proviruses are not on the reference autosomes included for GWA), as are all 16 reference

332  HERV-K known to be absent in some subjects (Figure 5).



333

**Figure 5** Genome positions associated with HERV-K *k*-mers
The blue dots in the left clustermap show the genome positions with association to *k*-mer clusters by GWA analysis.
The blue lines in the third right column show the 79 genome loci associated with *k*-mer clusters. The dark green lines
show the 8 known non-reference HERV-K on autosomes. The light green lines show the 7 *k*-mer cluster-associating
genome loci overlapping with the known non-reference polymorphic HERV-K on autosomes. The dark orange lines
show the 16 known reference polymorphic HERV-K on autosomes. The light orange lines show the 16 *k*-mer cluster-
associating genome loci overlapping with the known reference absent HERV-K on autosomes.
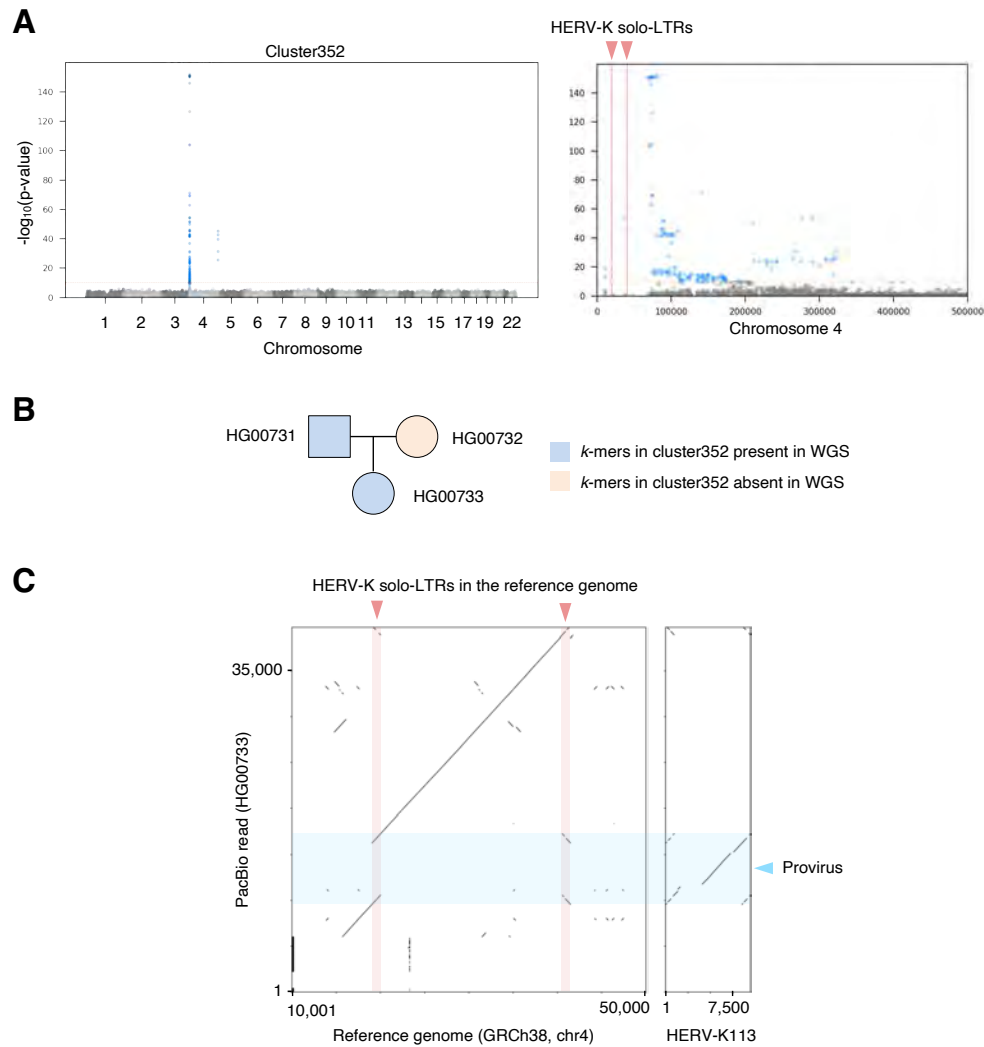
341

15

342        To assess whether LDfred could localize previously unknown polymorphisms, we

343    checked long-read sequencing datasets for reference HERV-K that are absent in some

344    subjects. We filtered an extensive catalog of SVs in three subjects, generated using multiple

345    sequencing technologies, for deletions that overlap with HERV-K (37). We found 24 reference

346    HERV-K elements with evidence of full or partial absence in at least one of the three subjects

347    (Supplementary Table 3). These 24 SVs spanned from 72 to 9,468-bp. We checked if these

348    HERV-K SVs were present in loci associated with $k$-mer clusters, and whether the

349    presence/absence pattern of the $k$-mers in each cluster was consistent with the

350    presence/absence of SV in the three subjects (see methods). Of the 24 HERV-K SVs, 9 were

351    concordantly detected (Supplementary Table 4); concordantly detected SVs tended to be longer

352    than those not concordantly detected. Notably, 4 out of 9 detected HERV-K SVs have not been

353    previously reported as polymorphic HERV-K. These 4 unreported SVs are attributable to

354    recombination between LTRs (Supplementary Figure 13), which is particularly difficult to find

355    using existing algorithms and short-read sequencing. This demonstrated that LDfred can

356    localize unknown HERV-K provirus polymorphisms, including provirus/solo-LTR polymorphisms.

357        SVs in complex or duplicated genome regions are also difficult to identify using short-

358    read data and available methods (38). To check the utility of our approach for this purpose, we

359    focused on HERV-K loci at chromosome ends, known to be complex genome regions (39). One

360    locus, associated with cluster352, is in the subtelomeric region of the short arm of chromosome

361    4. There are no HERV-K proviruses in this region, however there are 2 solo-LTRs, suggesting

362    the possibility of a provirus/solo-LTR polymorphism, or an additional non-reference HERV-K

363    insertion. We assessed whether either of these reference LTR loci sometimes contain a provirus

364    using long read sequencing data from a trio (40). WGS from the father and child contained

365    cluster352 $k$-mers, but the mother did not harbor any $k$-mers in cluster 352, suggesting that the

366    father and child could carry a non-reference provirus. We inspected reads mapped to the

367    subtelomere of chromosome 4 and found a read containing non-reference provirus at the region

368    corresponding to one of the solo-LTRs in the reference genome at this locus (Figure 6). Thus

369    LDfred can detect provirus/solo-LTR polymorphic HERV-K loci in complex regions including

370    subtelomeres.

371

**A**

Cluster352

HERV-K solo-LTRs

**B**

HG00731 — HG00732
HG00733

k-mers in cluster352 present in WGS
k-mers in cluster352 absent in WGS
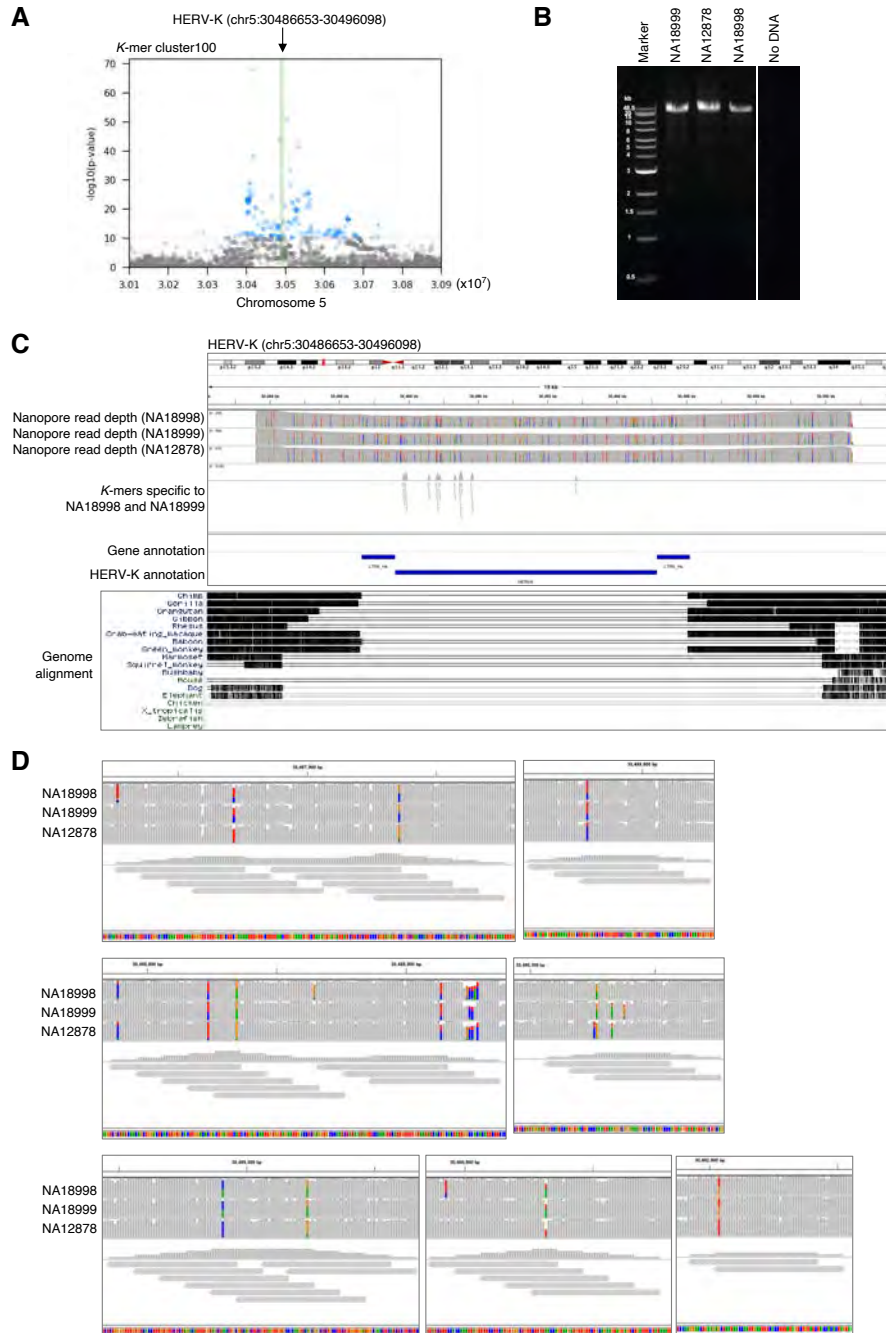
**C**

HERV-K solo-LTRs in the reference genome

Provirus

**Figure 6** *K*-mer-based method detects previously-unreported HERV-K polymorphism in subtelomere

A. SNVs associating with *k*-mers in cluster352. The right panel shows the region near the end of the p-arm of the chromosome 4. SNVs with p-value lower than 8.33e-11 are shown as blue dots. The red solid lines in the right panel show two reference HERV-K solo-LTR in the subtelomere region.

B. Presence and absence of *k*-mers of cluster352 in the public high-coverage short-read WGS of the Chinese trio.

C. HG00733 contains non-reference provirus. Left panel shows the dot matrix between the reference human genome and a PacBio read of HG00733. The right panel shows the dot matrix between the reference HERV-K113 and the PacBio reads of HG00733. Light blue and light red rectangles represent HERV-K provirus and solo-LTR, respectively.

Of the 79 loci identified here, five loci contain known non-reference presence/absence or provirus/solo-LTR polymorphisms and 14 contain reference HERV-K known to be absent in some subjects (including some merged loci containing more than one polymorphism). Six loci reflect provirus/solo-LTR polymorphisms that have not been previously reported, but which are readily demonstrated in long-read data (Figure 6, Supplementary Figure 13, 14). The remaining loci cannot be assessed using available long-read data, because the minor variants, as

17

390    determined by *k*-mer pattern, are not present in subjects with available data. To determine the

391    polymorphisms giving rise to the signal that allowed us to identify the remaining loci, we chose

392    three for which the regions flanking reference HERV-K in these regions allowed us to design

393    specific primers. Targeted long-read sequencing revealed differences in the HERV-K at these

394    loci consistent with the *k*-mer pattern differences between the individuals. The nature of the

395    variation at these loci was not structural; instead, it consisted of multiple SNVs (up to 11) linked

396    in a haplotype (Figure 7; Supplementary Figure 15, 16). This high degree of linked SNV

397    variation distinguishing HERV-K proviruses at the same locus is unexpected due to accrual of

398    substitutions; 11 mutations across these 4.9 kilobases of the HERV-K provirus (Figure 7) would

399    accumulate over 4.4 million years (41), and the linkage between them would be expected to be

400    degraded by crossover recombination events during that period. Instead, this more likely reflects

401    interlocus gene conversion via recombination, which has previously been described for HERV

402    on the basis of comparison of LTRs in different species (42), or introgression. Notably, the

403    specific HERV-K haplotypes present as minor variants in these loci are not detected by BLASTn

404    search of the hg38 reference genome. Thus LDfred can localize previously-uncharacterized

405    sources of non-structural HERV-K variation.

406

**Figure 7** Potential interlocus gene conversion in HERV-K localized by LDfred

A. Manhattan plot showing SNVs associating the *k*-mer cluster100. SNVs with p-value lower than 8.33e-11 are shown as blue dots. Green line shows the reference HERV-K provirus.

B. Amplification of the HERV-K provirus by PCR. HERV-K provirus with adjacent sequence was amplified and PCR products were separated by gel electrophoresis. DNA extracted from LCLs originating from NA18998, NA18999, and NA12878 were used as templates.

C. Upper panel: IGV view of long-read sequencing reads mapping to HERV-K. The PCR amplicons were sequenced using an Oxford Nanopore flongle flow cell and mapped to GRCh38. *k*-mers in *k*-mer detecting the HERV-K were also mapped to the PCR target regions. Lower panel: UCSC genome browser view showing the Multiz Alignment of 100 Vertebrates track.

D. Enlarged images of panel C. NA12878 carries two alleles of a non-reference HERV-K haplotype (which is not observed elsewhere in the reference genome) also present as a single allele in NA18998 and NA18999.

19

**Discussion**

420

421        This work provides a comprehensive picture of virus-derived structural variations in two

422    well-studied global WGS datasets. We found previously-missed germline structural variants

423    arising from HHV-6 and HERV-K, as well as virus integration in somatic cells due to natural

424    infection or contamination. The presence of SMRV integrations in LCLs introduces caveats in

425    analyzing these materials and the data derived from them. This is especially notable for cells

426    used in the 1kGP, from which only subjects from Utah (collected by CEPH) are SMRV-positive.

427    In most subjects, the number of virus-chromosome hybrid reads detected at each specific

428    genome locus is low, suggesting a mosaic of cells with different integration events. However, in

429    some datasets, such as HGDP01156 and HGDP01346, a substantial number of hybrid reads

430    are detected from individual genome loci, suggesting a high fraction of clonal cells bearing the

431    same virus integration event. In such cases, the SMRV insertion could influence nearby variant

432    calls, and could also influence clonal expansion during the course of LCL culture (43). SMRV is

433    transcribed in SMRV-positive LCLs, and this viral RNA could influence host transcription, for

434    example by triggering innate immune pattern recognition receptors. However, we observed no

435    significant change in expression of interferon-stimulated genes (ISGs). This may be related to

436    the concurrent presence of EBV, reported to counteract ISGs, in these cells (44). SMRV-

437    positive LCLs did express a few genes significantly differently than SMRV-negative cells, which

438    has implications for interpreting the results of studies using these cells and datasets (28). We

439    confirmed the presence of SMRV DNA in recently-distributed LCLs from these donors. While

440    biosafety regulations vary by locality, our results reinforce that even well-characterized LCLs

441    should be handled as potential sources of infectious viruses.

442        We found evidence of infection of LCL progenitors with HIV-1 or HLTV-1. This was

443    unexpected because B cells, the proximal progenitor of most LCLs, are not efficiently infected

444    by either of these viruses. However EBV transformation has recently been shown to permit

445    replication of some types of HIV-1 in B cells (31). Given the diversity of integration sites

446    observed, we suspect that a similar phenomenon may explain the presence of HTLV-1 in LCLs,

447    although we cannot exclude somatic mosaicism due to infection of hematopoetic stem cells

448    (45). We found no evidence of germline structural variants related to HIV-1 or HTLV-1. As has

449    recently been noted, HIV-1 is capable of infecting germ cells (46). Our result should be

450    interpreted to indicate that the plausible upper bounds of the global allele frequency of such

451    variants is ~0.01%; larger-scale projects, or projects sequencing populations with higher

452    prevalence of infection by these viruses, could apply the methods used here to discover these

453    rare variants, if they exist.

454    An endogenous form of Koala retrovirus has been considered a unique opportunity to
455    study a virus that is in the process of endogenization in mammals (47). Our report clarifies the
456    extent to which humans also harbor a virus that straddles the endogenous/exogenous divide,
457    HHV-6. Understanding the relationship between endogenous and exogenous HHV-6 is a critical
458    question to which the tools presented here can be usefully applied. The solo-DR form of
459    integrated HHV-6B is present in one 1kGP dataset, but was not detected in previous surveys of
460    these data and samples (32). The presence of unreported solo-DR integrated HHV-6 in such a
461    well-investigated database suggests that the prevalence and diversity of integrated HHV-6 has
462    likely been underestimated in other studies as well. Phylogenetic analysis suggests that this
463    solo-DR variant was potentially formed by partial excision of B4 or B8 clade endogenous HHV-
464    6B by recombination, leaving behind one DR as a "scar." This molecular event has been
465    proposed to lead to viral reactivation. In that context, it is notable that a few reportedly-
466    exogenous HHV-6B sequences are quite closely related to B4 and B8 endogenous HHV-6B,
467    while the majority of sampled exogenous HHV-6B are more divergent. We cannot exclude the
468    possibility that this solo-DR variant is related to a third independent integration by a virus related
469    to those giving rise to B4 and B8 endogenous HHV-6B. Furthermore, while the DR was
470    evidently present in a hemizygous state in the nuclei of LCLs from this subject, there is no
471    evidence that it is endogenous or "inherited" as is often used to describe such variants; we
472    consider it most accurate to describe it as chromosomally-integrated HHV-6 until another
473    identical-by-descent variant should happen to be observed, integrated in another human's
474    genome.
475    We recently reported the presence of the solo-DR form of endogenous HHV-6 in the
476    Japanese population. The solo-DR positive subject in 1kGP is of European ancestry,
477    demonstrating that solo-DR variants are present in populations of different ancestry. Integrated
478    HHV-6 has been associated with angina pectoris and pre-eclampsia, and in both cases a virus-
479    dependent mechanism has been postulated (48, 49). It is thus important to ask whether the
480    solo-DR form of endogenous HHV-6, or only full length HHV-6 integration, is associated with
481    human phenotypes and diseases. Screening additional databases using the tools developed
482    here can capture the complete diversity of HHV-6 integration in human chromosomes, leading
483    to a deeper understanding of its potential influence on human diseases. We also found four
484    previously unreported full-length endogenous HHV-6. Among these are independent
485    integrations, one HHV-6A and one HHV-6B, into chromosome 17p, which is reported to carry
486    the shortest telomere of the 46 chromosome arms (33). This mirrors our observation in the
487    Japanese population, in which two prevalent endogenous HHV-6 variants, one HHV-6A and one

488 HHV-6B, are integrated on chromosome 22q, which carries the second shortest telomere, on

489 average (33). One full-length endogenous HHV-6B falls within Clade B3, the integration site of

490 which is currently unknown, but will be able to be mapped using this LCL in the future.

491 We also explored the polymorphisms of HERV-K, which includes the most recently

492 integrated HERVs. The overlap of GWA peaks with genome loci known to harbor polymorphic

493 HERV-K suggests that our approach captures HERV-K polymorphisms that can be discovered

494 by other methods. In addition, our approach identifies regions with association to HERV-K $k$-

495 mers that were not previously reported to be polymorphic; these loci likely contain unreported

496 HERV-K polymorphisms. Using long-read sequencing data, we confirmed that six of these loci

497 indeed contain HERV-K structural variants. The $k$-mer-based method presented here does not

498 explicitly distinguish presence/absence-type polymorphisms from SNV polymorphisms in the

499 non-LTR portion of HERV-K. $K$-mer signals due to individual substitutions, as would begin to

500 accumulate after integration of a full length HERV-K provirus at a given locus, may also be

501 detected. This approach conceptually allows finding any sequence differences between repeat

502 element loci, and may be useful for other difficult-to-map repeat elements (38, 50). In total we

503 document the nature of the HERV-K polymorphisms explaining over 20 of the loci reported here,

504 yet many loci remain to be validated; increasing use of long-read sequencing should enable this

505 soon (51).

506 We set a threshold of clustering only presence/absence type $k$-mers found in at least 50

507 subjects, so HERV-K polymorphisms with allele frequencies below 0.02 should not be detected.

508 While we set this threshold arbitrarily and conservatively, this approach is limited in its ability to

509 localize very rare polymorphisms due to the nature of linkage disequilibrium and the decreased

510 statistical power of association testing using few polymorphism-bearing subjects. We treated $k$-

511 mer presence-absence pattern as a binary trait, yet retaining the continuous variation in these

512 patterns could approximate genotyping, potentially improving localization of polymorphisms in

513 repeats in the future. Previous studies reported that the majority of unfixed HERV-K in humans

514 are solo-LTR type (2). We defined $k$-mers using only reads mapped to the non-LTR regions of

515 HERV-K due to the high sequence similarity between HERV-K LTR and SVA retrotransposons.

516 This enabled us to detect new solo-LTR vs full provirus polymorphisms, but we could not detect

517 solo-LTR vs empty site polymorphisms. To understand polymorphisms of HERV-K

518 comprehensively, including presence/absence of solo-LTR, we will need to expand the $k$-mer-

519 based method; this will cross-detect SVA retrotransposons which are not virus-derived in this

520 same sense as the other variants considered within the scope of this report.

521        We used targeted long read sequencing to determine the HERV-K polymorphisms

522    present within several of the newly-identified loci. We observed divergent HERV-K haplotypes,

523    differing by 11 linked SNVs within 4,932 bp from the major allele in the most divergent, present

524    at the same locus. This degree of variation at a syntenic HERV-K integration site, absent in

525    other great ape genomes, is unexpected as a result of clock-like accumulation of mutations (52).

526    Among the potential explanations for this phenomenon, two are plausible and warrant

527    discussion. First, a process of non-allelic homologous recombination, most often referred to as

528    gene conversion in this context, could exchange the HERV-K haplotype at the locus *en bloc*

529    with that from another locus. This has often been invoked to explain differences between HERV

530    loci syntenic between species (42), and has been reported for HERV-K (53). However, the

531    potential "source" haplotype for such a recombination could not be identified in the hg38

532    reference genome (i.e. by BLASTn search using the variant *k*-mers). We thus cannot distinguish

533    whether the source HERV-K element was itself an insertion that has now been lost or fixed as a

534    solo-LTR, nor can we exclude the possibility of introgression of the HERV-K haplotypes from an

535    archaic source. In any case these results point to a previously unexplored cache of HERV-K

536    diversity in human genomes, and offer a new tool to guide its exploration. Considering that

537    infectious HERV-K sequences can be generated with via recombination between known HERV-

538    K elements (54), these previously hidden HERV-K polymorphisms are particularly relevant to

539    study in relation to human phenotypes. Ongoing and future large-scale population sequencing

540    projects will massively expand the data available to address viral contributions to human

541    genomes, and the tools presented here will enable integration of these analyses into the

542    planned output of these consortia (55).

543

544

545

546 **Methods**
547
548 *WGS datasets*
549 High-coverage WGS datasets from 1kGP were downloaded from the following URL:
550 'ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/'. High-
551 coverage WGS datasets from 1kGP Han Chinese trio were downloaded from the following URL:
552 ' https://www.internationalgenome.org/data-portal/data-collection/structural-variation'. High-
553 coverage WGS dataset from HGDP were downloaded from the following URL:
554 'ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/HGDP/data/'. The utilization of the
555 high-coverage WGS of multigenerational CEPH/Utah families (phs001872) are authorized by
556 the National Human Genome Research Institute through dbGaP for the following project: "The
557 prevalence, evolution, and health effects of polymorphic endogenous viral elements in human
558 populations."
559
560 *Preparation of reference virus genomes*
561     Reference virus genomes were downloaded from NCBI on April-6-2020. We
562 downloaded three files named
563 'ftp://ftp.ncbi.nlm.nih.gov/refseq/release/viral/viral.[1,2,3].1.genomic.fna.gz'. These three files
564 contained 12,182 virus genomes, including phages. These files were concatenated into one file
565 and used as reference virus genomes for further analysis.
566
567 *Virus detection and reconstruction from WGS*
568     Reads that did not map to the reference human genome were extracted from WGS
569 BAM or CRAM files using `samtools view -f 1 -F 3842 | samtools view -f 12 -F 3328 -`
570 command. Then, the unmapped reads were converted to FASTQ format using samtools fastq
571 command. FASTQ reads shorter than 20-nt were removed using a custom Python script and
572 excluded from downstream analysis. Retrieved sequences were then mapped to the reference
573 virus genomes using Hisat2 with `--mp 2,1 --no-spliced-alignment` options. After mapping, reads
574 estimated to be PCR duplicates were marked using picard MarkDuplicates command. Then, the
575 mapping depth and coverage of each virus was calculated using deeptools bamCoverage
576 command with `--binSize 1` option. Based on the depth and coverage, we searched for viruses
577 with abundant reads using a custom Python script. We labeled 'virus_exist' if 5% of a viral
578 genome was covered at more than 2x read depth. Virus genomes with a 'virus_exist' label were
579 then reconstructed by incorporating variations to the reference virus genomes. To reconstruct
580 viruses, variations in the reads mapping to virus genomes were called using gatk
581 HaplotypeCaller. The output vcf files were then normalized using bcftool norm command and
582 the reconstructed virus sequences were generated using bcftools consensus command.
583 Regions without any mapped reads were masked by 'N' using a custom Python script. The
584 workflow described here was compiled as a Python pipeline and available from the following
585 GitHub repository: 'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
586     We detected reads mapping to 634 viruses, including 553 phages, in total (Supplemental
587 figure 1 (the original heat map, showing all viruses by all people)). Phage are ubiquitous in the
588 human virome and thus should not necessarily be excluded as a potential source of horizontal
589 gene flow to humans (56). Phage-mapped reads were often found, however they were present

590     at low depth inconsistent with germline integration into human genomes, and were thus
591     excluded from further analysis. The same was true of most eukaryotic viruses, which showed
592     low average read depth, usually less than 1x across the entire length of the viral genome. This
593     may reflect virus infection in a small proportion of cells, contamination from other samples
594     sequenced on the same machine, or mis-mapping. As the primary goal of this study was to
595     detect potentially heritable virus-derived structural variants, these were not analyzed further.
596
597     *Endogenous HHV-6 detection and reconstruction from WGS*
598          We developed a bioinformatic pipeline specialized to detect and reconstruct full-length
599     endogenous HHV-6 as well as solo-DR form, because endogenous HHV-6 has a terminal direct
600     repeat sequence (DR), which is not appropriately reconstructed using the virus detection and
601     reconstruction pipeline described above. We extracted reads that did not map to the human
602     genome and mapped these reads to the reference HHV-6 using the same commands described
603     above. Rather than all viral sequences for HHV-6 reconstruction, we used full-length exogenous
604     HHV-6 genomes NC_001664.4 and NC_000898.1 as HHV-6A and HHV-6B, respectively. We
605     judged whether a WGS dataset contains abundant HHV-6 reads using the same cutoff
606     described above. When abundant HHV-6 reads were detected, we reconstructed the full-length
607     HHV-6 sequence using the same reconstruction protocol as described above. The DR region of
608     a reconstructed full-length genome is not accurate, because reads mapping to DR are mapped
609     to both left DR and right DR. The reads with multimapping have the mapping score 0, being
610     excluded from downstream variant calling. To accurately reconstruct DR, all reads mapping to
611     HHV-6 genomes were re-mapped to DR-only. For this reconstruction, we used nucleotides 1-
612     1089 of NC_001664.4 and 1-8793 of NC_000898.1 as DR-only sequences of HHV-6A and
613     HHV-6B, respectively. The workflow described here was compiled as a Python pipeline and
614     available from the following GitHub repository:
615     'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
616          For validation of the accuracy of endogenous HHV-6 reconstruction, we used a 35x
617     coverage WGS dataset from subject NA18999, a Japanese subject known to contain full-length
618     endogenous HHV-6A (Supplementary Figure 7). The reconstructed sequence covered 96% of
619     the reference genome (U1102) with 96.7% similarity. Phylogenetic analysis of the reconstructed
620     sequence with sequences from this subject previously determined by Sanger sequencing
621     demonstrate that the reconstructed sequence is very close to that determined by Sanger
622     sequencing. To understand the influence of WGS depth on the accuracy of endogenous HHV-6
623     reconstruction, we downsampled the dataset to approximately 30x, 20x, 15x, 10x, and 5x
624     autosome depths using the `picard DownsampleSam` function. Our pipeline detected
625     endogenous HHV-6 at all read depths, and had accuracy near that of Sanger sequencing when
626     the depth was higher than 15x. This demonstrates that, from moderate- to high-depth WGS
627     datasets, our pipeline can reconstruct relatively accurate endogenous HHV-6 sequences
628     suitable for phylogenetic analysis.
629
630     *Visualization of virus-chromosome hybrid reads*
631          WGS reads that failed to map to the reference human genome were mapped to viruses
632     using the pipeline described above. To detect virus-chromosome hybrid reads, read pairs with
633     one read mapped to a virus and the paired read not mapped to the same virus were retrieved

634  using `samtools view -f 8` command. Then, the unmapped mate reads were mapped to the
635  reference human genome GRCh38DH using `blastn -evalue 1e-15 -culling_limit 2 -
636  qcov_hsp_perc 90 -perc_identity 95 -word_size 11` command. Then, reads uniquely mapped to
637  the human genome were retrieved and the mapped positions of the hybrid reads on the virus
638  genomes and the human genome were visualized using a custom Python script. Because both
639  SMRV and HTLV-1 have LTRs, reads mapped to 3'LTR were re-mapped to 5' LTR. The script
640  used for visualization is available from the following GitHub repository:
641  'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
642
643  *Phylogenetic analysis of endogenous HHV-6*
644          To reconstruct phylogenetic trees of U regions, we used full-length genomes
645  reconstructed by the endogenous HHV-6 reconstruction pipeline described above. The
646  reconstructed sequences were aligned with known endogenous and exogenous HHV-6 using
647  `mafft --auto` command. To exclude the regions thought to have low reconstruction accuracy,
648  we removed DR and repeat sequences annotated in NCBI from the alignment. We removed
649  regions corresponding to nucleotides 0-8089, 127548-128233, 131076-131854, 140075-
650  140951, and 151288-159378 of HHV-6A NC_001664.4 and 0-8793, 9314-9510, 129045-
651  129681, 133500-133863, 133981-134076, 140081-142691, and 153321-162114 of HHV-6B
652  NC_000898.1 (all start positions here are 0-based numbering and end positions are 1-based
653  numbering) using a custom Python script. Then, phylogenetic trees were inferred by the
654  maximum likelihood method with the complete deletion option using MEGA X software. The
655  Kimura 2-parameter model was used. The reliability of each internal branch was assessed by
656  100 bootstrap resamplings. The phylogenetic trees were visualized using ETEtoolkit.
657          To reconstruct phylogenetic trees of DR regions, we used DR reconstructed by the
658  endogenous HHV-6 reconstruction pipeline described above. The reconstructed sequences are
659  aligned with known endogenous and exogenous HHV-6 using `mafft --auto` command. To
660  exclude the regions thought to have low reconstruction accuracy, we removed simple repeats
661  and low complexity sequences from the alignment. To define simple repeats and low complexity
662  sequences, we used RepeatMasker. We masked the reference HHV-6 (NC_001664.4 and
663  NC_000898.1) with a `repeatmasker -s -no_is` command. We removed the regions
664  corresponding to nucleotides 0-376, 1682-1730, 2302-2367, 2369-2451, 2692-2733, 3149-
665  3181, 3433-3502, 3626-3670, 7483-7519, 7655-8089 of HHV-6A NC_001664.4 and 0-393,
666  1926-2011, 2674-2717, 3013-3067, 3670-3713, 3959-3988, 8248-8793 of HHV-6B
667  NC_000898.1 (all start positions here are 0-based numbering and end positions are 1-based
668  numbering) using a custom Python script. Then, phylogenetic trees were inferred and visualized
669  as described above. The scripts used for phylogenetic analysis and the newick files of the
670  phylogenetic trees are available from the following GitHub repository:
671  'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
672
673  *Processing of RNA-sequencing datasets*
674          The SRA files of the Geuvadis RNA-seq dataset were downloaded from NCBI using
675  the `prefetch` command in the NCBI SRA-tools. We used 159 datasets derived from LCL of
676  Utah residents. The downloaded SRA files were converted to FASTQ files using `fasterq-dump`
677  command in the NCBI SRA-tools with `-S` option. Paired-reads were then filtered using fastp

678   software with `-l 20 -3 -W 4 -M 20 -t 1 -T 1 -x` options. The filtered paired-reads were then
679   mapped to the human genome 'GRCh38.p13.genome.fa' downloaded from GenCode. For
680   mapping, we used STAR software with `--quantMode GeneCounts --twopassMode Basic --
681   outFilterType BySJout --outFilterMultimapNmax 20 --alignSJoverhangMin 8 --
682   alignSJDBoverhangMin 1 --outFilterMismatchNmax 999 --outFilterMismatchNoverReadLmax
683   0.04 --alignIntronMin 20 --alignIntronMax 1000000 --alignMatesGapMax 1000000` options. We
684   provided the human gene annotation 'gencode.v33.annotation.gtf' downloaded from GenCode
685   when indexing the reference human genome using STAR.
686
687   *Differential gene analysis*
688           The 159 RNA-seq datasets were derived from 90 LCLs. Because some LCLs were
689   represented by two different RNA-seq datasets, we merged the count tables originating from
690   LCLs from the same donor. To remove low-expression genes from differential gene expression
691   analysis (DE analysis), we calculated the average FPKM in 90 datasets and genes with FPKM
692   lower than 1 were excluded from the downstream analysis. 43,585 genes were removed by this
693   filtering, leaving 17,077 genes. 5 LCLs with a very low number of SMRV WGS reads (NA12286,
694   NA12287, NA11930, NA12760, and NA11840), as these could potentially be derived from other
695   SMRV-positive samples sequenced on the same lane, were excluded from the downstream
696   analysis. The count tables generated by STAR were used for DE analysis. DE analysis was
697   performed using the DESeq function in the DESeq2 package. For visualization, count tables
698   were normalized by the counts function in the DESeq2 package with a `normalized=TRUE`
699   option. Genes with p-value lower than 0.05 and changed by at least 2-fold were defined as
700   genes with significant expression changes.
701
702   *HERV-K* k-*mer counting*
703           If a presence/absence-type polymorphic HERV-K contains a unique region which
704   distinguishes it from the other HERV-K loci, the presence/absence pattern of this HERV-K in
705   humans should match to the presence/absence pattern of WGS reads originating from the
706   unique region. To comprehensively find such *k*-mers, we exploited *k*-mer hashing of WGS
707   reads. We first mapped WGS reads to a reference HERV-K (HERV-K113, NC_022518.1) and
708   hashed mapped reads into *k*-mers. We excluded the LTR region from this analysis, because the
709   LTR of HERV-K has a high similarity to SVA. Because of this exclusion criteria, our method
710   captures *k*-mers derived from the HERV-K provirus, but does not detect polymorphisms of solo-
711   LTRs.
712           FASTQ files of 2,504 1kGP high-coverage samples were downloaded from the
713   following URL: 'ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/'. FASTQ files of the Han
714   Chinese trio were downloaded from the following URL:
715   'https://www.internationalgenome.org/data-portal/data-collection/structural-variation'. The
716   FASTQ reads were mapped to the HERV-K reference sequence using Hisat2 with `--mp 2,1 --
717   no-spliced-alignment` options and stored as BAM files. To exclude LTR regions for analysis, the
718   reads mapped to 968-8504 of the HERV-K113 genome (start position is 0-based numbering and
719   end position is 1-based numbering) were used for downstream analysis. To reduce the
720   computational burden of *k*-mer counting, the HERV-K113 genome was split into 50-bp bins with
721   a 10-bp sliding window. Then, sequences of the mapped reads corresponding to the HERV-K

722    50-bp bins were listed from the BAM files using a custom Python script. We defined those
723    sequences as HERV-K *k*-mers. This detected 460,491 different HERV-K *k*-mers from 2,504
724    subjects in 1kGP. The occurrence of each HERV-K *k*-mer was counted by each sample using a
725    custom Python script. The workflow described here was compiled as a Python pipeline and
726    available from the following GitHub repository:
727    'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
728        To normalize the *k*-mer occurrence by the depth of human autosomes, we calculated
729    the chromosome depths using `samtools coverage` command. For this calculation, we used
730    CRAM files of the 2,504 high-coverage WGS provided from 1kGP (downloaded as described in
731    the '*WGS datasets*' section). Then, we calculated the mean depth of chromosome 1 to 22,
732    which we refer to as the autosome depth. The *k*-mer occurrence of each dataset was divided by
733    the calculated autosome and used as a normalized value.
734

735    *Definition of high-frequency presence/absence-type HERV-K k-mers*
736        To perform GWA analysis using polymorphic HERV-K *k*-mers, we decided to focus on
737    HERV-K *k*-mers above a certain frequency threshold. Very rare *k*-mers would often be
738    individual- or population- specific and thus be less informative for findingg association with
739    SNVs from the trans-ethnic datasets. Therefore, we discarded presence/absence-type *k*-mers
740    which were detected in less than 50 subjects (n=8,642) and we defined the remaining ones as
741    high-frequency HERV-K *k*-mers.
742

743    *Clustering of HERV-K k-mers*
744        To perform hierarchical clustering of the frequencies of the presence/absence-type
745    HERV-K *k*-mers by the 26 human populations, the mean *k*-mer frequencies in each population
746    was first calculated. Then the mean *k*-mer frequencies of the 26 populations were clustered with
747    Ward's method in the clutermap function in the seaborn Python package.
748        To perform hierarchical clustering of the Pearson correlation coefficient of the high-
749    frequency presence/absence-type HERV-K *k*-mers, we generated an all-by-all Pearson
750    correlation coefficient matrix for presence/absence patterns of *k*-mers and performed
751    hierarchical clustering of *k*-mers. The Pearson correlation coefficient was calculated by the corr
752    function in the Python Pandas package. The *k*-mers were then clustered by Ward's method in
753    the clutermap function in the seaborn Python package.
754        Prior to GWA analysis, we formally defined clusters using DBSCAN. We clustered the
755    all-by-all Pearson correlation coefficient matrix for presence/absence patterns of *k*-mers. We
756    used the DBSCAN function in the Python scikit-learn package. Any mutation in the HERV-K
757    provirus should actually result in 5 different *k*-mers, because we listed up *k*-mers corresponding
758    to the reference HERV-K sequences scanned by 50-bp bins with a 10-bp window. Therefore,
759    we used 5 for the `min_samples` parameter. We used 2.5 for the `eps` parameter. To determine
760    appropriate epsilon, we performed DBSCAN using 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5, and 4.0. The
761    clustering results were visually cross-referenced with the result of hierarchical clustering, and an
762    epsilon of 2.5 was chosen because it showed good concordance with the results of the
763    hierarchical clustering while defining few enough clusters to allow GWA analysis using a
764    reasonable computational load. The largest cluster defined by DBSCAN contained 175 unique
765    *k*-mers.

766
767 *GWA analysis*
768 The presence and absence of HERV-K *k*-mers in *k*-mer clusters defined by DBSCAN were then
769 converted to categorical values. To reduce the computational cost of GWA analysis, we
770 generated a consensus presence/absence pattern of *k*-mers in each cluster defined by
771 DBSCAN. If a WGS dataset contained more than 80% of the *k*-mers in a *k*-mer cluster, the
772 dataset was considered as a *k*-mer-positive dataset and labeled as 1, while a WGS dataset
773 contained no or 80% or less number of *k*-mers in the *k*-mer cluster, the dataset was considered
774 as a *k*-mer-negative dataset and labeled as 0. These presence/absence binary categorical
775 values were used for the GWA analysis. For SNV annotations, we used GRCh38_v1a
776 downloaded from
777 'ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000_genomes_project/release/20181
778 203_biallelic_SNV/'. We first removed SNVs with low frequency (< 1%), those violating Hardy-
779 Weinberg equilibrium (1e-05), and those with high missing call rate (> 5%) using plink2 software
780 with `--geno 0.05 --hwe 0.00001 --maf 0.01` options. Then the SNVs were pruned using plink2
781 software with `--maf 0.05 --indep-pairwise 100kb 0.5` options. PCA was performed using plink2
782 software using pruned SNVs. Association analysis was performed using plink2 with covariates
783 and binary categorical values representing the presence and absence status of *k*-mers. The sex
784 of WGS datasets and the eigen vectors generated by PCA were used as covariates. Because
785 we performed 599 association tests, we set 8.33e-11 as the genome-wide significant p-value
786 threshold.
787
788 *Detection of HERV-K k*-mer-*associated loci*
789 To determine genome loci associated with the HERV-K *k*-mer clusters in GWA
790 analysis, we first defined genome regions with association by each *k*-mer cluster. If two SNVs
791 with significant p-values were within 1 Mb of one another, we considered those two SNVs to be
792 within the same *k*-mer-associated locus. Otherwise, we considered the two SNVs to be in two
793 separate *k*-mer-associated loci. If a *k*-mer-associated locus harbored 10 or more SNVs, we
794 considered the genome region as a *k*-mer-associated locus, and the largest continuous region
795 containing SNVs with association p values below the 8.33e-11 threshold were defined as a *k*-
796 mer-associated loci. We detected 589 loci from 503 *k*-mer clusters. Because the same locus
797 was detected in multiple GWA analyses with different *k*-mer clusters, we merged 589 genome
798 regions using BEDTools merge command. Finally, we obtained 79 genome loci associated with
799 HERV-K *k*-mers. The HERV-K *k*-mer-associated loci are available from the following GitHub
800 repository: 'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
801
802 *Evaluation of LDfred to detect unknown HERV-K polymorphisms*
803 To assess the sensitivity of LDfred to detect previously unknown polymorphisms, we
804 used structural variations (SVs) in three subjects (NA12878, NA19434, HG00268) called by
805 Audano et al. We extracted deletions that intersect with HERV-K annotated by RepeatMasker
806 using the repeat library version 24.01 from Repbase. We detected 24 deletions in at least one in
807 the three subjects. These 24 deletions ranged from 72 to 9,468-bp. We checked if these HERV-
808 K SVs were located within loci associated with *k*-mer clusters identified by LDfred. Seventeen
809 out of 24 were located with loci associated with *k*-mer clusters. Next we checked the

810 consistency of the presence-absence patterns between *k*-mers and the deletions. When the *k*-
811 mer presence-absence pattern of a *k*-mer cluster and the presence-absence pattern of the
812 overlapping deletions were exactly the same, we considered that the LDfred result accurately
813 reflected the HERV-K polymorphism. For example, if a presence pattern of *k*-mers in a *k*-mer
814 cluster is (NA12878 = +, NA19434 = -, HG00268 = +) and the presence of HERV-K in the
815 associated locus has the same pattern (NA12878 = +, NA19434 = -, HG00268 = +), we
816 considered that LDfred detected the HERV-K polymorphism. On the other hand, if a presence
817 pattern of *k*-mers in a *k*-mer cluster was (NA12878 = +, NA19434 = -, HG00268 = +) and the
818 presence of HERV-K in the associated locus has different pattern (NA12878 = -, NA19434 = -,
819 HG00268 = +), we considered that the LDfred result did not accurately reflect the HERV-K
820 polymorphism. We detected 9 loci associating with *k*-mer clusters which harbor polymorphic
821 HERV-K with consistent presence-absence patterns.
822
823 *Dot matrix analysis*
824 The PacBio alignments to the human genome were downloaded from the following URL:
825 'http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/hgsv_sv_discovery/working/'. To
826 generate a dot matrix between a PacBio sequence and the reference human genome, first,
827 similar sequences between two input sequences were detected and aligned using blastn with `-
828 evalue 1 -word_size 7 -dust no` options. Then, the alignment was visualized by a custom
829 Python script. The script used for this analysis is available from the following GitHub repository:
830 'https://github.com/GenomeImmunobiology/Kojima_et_al_2020'.
831
832 *Detection of SMRV by PCR*
833 The following cell lines were obtained from the NIGMS Human Genetic Cell Repository
834 at the Coriell Institute for Medical Research: GM18998, GM18999, GM12878, GM12399, and
835 GM11920. To confirm the existence of SMRV DNA in LCLs, we designed SMRV-specific
836 primers and amplified SMRV DNA by PCR. Total DNA extracted from GM12399, GM11920,
837 GM18998 were used as PCR templates. PCR primers used are listed in the Supplementary
838 Table 5.
839
840 *Amplification of HERV-K by PCR and mapping to the human genome*
841 Genome regions containing HERV-K in interest were amplified by PCR. Total DNA
842 extracted from GM18998, GM18999, GM12878 were used as PCR templates. The amplicons
843 were barcoded and sequenced using an Oxford Nanopore flongle flow cell. We obtained
844 18,702, 30,651, and 18,604 reads from each subject which passed standard minKNOW v3.6.5
845 quality control from these LCLs, respectively. The reads were mapped to GRCh38DH by bwa
846 mem with the `-Y -K 1000000 -x ont2d` option. Because the HERV-K sequences could
847 potentially be mis-aligned to multiple HERV-K loci, reads harboring sequences which mapping
848 to the non-HERV-K regions at the termini of each PCR amplicon were extracted using a custom
849 script. Finally, we obtained 2,928, 6,676, and 5,660 mapped reads, respectively. PCR primers
850 used are listed in the Supplementary Table 5. A mutation rate of $0.5 \times 10^{-9}$ substitutions per
851 base, per year was assumed to estimate the age of the novel haplotypes (41).
852
853 *Software versions*

854    Python 3.7.4

855    scikit-learn 0.22.1

856    biopython 1.74

857    pandas 0.25.1

858    seaborn 0.10.1

859    pysam 0.15.2

860    MEGA X 10.0.5

861    MAFFT v7.407

862    ete3 3.1.2

863    STAR 2.7.3a

864    R 3.6.1

865    DESeq2 1.22.2

866    BLAST 2.9.0+

867    samtools 1.10

868    bedtools v2.29.2

869    bcftools 1.9

870    Hisat2 version 2.2.0

871    PLINK v2.00a2.3LM

872    prefetch 2.9.3

873    fasterq-dump 2.9.6

874    bamCoverage 3.4.1

875    RepeatMasker 4.0.9

876

877

878

892    **Reference**
893

894    1.    Prüfer K, de Filippo C, Grote S, Mafessoni F, Korlević P, Hajdinjak M, Vernot B, Skov L,
895          Hsieh P, Peyrégne S, Reher D, Hopfe C, Nagel S, Maricic T, Fu Q, Theunert C, Rogers
896          R, Skoglund P, Chintalapati M, Dannemann M, Nelson BJ, Key FM, Rudan P, Kućan Ž,
897          Gušić I, Golovanova L V, Doronichev VB, Patterson N, Reich D, Eichler EE, Slatkin M,
898          Schierup MH, Andrés AM, Kelso J, Meyer M, Pääbo S. 2017. A high-coverage
899          Neandertal genome from Vindija Cave in Croatia. Science 358:655–658.
900    2.    Wildschutte JH, Williams ZH, Montesion M, Subramanian RP, Kidd JM, Coffin JM. 2016.
901          Discovery of unfixed endogenous retrovirus insertions in diverse human populations.
902          Proc Natl Acad Sci 113:E2326 LP-E2334.
903    3.    Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K,
904          Doyle M, FitzHugh W, Funke R, Gage D, Harris K, Heaford A, Howland J, Kann L,
905          Lehoczky J, LeVine R, McEwan P, McKernan K, Meldrim J, Mesirov JP, Miranda C,
906          Morris W, Naylor J, Raymond C, Rosetti M, Santos R, Sheridan A, Sougnez C, Stange-
907          Thomann Y, Stojanovic N, Subramanian A, Wyman D, Rogers J, Sulston J, Ainscough R,
908          Beck S, Bentley D, Burton J, Clee C, Carter N, Coulson A, Deadman R, Deloukas P,
909          Dunham A, Dunham I, Durbin R, French L, Grafham D, Gregory S, Hubbard T, Humphray
910          S, Hunt A, Jones M, Lloyd C, McMurray A, Matthews L, Mercer S, Milne S, Mullikin JC,
911          Mungall A, Plumb R, Ross M, Shownkeen R, Sims S, Waterston RH, Wilson RK, Hillier
912          LW, McPherson JD, Marra MA, Mardis ER, Fulton LA, Chinwalla AT, Pepin KH, Gish
913          WR, Chissoe SL, Wendl MC, Delehaunty KD, Miner TL, Delehaunty A, Kramer JB, Cook
914          LL, Fulton RS, Johnson DL, Minx PJ, Clifton SW, Hawkins T, Branscomb E, Predki P,
915          Richardson P, Wenning S, Slezak T, Doggett N, Cheng JF, Olsen A, Lucas S, Elkin C,
916          Uberbacher E, Frazier M, Gibbs RA, Muzny DM, Scherer SE, Bouck JB, Sodergren EJ,
917          Worley KC, Rives CM, Gorrell JH, Metzker ML, Naylor SL, Kucherlapati RS, Nelson DL,
918          Weinstock GM, Sakaki Y, Fujiyama A, Hattori M, Yada T, Toyoda A, Itoh T, Kawagoe C,
919          Watanabe H, Totoki Y, Taylor T, Weissenbach J, Heilig R, Saurin W, Artiguenave F,
920          Brottier P, Bruls T, Pelletier E, Robert C, Wincker P, Smith DR, Doucette-Stamm L,
921          Rubenfield M, Weinstock K, Lee HM, Dubois J, Rosenthal A, Platzer M, Nyakatura G,
922          Taudien S, Rump A, Yang H, Yu J, Wang J, Huang G, Gu J, Hood L, Rowen L, Madan A,
923          Qin S, Davis RW, Federspiel NA, Abola AP, Proctor MJ, Myers RM, Schmutz J, Dickson
924          M, Grimwood J, Cox DR, Olson M V, Kaul R, Raymond C, Shimizu N, Kawasaki K,
925          Minoshima S, Evans GA, Athanasiou M, Schultz R, Roe BA, Chen F, Pan H, Ramser J,
926          Lehrach H, Reinhardt R, McCombie WR, de la Bastide M, Dedhia N, Blöcker H,
927          Hornischer K, Nordsiek G, Agarwala R, Aravind L, Bailey JA, Bateman A, Batzoglou S,
928          Birney E, Bork P, Brown DG, Burge CB, Cerutti L, Chen HC, Church D, Clamp M, Copley
929          RR, Doerks T, Eddy SR, Eichler EE, Furey TS, Galagan J, Gilbert JG, Harmon C,
930          Hayashizaki Y, Haussler D, Hermjakob H, Hokamp K, Jang W, Johnson LS, Jones TA,
931          Kasif S, Kaspryzk A, Kennedy S, Kent WJ, Kitts P, Koonin E V, Korf I, Kulp D, Lancet D,
932          Lowe TM, McLysaght A, Mikkelsen T, Moran J V, Mulder N, Pollara VJ, Ponting CP,
933          Schuler G, Schultz J, Slater G, Smit AF, Stupka E, Szustakowki J, Thierry-Mieg D,
934          Thierry-Mieg J, Wagner L, Wallis J, Wheeler R, Williams A, Wolf YI, Wolfe KH, Yang SP,
935          Yeh RF, Collins F, Guyer MS, Peterson J, Felsenfeld A, Wetterstrand KA, Patrinos A,
936          Morgan MJ, de Jong P, Catanese JJ, Osoegawa K, Shizuya H, Choi S, Chen YJ,
937          Szustakowki J. 2001. Initial sequencing and analysis of the human genome. Nature
938          409:860–921.
939    4.    Li W, Lin L, Malhotra R, Yang L, Acharya R, Poss M. 2019. A computational framework to
940          assess genome-wide distribution of polymorphic human  endogenous retrovirus-K In
941          human populations. PLoS Comput Biol 15:e1006564.

942    5.       Macfarlane CM, Badge RM. 2015. Genome-wide amplification of proviral sequences
943             reveals new polymorphic HERV-K(HML-2) proviruses in humans and chimpanzees that
944             are absent from genome assemblies. Retrovirology 12:35.

945    6.       Thomas J, Perron H, Feschotte C. 2018. Variation in proviral content among human
946             genomes mediated by LTR recombination. Mob DNA 9:36.

947    7.       Belshaw R, Dawson ALA, Woolven-Allen J, Redding J, Burt A, Tristem M. 2005.
948             Genomewide screening reveals high levels of insertional polymorphism in the human
949             endogenous retrovirus family HERV-K(HML2): implications for present-day activity. J
950             Virol 79:12507–12514.

951    8.       Subramanian RP, Wildschutte JH, Russo C, Coffin JM. 2011. Identification,
952             characterization, and comparative genomic distribution of the HERV-K (HML-2) group of
953             human endogenous retroviruses. Retrovirology 8:90.

954    9.       Bhardwaj N, Montesion M, Roy F, Coffin JM. 2015. Differential expression of HERV-K
955             (HML-2) proviruses in cells and virions of the teratocarcinoma cell line Tera-1. Viruses
956             7:939–968.

957    10.    Horie M, Honda T, Suzuki Y, Kobayashi Y, Daito T, Oshida T, Ikuta K, Jern P, Gojobori T,
958             Coffin JM, Tomonaga K. 2010. Endogenous non-retroviral RNA virus elements in
959             mammalian genomes. Nature 463:84–7.

960    11.    Zhang E, Bell AJ, Wilkie GS, Suárez NM, Batini C, Veal CD, Armendáriz-Castillo I,
961             Neumann R, Cotton VE, Huang Y, Porteous DJ, Jarrett RF, Davison AJ, Royle NJ. 2017.
962             Inherited Chromosomally Integrated Human Herpesvirus 6 Genomes Are Ancient, Intact,
963             and Potentially Able To Reactivate from Telomeres. J Virol 91.

964    12.    Liu X, Kosugi S, Koide R, Kawamura Y, Ito J, Miura H, Matoba N, Matsuzaki M, Fujita M,
965             Kamada AJ, Nakagawa H, Tamiya G, Matsuda K, Murakami Y, Kubo M, Aswad A, Sato
966             K, Momozawa Y, Ohashi J, Terao C, Yoshikawa T, Parrish NF, Kamatani Y. 2020.
967             Endogenization and excision of human herpesvirus 6 in human genomes. PLoS Genet
968             16:e1008915.

969    13.    Aswad A, Aimola G, Wight D, Roychoudhury P, Zimmermann C, Hill J, Lassner D, Xie H,
970             Huang M-L, Parrish NF, Schultheiss H-P, Venturini C, Lager S, Smith GCS, Charnock-
971             Jones DS, Breuer J, Greninger AL, Kaufer BB. 2020. Evolutionary history of endogenous
972             Human Herpesvirus 6 reflects human migration out of Africa. Mol Biol Evol.

973    14.    Weismann A. 1893. The germ-plasm: a theory of heredity. Scribner's.

974    15.    Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y,
975             Ye K, Jun G, Fritz MH-Y, Konkel MK, Malhotra A, Stütz AM, Shi X, Casale FP, Chen J,
976             Hormozdiari F, Dayama G, Chen K, Malig M, Chaisson MJP, Walter K, Meiers S, Kashin
977             S, Garrison E, Auton A, Lam HYK, Mu XJ, Alkan C, Antaki D, Bae T, Cerveira E, Chines
978             P, Chong Z, Clarke L, Dal E, Ding L, Emery S, Fan X, Gujral M, Kahveci F, Kidd JM,
979             Kong Y, Lameijer E-W, McCarthy S, Flicek P, Gibbs RA, Marth G, Mason CE, Menelaou
980             A, Muzny DM, Nelson BJ, Noor A, Parrish NF, Pendleton M, Quitadamo A, Raeder B,
981             Schadt EE, Romanovitch M, Schlattl A, Sebra R, Shabalin AA, Untergasser A, Walker JA,
982             Wang M, Yu F, Zhang C, Zhang J, Zheng-Bradley X, Zhou W, Zichner T, Sebat J, Batzer
983             MA, McCarroll SA, Mills RE, Gerstein MB, Bashir A, Stegle O, Devine SE, Lee C, Eichler
984             EE, Korbel JO. 2015. An integrated map of structural variation in 2,504 human genomes.
985             Nature 526:75–81.

986    16.    Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera A V, Lowther
987             C, Gauthier LD, Wang H, Watts NA, Solomonson M, O'Donnell-Luria A, Baumann A,
988             Munshi R, Walker M, Whelan CW, Huang Y, Brookings T, Sharpe T, Stone MR, Valkanas
989             E, Fu J, Tiao G, Laricchia KM, Ruano-Rubio V, Stevens C, Gupta N, Cusick C, Margolin
990             L, Taylor KD, Lin HJ, Rich SS, Post WS, Chen Y-DI, Rotter JI, Nusbaum C, Philippakis A,
991             Lander E, Gabriel S, Neale BM, Kathiresan S, Daly MJ, Banks E, MacArthur DG,

992   Talkowski ME. 2020. A structural variation reference for medical and population genetics.
993   Nature 581:444–451.

994 17. Almarri MA, Bergström A, Prado-Martinez J, Yang F, Fu B, Dunham AS, Chen Y, Hurles
995   ME, Tyler-Smith C, Xue Y. 2020. Population Structure, Stratification, and Introgression of
996   Human Structural Variation. Cell 182:189-199.e15.

997 18. Moustafa A, Xie C, Kirkness E, Biggs W, Wong E, Turpaz Y, Bloom K, Delwart E, Nelson
998   KE, Venter JC, Telenti A. 2017. The blood DNA virome in 8,000 humans. PLoS Pathog
999   13:e1006292.

1000 19. Liu S, Huang S, Chen F, Zhao L, Yuan Y, Francis SS, Fang L, Li Z, Lin L, Liu R, Zhang Y,
1001   Xu H, Li S, Zhou Y, Davies RW, Liu Q, Walters RG, Lin K, Ju J, Korneliussen T, Yang
1002   MA, Fu Q, Wang J, Zhou L, Krogh A, Zhang H, Wang W, Chen Z, Cai Z, Yin Y, Yang H,
1003   Mao M, Shendure J, Wang J, Albrechtsen A, Jin X, Nielsen R, Xu X. 2018. Genomic
1004   Analyses from Non-invasive Prenatal Testing Reveal Genetic Associations, Patterns of
1005   Viral Infections, and Chinese Population History. Cell 175:347-359.e14.

1006 20. Bergström A, McCarthy SA, Hui R, Almarri MA, Ayub Q, Danecek P, Chen Y, Felkel S,
1007   Hallast P, Kamm J, Blanché H, Deleuze J-F, Cann H, Mallick S, Reich D, Sandhu MS,
1008   Skoglund P, Scally A, Xue Y, Durbin R, Tyler-Smith C. 2020. Insights into human genetic
1009   variation and population history from 929 diverse genomes. Science 367.

1010 21. Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL,
1011   McCarthy S, McVean GA, Abecasis GR. 2015. A global reference for human genetic
1012   variation. Nature 526:68–74.

1013 22. Sun R, Grogan E, Shedd D, Bykovsky AF, Kushnaryov VM, Grossberg SE, Miller G.
1014   1995. Transmissible retrovirus in Epstein-Barr virus-producer B95-8 cells. Virology
1015   209:374–383.

1016 23. Cheval J, Muth E, Gonzalez G, Coulpier M, Beurdeley P, Cruveiller S, Eloit M. 2019.
1017   Adventitious Virus Detection in Cells by High-Throughput Sequencing of Newly
1018   Synthesized RNAs: Unambiguous Differentiation of Cell Infection from Carryover of Viral
1019   Nucleic Acids. mSphere 4.

1020 24. Rhim JS, Schell K, Creasy B, Case W. 1969. Biological characteristics and viral
1021   susceptibility of an African green monkey kidney cell line (Vero). Proc Soc Exp Biol Med
1022   Soc Exp Biol Med (New York, NY) 132:670–678.

1023 25. Feusier J, Watkins WS, Thomas J, Farrell A, Witherspoon DJ, Baird L, Ha H, Xing J,
1024   Jorde LB. 2019. Pedigree-based estimation of human mobile element retrotransposition
1025   rates. Genome Res 29:1567–1577.

1026 26. Sasani TA, Pedersen BS, Gao Z, Baird L, Przeworski M, Jorde LB, Quinlan AR. 2019.
1027   Large, three-generation human families reveal post-zygotic mosaicism and variability in
1028   germline mutation accumulation. Elife 8.

1029 27. Oda T, Ikeda S, Watanabe S, Hatsushika M, Akiyama K, Mitsunobu F. 1988. Molecular
1030   cloning, complete nucleotide sequence, and gene structure of the provirus genome of a
1031   retrovirus produced in a human lymphoblastoid cell line. Virology 167:468–476.

1032 28. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas MA,
1033   Gonzàlez-Porta M, Kurbatova N, Griebel T, Ferreira PG, Barann M, Wieland T, Greger L,
1034   van Iterson M, Almlöf J, Ribeca P, Pulyakhina I, Esser D, Giger T, Tikhonov A, Sultan M,
1035   Bertier G, MacArthur DG, Lek M, Lizano E, Buermans HPJ, Padioleau I, Schwarzmayr T,
1036   Karlberg O, Ongen H, Kilpinen H, Beltran S, Gut M, Kahlem K, Amstislavskiy V, Stegle O,
1037   Pirinen M, Montgomery SB, Donnelly P, McCarthy MI, Flicek P, Strom TM, Lehrach H,
1038   Schreiber S, Sudbrak R, Carracedo A, Antonarakis SE, Häsler R, Syvänen A-C, van
1039   Ommen G-J, Brazma A, Meitinger T, Rosenstiel P, Guigó R, Gut IG, Estivill X,
1040   Dermitzakis ET. 2013. Transcriptome and genome sequencing uncovers functional
1041   variation in humans. Nature 501:506–511.

1042   29.   Tan K-T, Ding L-W, Sun Q-Y, Lao Z-T, Chien W, Ren X, Xiao J-F, Loh XY, Xu L, Lill M,
1043         Mayakonda A, Lin D-C, Yang H, Koeffler HP. 2018. Profiling the B/T cell receptor
1044         repertoire of lymphocyte derived cell lines. BMC Cancer 18:940.
1045   30.   Sebastian NT, Zaikos TD, Terry V, Taschuk F, McNamara LA, Onafuwa-Nuga A, Yucha
1046         R, Signer RAJ, Riddell JI V, Bixby D, Markowitz N, Morrison SJ, Collins KL. 2017. CD4 is
1047         expressed on a heterogeneous subset of hematopoietic progenitors, which  persistently
1048         harbor CXCR4 and CCR5-tropic HIV proviral genomes in vivo. PLoS Pathog
1049         13:e1006509.
1050   31.   McHugh D, Myburgh R, Caduff N, Spohn M, Kok YL, Keller CW, Murer A, Chatterjee B,
1051         Rühl J, Engelmann C, Chijioke O, Quast I, Shilaih M, Strouvelle VP, Neumann K, Menter
1052         T, Dirnhofer S, Lam JK, Hui KF, Bredl S, Schlaepfer E, Sorce S, Zbinden A, Capaul R,
1053         Lünemann JD, Aguzzi A, Chiang AK, Kempf W, Trkola A, Metzner KJ, Manz MG,
1054         Grundhoff A, Speck RF, Münz C. 2020. EBV renders B cells susceptible to HIV-1 in
1055         humanized mice. Life Sci alliance 3.
1056   32.   Telford M, Navarro A, Santpere G. 2018. Whole genome diversity of inherited
1057         chromosomally integrated HHV-6 derived from healthy individuals of diverse geographic
1058         origin. Sci Rep 8:3472.
1059   33.   Martens UM, Zijlmans JM, Poon SS, Dragowska W, Yui J, Chavez EA, Ward RK,
1060         Lansdorp PM. 1998. Short telomeres on human chromosome 17p. Nat Genet 18:76–80.
1061   34.   Payer LM, Steranka JP, Yang WR, Kryatova M, Medabalimi S, Ardeljan D, Liu C, Boeke
1062         JD, Avramopoulos D, Burns KH. 2017. Structural variants caused by Alu insertions are
1063         associated with risks for many  human diseases. Proc Natl Acad Sci U S A 114:E3984–
1064         E3992.
1065   35.   Wallace AD, Wendt GA, Barcellos LF, de Smith AJ, Walsh KM, Metayer C, Costello JF,
1066         Wiemels JL, Francis SS. 2018. To ERV Is Human: A Phenotype-Wide Scan Linking
1067         Polymorphic Human Endogenous  Retrovirus-K Insertions to Complex Phenotypes. Front
1068         Genet 9:298.
1069   36.   Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, McMahon A,
1070         Morales J, Mountjoy E, Sollis E, Suveges D, Vrousgou O, Whetzel PL, Amode R, Guillen
1071         JA, Riat HS, Trevanion SJ, Hall P, Junkins H, Flicek P, Burdett T, Hindorff LA,
1072         Cunningham F, Parkinson H. 2019. The NHGRI-EBI GWAS Catalog of published
1073         genome-wide association studies, targeted  arrays and summary statistics 2019. Nucleic
1074         Acids Res 47:D1005–D1012.
1075   37.   Audano PA, Sulovari A, Graves-Lindsay TA, Cantsilieris S, Sorensen M, Welch AE,
1076         Dougherty ML, Nelson BJ, Shah A, Dutcher SK, Warren WC, Magrini V, McGrath SD, Li
1077         YI, Wilson RK, Eichler EE. 2019. Characterizing the Major Structural Variant Alleles of the
1078         Human Genome. Cell 176:663-675.e19.
1079   38.   De Coster W, Van Broeckhoven C. 2019. Newest Methods for Detecting Structural
1080         Variations. Trends Biotechnol 37:973–982.
1081   39.   Linardopoulou E V, Williams EM, Fan Y, Friedman C, Young JM, Trask BJ. 2005. Human
1082         subtelomeres are hot spots of interchromosomal recombination and segmental
1083         duplication. Nature 437:94–100.
1084   40.   Chaisson MJP, Sanders AD, Zhao X, Malhotra A, Porubsky D, Rausch T, Gardner EJ,
1085         Rodriguez OL, Guo L, Collins RL, Fan X, Wen J, Handsaker RE, Fairley S, Kronenberg
1086         ZN, Kong X, Hormozdiari F, Lee D, Wenger AM, Hastie AR, Antaki D, Anantharaman T,
1087         Audano PA, Brand H, Cantsilieris S, Cao H, Cerveira E, Chen C, Chen X, Chin C-S,
1088         Chong Z, Chuang NT, Lambert CC, Church DM, Clarke L, Farrell A, Flores J, Galeev T,
1089         Gorkin DU, Gujral M, Guryev V, Heaton WH, Korlach J, Kumar S, Kwon JY, Lam ET, Lee
1090         JE, Lee J, Lee W-P, Lee SP, Li S, Marks P, Viaud-Martinez K, Meiers S, Munson KM,
1091         Navarro FCP, Nelson BJ, Nodzak C, Noor A, Kyriazopoulou-Panagiotopoulou S, Pang
1092         AWC, Qiu Y, Rosanio G, Ryan M, Stütz A, Spierings DCJ, Ward A, Welch AE, Xiao M,

1093       Xu W, Zhang C, Zhu Q, Zheng-Bradley X, Lowy E, Yakneen S, McCarroll S, Jun G, Ding
1094       L, Koh CL, Ren B, Flicek P, Chen K, Gerstein MB, Kwok P-Y, Lansdorp PM, Marth GT,
1095       Sebat J, Shi X, Bashir A, Ye K, Devine SE, Talkowski ME, Mills RE, Marschall T, Korbel
1096       JO, Eichler EE, Lee C. 2019. Multi-platform discovery of haplotype-resolved structural
1097       variation in human  genomes. Nat Commun 10:1784.
1098  41.    Scally A, Durbin R. 2012. Revising the human mutation rate: implications for
1099       understanding human evolution. Nat Rev Genet. England.
1100  42.    Hughes JF, Coffin JM. 2005. Human endogenous retroviral elements as indicators of
1101       ectopic recombination events  in the primate genome. Genetics 171:1183–1194.
1102  43.    Volleth M, Zenker M, Joksic I, Liehr T. 2020. Long-term Culture of EBV-induced Human
1103       Lymphoblastoid Cell Lines Reveals Chromosomal  Instability. J Histochem Cytochem  Off
1104       J  Histochem Soc 68:239–251.
1105  44.    Nanbo A, Inoue K, Adachi-Takasawa K, Takada K. 2002. Epstein-Barr virus RNA confers
1106       resistance to interferon-alpha-induced apoptosis in  Burkitt's lymphoma. EMBO J 21:954–
1107       965.
1108  45.    Furuta R, Yasunaga J-I, Miura M, Sugata K, Saito A, Akari H, Ueno T, Takenouchi N,
1109       Fujisawa J-I, Koh K-R, Higuchi Y, Mahgoub M, Shimizu M, Matsuda F, Melamed A,
1110       Bangham CR, Matsuoka M. 2017. Human T-cell leukemia virus type 1 infects multiple
1111       lineage hematopoietic cells in  vivo. PLoS Pathog 13:e1006722.
1112  46.    Mahé D, Matusali G, Deleage C, Alvarenga RLLS, Satie A-P, Pagliuzza A, Mathieu R,
1113       Lavoué S, Jégou B, de França LR, Chomont N, Houzet L, Rolland AD, Dejucq-Rainsford
1114       N. 2020. Potential for virus endogenization in humans through testicular germ cell
1115       infection: the case of HIV. bioRxiv 2020.06.04.135657.
1116  47.    Stoye JP. 2006. Koala retrovirus: a genome invasion in real time. Genome Biol 7:241.
1117  48.    Gaccioli F, Lager S, de Goffau MC, Sovio U, Dopierala J, Gong S, Cook E, Sharkey A,
1118       Moffett A, Lee WK, Delles C, Venturini C, Breuer J, Parkhill J, Peacock SJ, Charnock-
1119       Jones DS, Smith GCS. 2020. Fetal inheritance of chromosomally integrated human
1120       herpesvirus 6 predisposes the mother to pre-eclampsia. Nat Microbiol.
1121  49.    Gravel A, Dubuc I, Morissette G, Sedlak RH, Jerome KR, Flamand L. 2015. Inherited
1122       chromosomally integrated human herpesvirus 6 as a predisposing risk factor for the
1123       development of angina pectoris. Proc Natl Acad Sci 112:8058 LP – 8063.
1124  50.    Zhou W, Emery SB, Flasch DA, Wang Y, Kwan KY, Kidd JM, Moran J V, Mills RE. 2020.
1125       Identification and characterization of occult human-specific LINE-1 insertions using  long-
1126       read sequencing technology. Nucleic Acids Res 48:1146–1163.
1127  51.    Ewing AD, Smits N, Sanchez-Luque FJ, Faivre J, Brennan PM, Richardson SR,
1128       Cheetham SW, Faulkner GJ. 2020. Nanopore Sequencing Enables Comprehensive
1129       Transposable Element Epigenomic Profiling. Mol Cell.
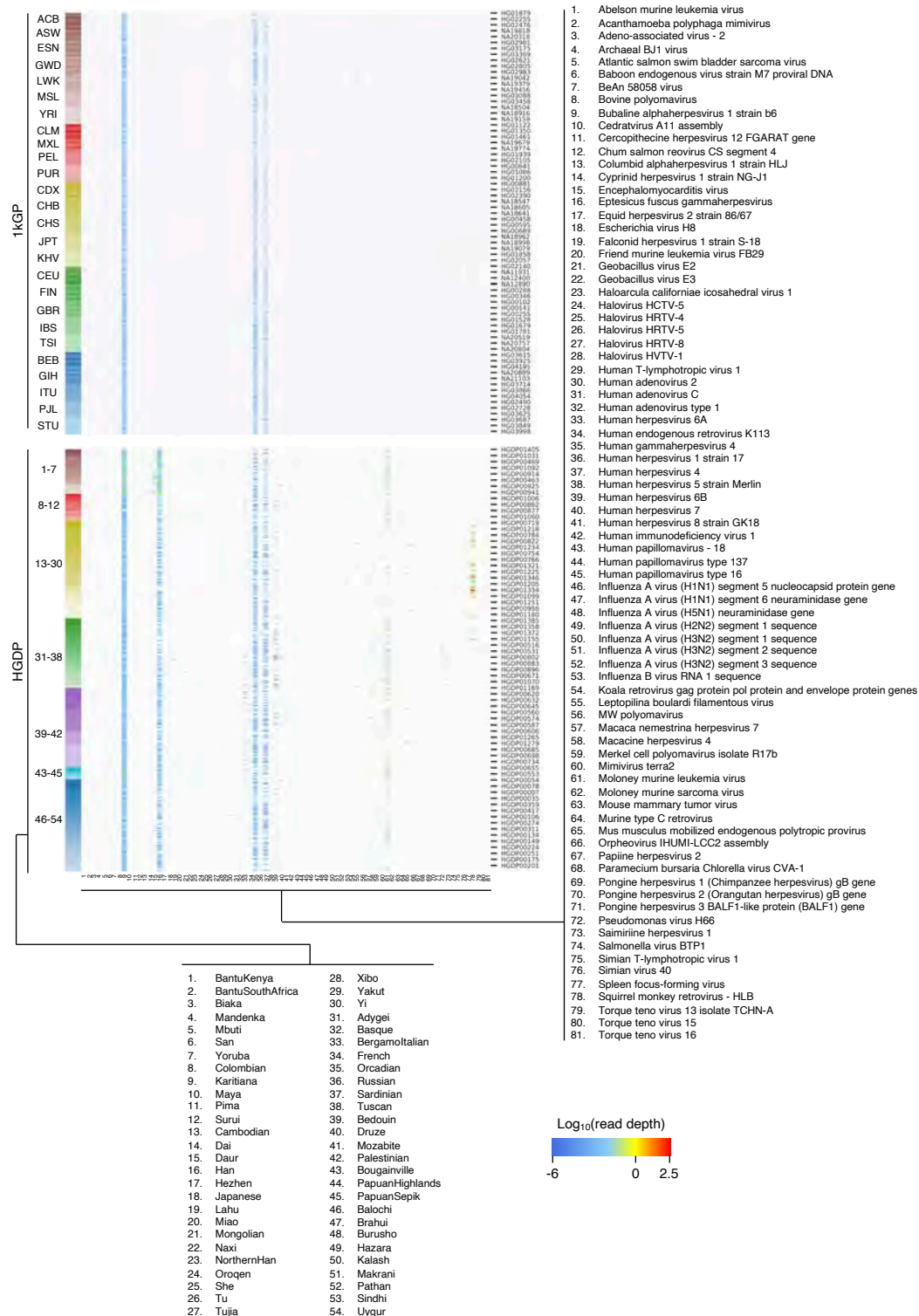1130  52.    Jha AR, Pillai SK, York VA, Sharp ER, Storm EC, Wachter DJ, Martin JN, Deeks SG,
1131       Rosenberg MG, Nixon DF, Garrison KE. 2009. Cross-sectional dating of novel haplotypes
1132       of HERV-K 113 and HERV-K 115 indicate  these proviruses originated in Africa before
1133       Homo sapiens. Mol Biol Evol 26:2617–2626.
1134  53.    Turner G, Barbulescu M, Su M, Jensen-Seaman MI, Kidd KK, Lenz J. 2001. Insertional
1135       polymorphisms of full-length endogenous retroviruses in humans. Curr Biol 11:1531–
1136       1535.
1137  54.    Dewannieux M, Harper F, Richaud A, Letzelter C, Ribet D, Pierron G, Heidmann T. 2006.
1138       Identification of an infectious progenitor for the multiple-copy HERV-K human
1139       endogenous retroelements. Genome Res 16:1548–1556.
1140  55.    Green ED, Gunter C, Biesecker LG, Di Francesco V, Easter CL, Feingold EA, Felsenfeld
1141       AL, Kaufman DJ, Ostrander EA, Pavan WJ, Phillippy AM, Wise AL, Dayal JG, Kish BJ,
1142       Mandich A, Wellington CR, Wetterstrand KA, Bates SA, Leja D, Vasquez S, Gahl WA,
1143       Graham BJ, Kastner DL, Liu P, Rodriguez LL, Solomon BD, Bonham VL, Brody LC,
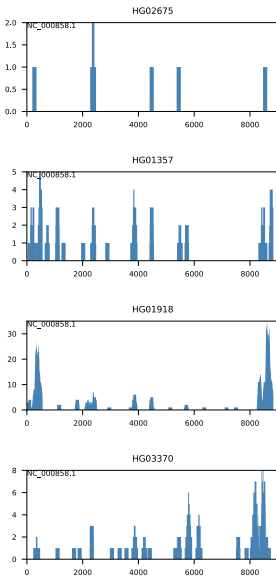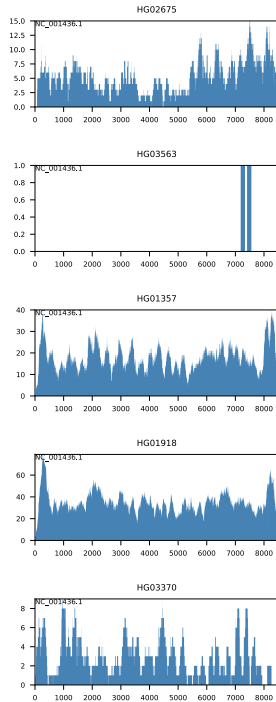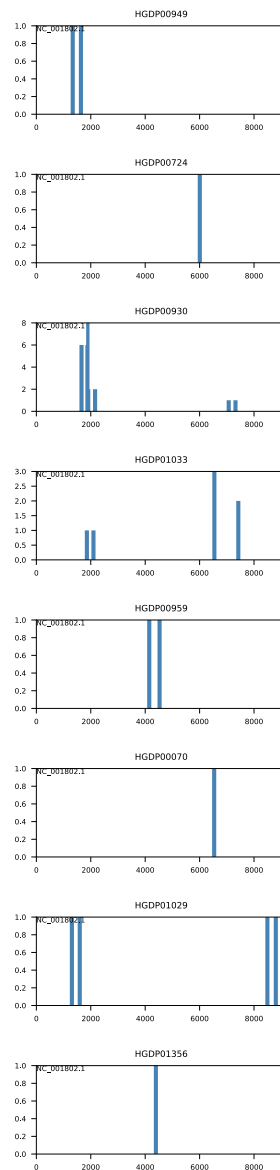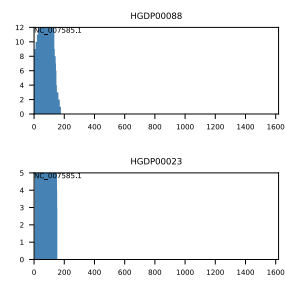
1144         Hutter CM, Manolio TA. 2020. Strategic vision for improving human health at The
1145         Forefront of Genomics. Nature 586:683–692.
1146   56.   Kumata R, Ito J, Takahashi K, Suzuki T, Sato K. 2020. A tissue level atlas of the healthy
1147         human virome. BMC Biol 18:55.
1148
1149

1kGP

ACB
ASW
ESN
GWD
LWK
MSL
YRI
CLM
MXL
PEL
PUR
CDX
CHB
CHS
JPT
KHV
CEU
FIN
GBR
IBS
TSI
BEB
GIH
ITU
PJL
STU

HGDP

1-7
8-12
13-30
31-38
39-42
43-45
46-54

1. Abelson murine leukemia virus
2. Acanthamoeba polyphaga mimivirus
3. Adeno-associated virus - 2
4. Archaeal BJ1 virus
5. Atlantic salmon swim bladder sarcoma virus
6. Baboon endogenous virus strain M7 proviral DNA
7. BeAn 58058 virus
8. Bovine polyomavirus
9. Bubaline alphaherpesvirus 1 strain b6
10. Cedratvirus A11 assembly
11. Cercopithecine herpesvirus 12 FGARAT gene
12. Chum salmon reovirus CS segment 4
13. Columbid alphaherpesvirus 1 strain HLJ
14. Cyprinid herpesvirus 1 strain NG-J1
15. Encephalomyocarditis virus
16. Eptesicus fuscus gammaherpesvirus
17. Equid herpesvirus 2 strain 86/67
18. Escherichia virus H8
19. Falconid herpesvirus 1 strain S-18
20. Friend murine leukemia virus FB29
21. Geobacillus virus E2
22. Geobacillus virus E3
23. Haloarcula californiae icosahedral virus 1
24. Halovirus HCTV-5
25. Halovirus HRTV-4
26. Halovirus HRTV-5
27. Halovirus HRTV-8
28. Halovirus HVTV-1
29. Human T-lymphotropic virus 1
30. Human adenovirus 2
31. Human adenovirus C
32. Human adenovirus type 1
33. Human herpesvirus 6A
34. Human endogenous retrovirus K113
35. Human gammaherpesvirus 4
36. Human herpesvirus 1 strain 17
37. Human herpesvirus 4
38. Human herpesvirus 5 strain Merlin
39. Human herpesvirus 6B
40. Human herpesvirus 7
41. Human herpesvirus 8 strain GK18
42. Human immunodeficiency virus 1
43. Human papillomavirus - 18
44. Human papillomavirus type 137
45. Human papillomavirus type 16
46. Influenza A virus (H1N1) segment 5 nucleocapsid protein gene
47. Influenza A virus (H1N1) segment 6 neuraminidase gene
48. Influenza A virus (H5N1) neuraminidase gene
49. Influenza A virus (H2N2) segment 1 sequence
50. Influenza A virus (H3N2) segment 1 sequence
51. Influenza A virus (H3N2) segment 2 sequence
52. Influenza A virus (H3N2) segment 3 sequence
53. Influenza B virus RNA 1 sequence
54. Koala retrovirus gag protein pol protein and envelope protein genes
55. Leptopilina boulardi filamentous virus
56. MW polyomavirus
57. Macaca nemestrina herpesvirus 7
58. Macacine herpesvirus 4
59. Merkel cell polyomavirus isolate R17b
60. Mimivirus terra2
61. Moloney murine leukemia virus
62. Moloney murine sarcoma virus
63. Mouse mammary tumor virus
64. Murine type C retrovirus
65. Mus musculus mobilized endogenous polytropic provirus
66. Orpheovirus IHUMI-LCC2 assembly
67. Papiine herpesvirus 2
68. Paramecium bursaria Chlorella virus CVA-1
69. Pongine herpesvirus 1 (Chimpanzee herpesvirus) gB gene
70. Pongine herpesvirus 2 (Orangutan herpesvirus) gB gene
71. Pongine herpesvirus 3 BALF1-like protein (BALF1) gene
72. Pseudomonas virus H66
73. Saimiriine herpesvirus 1
74. Salmonella virus BTP1
75. Simian T-lymphotropic virus 1
76. Simian virus 40
77. Spleen focus-forming virus
78. Squirrel monkey retrovirus - HLB
79. Torque teno virus 13 isolate TCHN-A
80. Torque teno virus 15
81. Torque teno virus 16

| 1. | BantuKenya | 28. | Xibo |
| 2. | BantuSouthAfrica | 29. | Yakut |
| 3. | Biaka | 30. | Yi |
| 4. | Mandenka | 31. | Adygei |
| 5. | Mbuti | 32. | Basque |
| 6. | San | 33. | BergamoItalian |
| 7. | Yoruba | 34. | French |
| 8. | Colombian | 35. | Orcadian |
| 9. | Karitiana | 36. | Russian |
| 10. | Maya | 37. | Sardinian |
| 11. | Pima | 38. | Tuscan |
| 12. | Surui | 39. | Bedouin |
| 13. | Cambodian | 40. | Druze |
| 14. | Dai | 41. | Mozabite |
| 15. | Daur | 42. | Palestinian |
| 16. | Han | 43. | Bougainville |

Log$_{10}$(read depth)

-6          0    2.5

1150
1151

**Supplementary Figure 1** Virus search from 1kGP and HGDP WGS
Heatmap shows the read depth of viruses with at least one read in at least one dataset from
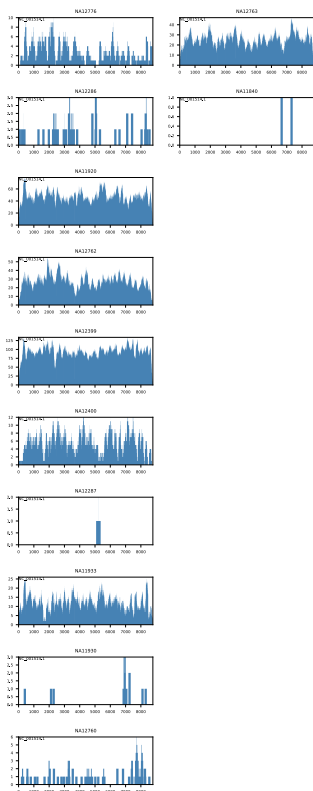2,504 1kGP as well as 808 HDGP datasets.

1155

**A** Simian T-lymphotropic virus 1(+) samples in 1kGP

**B** HTLV-1(+) samples in 1kGP

**C** Human immunodeficiency virus 1(+) samples in HGDP

**D** Chum salmon reovirus CS segment 4 (+) samples in HGDP

1156
1157 **Supplementary Figure 2** Read depth of STLV-1, HIV-1, and chum salmon reovirus
1158 Depth of reads mapping to STLV-1 (A), HTLV-1 (B), HIV-1 (C), and chum salmon reovirus (D)
1159 are shown. X-axis and Y-axis show the genome position of indicated virus and the depth of
1160 reads mapping to the indicated virus, respectively. The name of each dataset is shown at the
1161 top of the panel.
1162

**A** 1kGP, SMRV(+) samples



**B** HGDP, SMRV(+) samples



1163
1164 **Supplementary Figure 3** Read depth of SMRV and HTLV-1
1165 Depth of reads from 1kGP (A) and HGDP (B) datasets mapping to SMRV-H are shown. X-axis
1166 and Y-axis show the genome position of indicated virus and the depth of reads mapping to the
1167 indicated virus, respectively. The name of each dataset is shown at the top of the panel.
1168

**A**

1kGP SMRV hybrid reads



**B**

1kGP HTLV-1 hybrid reads



**C**

HGDP SMRV hybrid reads



1169

1170 **Supplementary Figure 4** Virus-chromosome hybrid reads

1171 WGS read pairs which are mapped to both the virus genome and the human genome are

1172 shown. Gray bar in the top of each panel shows the SMRV-H (A), HTLV-1 (B), and SMRV-H

1173 (C). LTR are shown as dark gray rectangles. Gray bars in the bottom of each panel show the

1174 chromosome 1 to 22, X, and Y (from left to right). Read-1 and Read-2 of a read-pair are

1175 connected with a line. Reads mapping to forward and reverse directions are shown as blue and

1176 red dots, respectively. The reads mapping to left LTR was kept when a read was multi-mapped

1177 to both left and right LTR. The genome position of reads mapping only to right LTR were

1178 replaced to the left LTR. The name of each dataset is shown at the top of the panel.

1179

**A**



**B**



**C**



1180
1181  **Supplementary Figure 5** Differential gene analysis between SMRV-positive and -negative
1182  samples
1183  A. Depth of RNA-seq reads mapping to SMRV-H. Geuvadis RNA-seq datasets were used for
1184  the analysis. All datasets are shown with the same scale of the Y-axis.
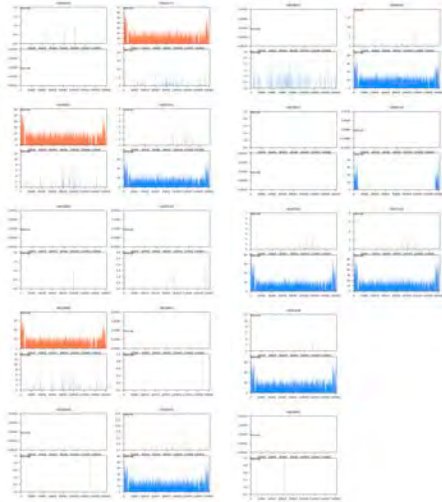1185  B. Correlation of the abundance of reads mapping to SMRV between WGS and RNA-seq. LCLs
1186  producing both WGS and RNA-seq were used for this analysis.
1187  C. MA plot showing the differences of gene expression and the normalized read counts. Two
1188  genes with differential expression are shown as blue dots.
1189

1kGP, HHV-6(+) samples
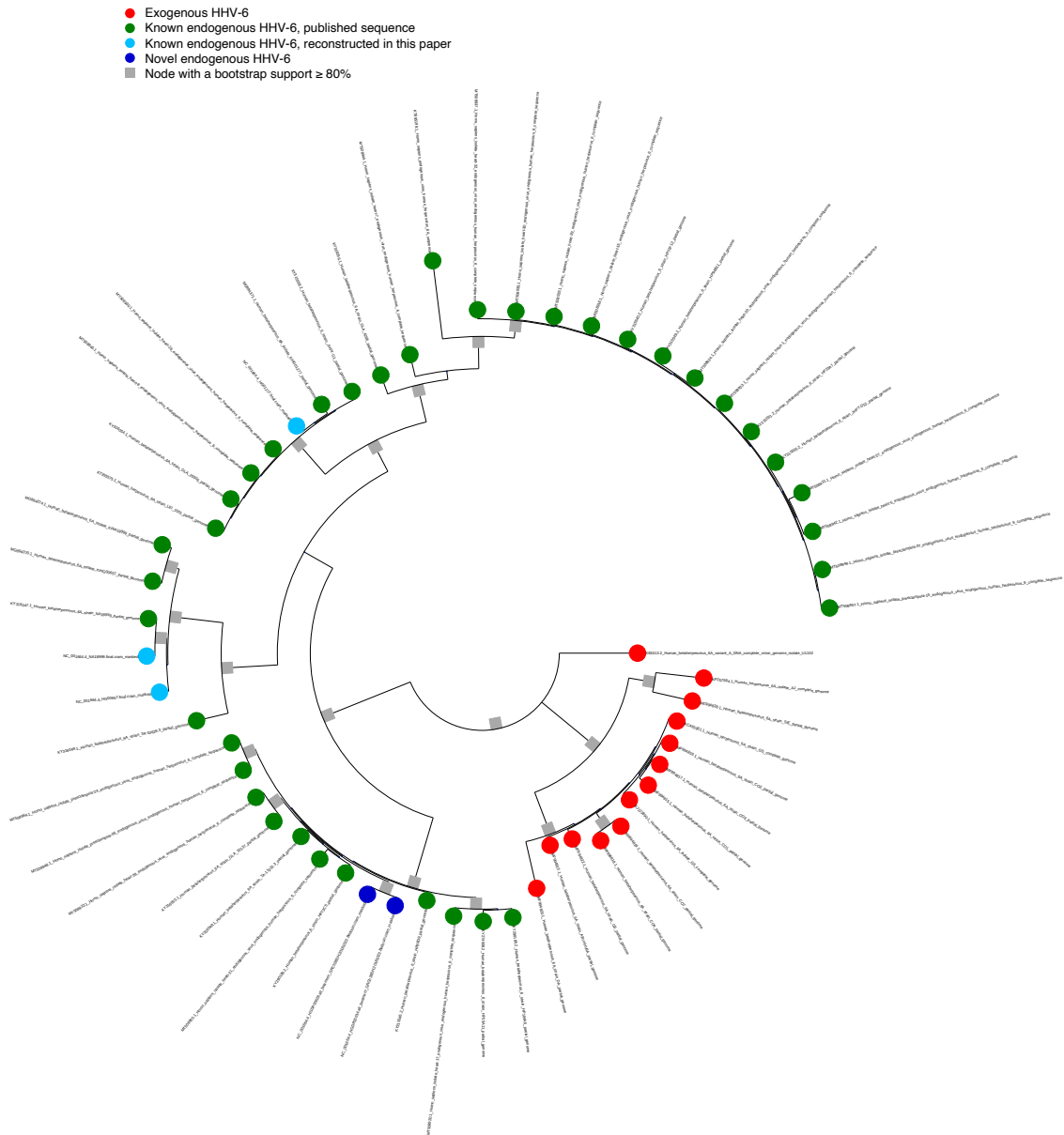


HGDP, HHV-6(+) samples



1190
1191 **Supplementary Figure 6** Read depth of HHV-6A and HHV-6B
1192 Depth of reads mapping to HHV-6A and HHV-6B are shown. Read depth of HHV-6A and HHV-
1193 6B are shown as orange and blue coverage plots, respectively. X-axis and Y-axis show the
1194 genome position of indicated virus and the depth of reads mapping to the indicated virus,
1195 respectively. The name of each dataset is shown at the top of the panel.
1196

**A**

```
                                                        NC 001664.4 (U1102)
              ┌──────────────────── Reconstructed from 5x
              │                          ┌─ Reconstructed from 10x
              │                          │  ┌─ MG894374.1
              │                     100 ─┤  │
              │                          │  ├─ Reconstructed from 15x
              │                    59 ───┤  │
              │                          │  ├─ Reconstructed from 20x
              │                          │  │
              │                   100 ───┤  ├─ Reconstructed from 30x
              │                          │  │
              │                          │  ├─ Reconstructed from 35x
              │                          │  └─ KY316047.1 (Sanger sequencing)

     ├──────────────┤
          0.0020
```

**B**

| WGS depth | Reconstructed length (nt) | Reconstructed length (percent to U1102) |
|---|---|---|
| 5x | 141,178 | 88.5806071 |
| 10x | 152,679 | 95.796785 |
| 15x | 153,468 | 96.2918345 |
| 20x | 154,062 | 96.6645334 |
| 30x | 154,634 | 97.0234286 |
| 35x | 154,238 | 96.7749627 |
| -------- | -------- | -------- |
| U1102 | 159,378 | 100 |

1197

1198 **Supplementary Figure 7** Accuracy of endogenous HHV-6 sequence reconstruction

1199 A. Phylogenetic analysis of reconstructed HHV-6A sequences. HHV-6A in NA18999

1200 reconstructed from 35x, 30x, 20x, 15x, 10x, and 5x autosome depths are aligned with the

1201 reference HHV-6A (U1102) as well as the published sequences of HHV-6A Sanger sequenced

1202 using NA18999 DNA (KY316047.1) or reconstructed from a WGS of NA18999 (MG894374.1).

1203 B. The lengths of endogenous HHV-6A reconstructed from various WGS read depths.

1204

Legend:
- Exogenous HHV-6
- Known endogenous HHV-6, published sequence
- Known endogenous HHV-6, reconstructed in this paper
- Novel endogenous HHV-6
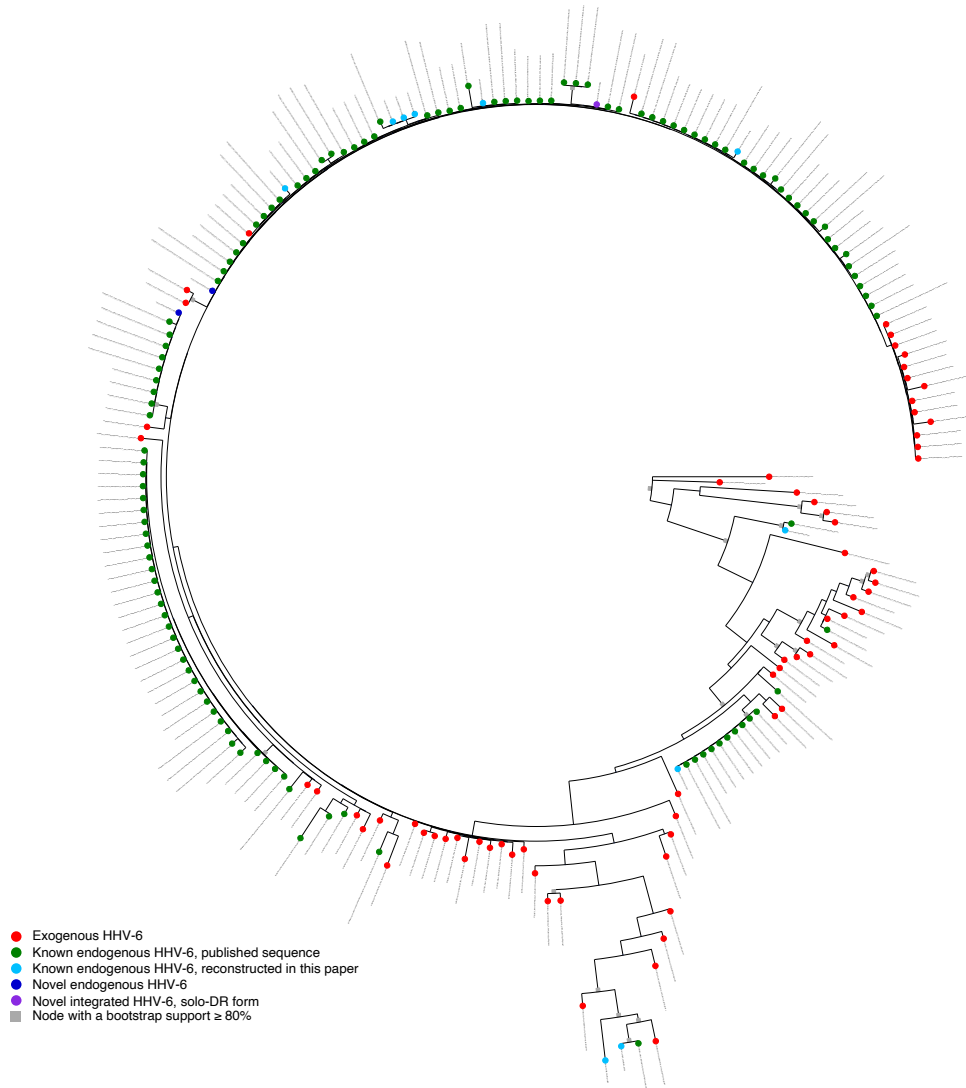- Node with a bootstrap support ≥ 80%

1205
1206 **Supplementary Figure 8** Phylogenetic tree of HHV-6A U with leaf names
1207 Phylogenetic trees inferred from U regions of HHV-6A. The publicly available sequences of
1208 endogenous and exogenous HHV-6A as well as ones reconstructed in the present study were
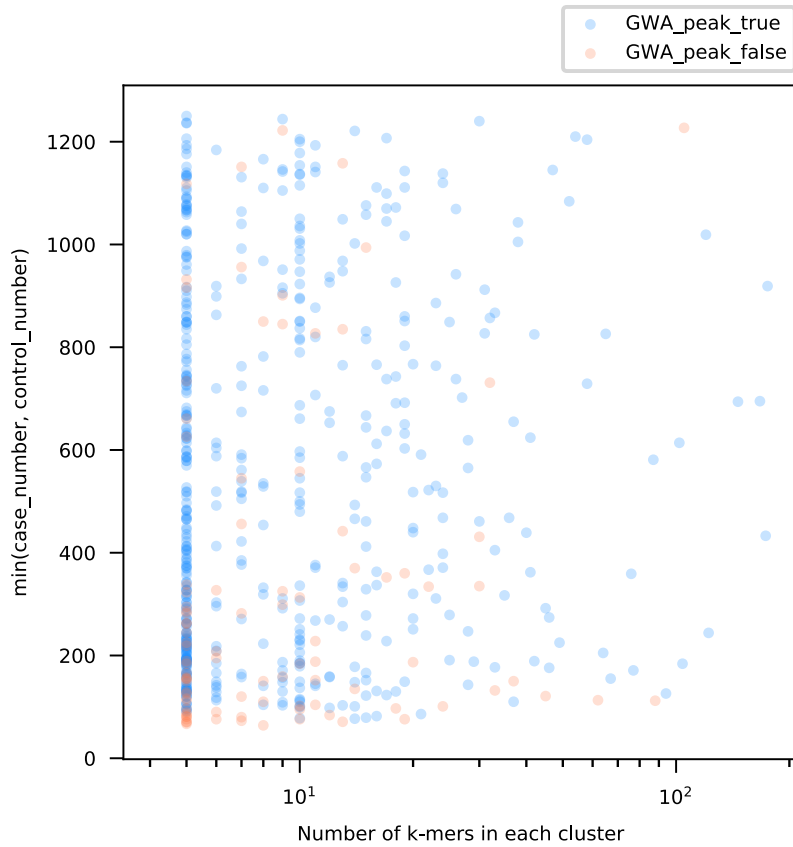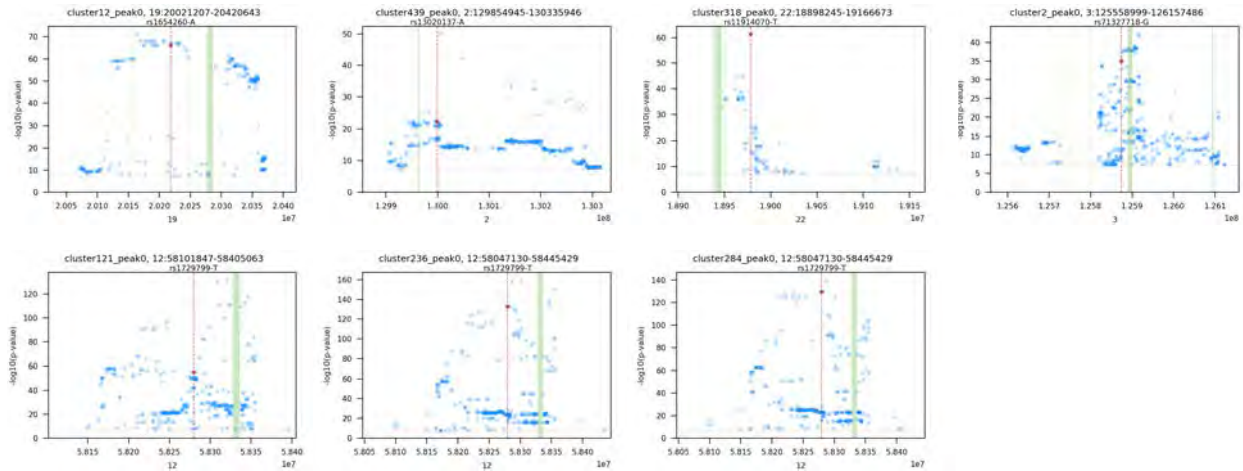1209 used. The tree shown here is the same tree as Figure 3B left panel, except for showing the leaf
1210 names.
1211

**Legend:**
- Exogenous HHV-6
- Known endogenous HHV-6, published sequence
- Known endogenous HHV-6, reconstructed in this paper
- Novel endogenous HHV-6
- Novel integrated HHV-6, solo-DR form
- Node with a bootstrap support ≥ 80%

1219

1220 **Supplementary Figure 10** Phylogenetic tree of HHV-6B DR with leaf names

1221 Phylogenetic trees inferred from DR regions of HHV-6B. The publicly available sequences of

1222 endogenous and exogenous HHV-6B as well as ones reconstructed in the present study were

1223 used. The tree shown here is the same tree as Figure 3C, except for showing the leaf names.
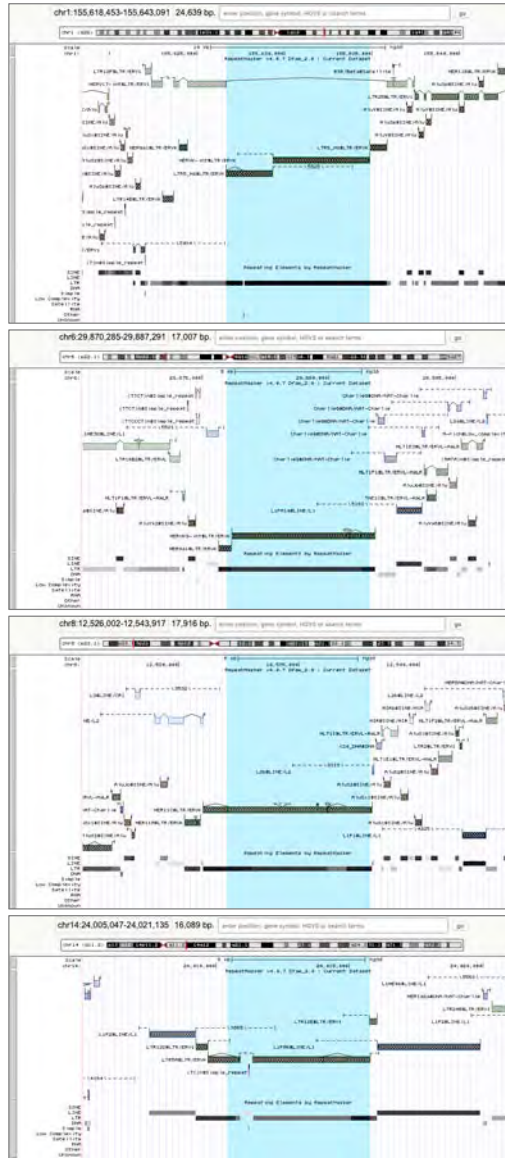
1224

1225

**Supplementary Figure 11** Numbers of *k*-mers in *k*-mer clusters
1227 Scatter plot shows the 597 *k*-mer clusters as dots. X-axis shows the number of *k*-mers in *k*-mer
1228 clusters. Y-axis shows the number of either case or control used for GWA analysis, whichever is
1229 smaller. Blue dots represent *k*-mer clusters with SNVs with association, while red dots show
1230 ones without any association to SNVs.
1231

1232
1233 **Supplementary Figure 12** SNVs in the NHGRI-EBI GWAS catalog overlapping with the HERV-
1234 K *k*-mer LD regions
1235 Manhattan plots showing SNVs in the GWAS catalog overlapping with the indicated *k*-mer
1236 clusters. SNVs with p-value lower than 5e-08 are shown. Green lines show the reference
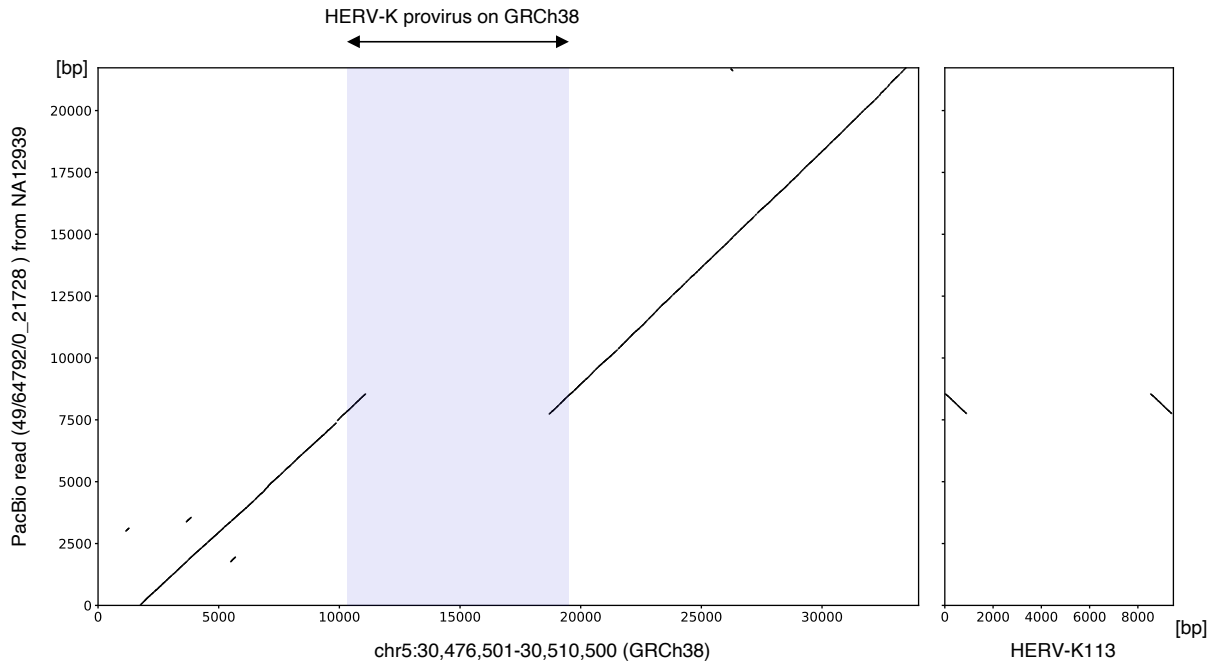1237 HERV-K provirus. Red dots show the lead SNVs listed in the NHGRI-EBI GWAS catalog.
1238
1239

**A**



Absent in at least one sample.

**B**

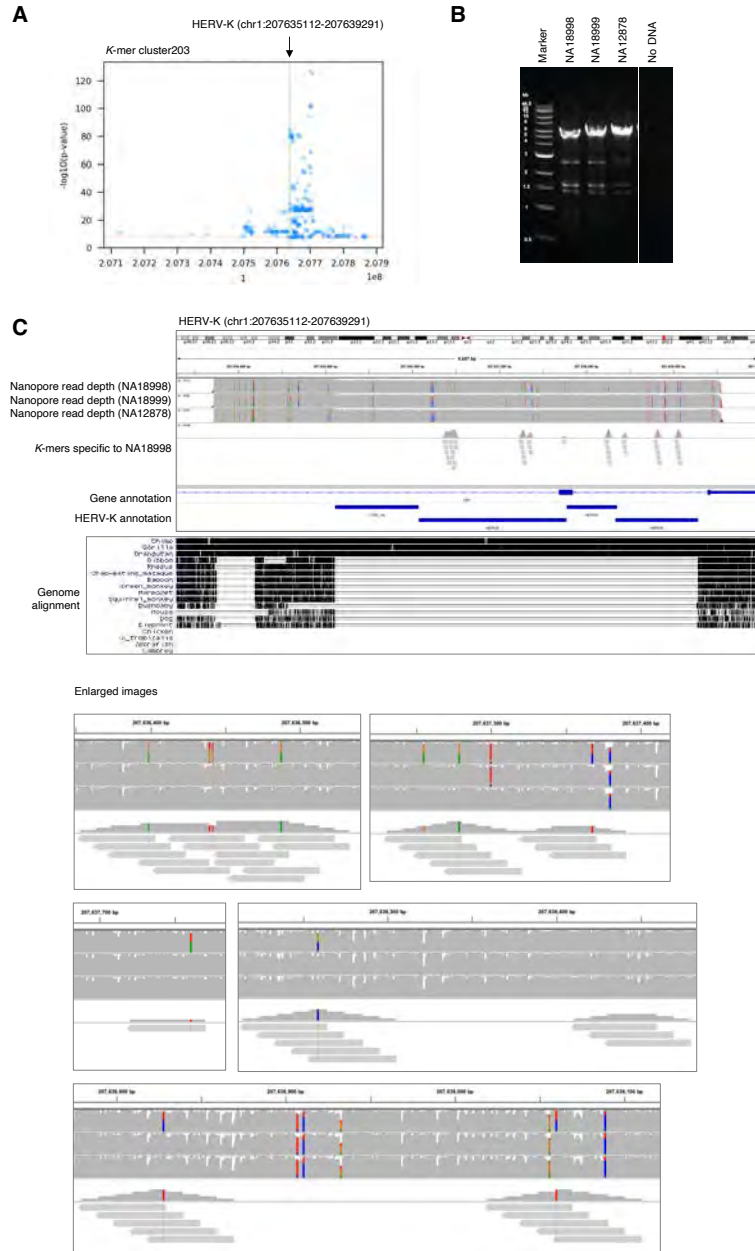| Deletion in Audano et al. | Sample(s) with deletion in Audano et al. | Samples lacking k-mers | K-mer cluster(s) |
|---|---|---|---|
| chr1:155626666-155634878 | NA19434 | NA19434 | cluster76 |
| chr6:29875954-29881622 | HG00268;NA12878;NA19434 | HG00268;NA12878;NA19434 | cluster520 |
| chr8:12531974-12537945 | NA19434 | NA19434 | cluster191;cluster336 |
| chr14:24010410-24015772 | NA12878;NA19434 | NA12878;NA19434 | cluster368 |

1240
1241 **Supplementary Figure 13** Provirus/solo-LTR-type HERV-K polymorphism captured by LDfred
1242 A. UCSC genome browser view showing Four polymorphic HERV-K. Blue regions are detected
1243 as sequence deletions in Audano et al.
1244 B. Cross-reference between deletions exist within HERV-K detected in Audano et al. and *k*-mer
1245 clusters detected by LDfred.
1246

1247
1248    **Supplementary Figure 14** Provirus/solo-LTR-type HERV-K polymorphism captured by LDfred
1249    A PacBio read showing the absence of a HERV-K on chromosome 5 in NA12939. Left dot
1250    matrix shows the alignment between the partial sequence of chromosome 5 and PacBio read
1251    from NA12939 sequenced in Chaisson et al. Right dot matrix shows the alignment between
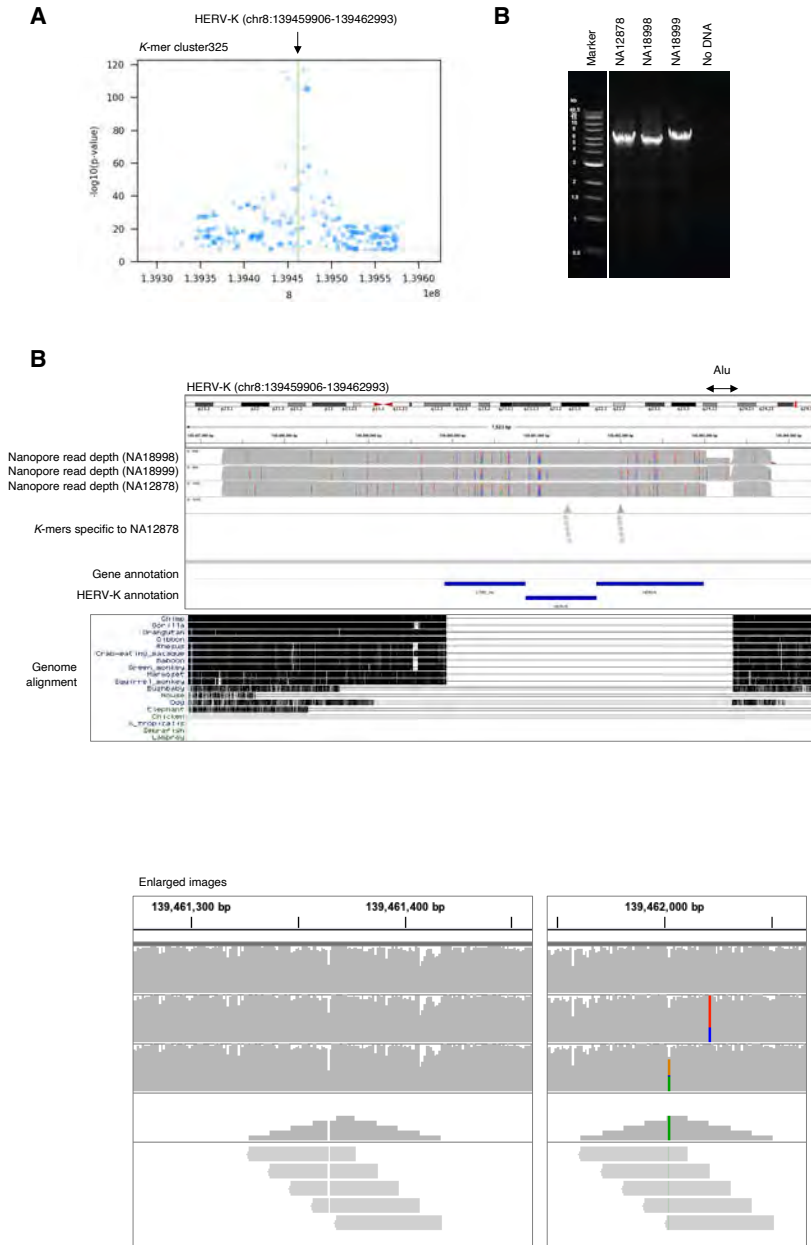1252    HERV-K113 and the PacBio read. Blue region shows the provirus on chromosome 5.
1253

1254

1255 **Supplementary Figure 15** HERV-K SNVs captured by LDfred

1256 A. Manhattan plot showing SNVs associating the *k*-mer cluster203. SNVs with p-value lower

1257 than 5e-08 are shown. Green line shows the reference HERV-K provirus.

1258 B. Amplification of the HERV-K provirus by PCR. HERV-K provirus with adjacent sequence was

1259 amplified and PCR products were separated by gel electrophoresis. DNA extracted from LCLs

1260 originating from NA18998, NA18999, and NA12878 were used as templates.

1261 C. Upper panel: IGV view of long-read sequencing reads mapping to HERV-K. The PCR

1262 amplicons were sequenced using an Oxford Nanopore flongle flow cell and mapped to

1263 GRCh38. *k*-mers in *k*-mer detecting the HERV-K were also mapped to the PCR target regions.

1264 Lower panel: UCSC genome browser view showing the Multiz Alignment of 100 Vertebrates

1265 track.

1266

**A.** HERV-K (chr8:139459906-139462993)

*K*-mer cluster325

**B.**

**B.**

HERV-K (chr8:139459906-139462993)

Alu

Nanopore read depth (NA18998)
Nanopore read depth (NA18999)
Nanopore read depth (NA12878)

*K*-mers specific to NA12878

Gene annotation
HERV-K annotation

Genome alignment

Enlarged images

139,461,300 bp    139,461,400 bp

139,462,000 bp

1267
1268 **Supplementary Figure 16** HERV-K SNVs captured by LDfred
1269 A. Manhattan plot showing SNVs associating the *k*-mer cluster325. SNVs with p-value lower
1270 than 5e-08 are shown. Green line shows the reference HERV-K provirus.
1271 B. Amplification of the HERV-K provirus by PCR. HERV-K provirus with adjacent sequence was
1272 amplified and PCR products were separated by gel electrophoresis. DNA extracted from LCLs
1273 originating from NA18998, NA18999, and NA12878 were used as templates.
1274 C. Upper panel: IGV view of long-read sequencing reads mapping to HERV-K. The PCR
1275 amplicons were sequenced using an Oxford Nanopore flongle flow cell and mapped to
1276 GRCh38. *k*-mers in *k*-mer detecting the HERV-K were also mapped to the PCR target regions.
1277 Lower panel: UCSC genome browser view showing the Multiz Alignment of 100 Vertebrates
1278 track.
1279