1  **The evolutionary history of a gammaretrovirus currently colonizing the mule deer**

2  **genome is marked by extensive recombination**

3

4

5  Lei Yang[1,2], Raunaq Malhotra[3], Rayan Chikhi[2,3,4], Daniel Elleder[1,5], Theodora Kaiser[1],

6  Jesse Rong[3], Paul Medvedev[2,3,4] and Mary Poss[1,2*]

7

8

9  [1] Department of Biology,

10  [2] Center for Comparative Genomics and Bioinformatics,

11  [3] Department of Computer Science and Engineering,

12  [4] Department of Biochemistry and Molecular Biology,

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

14  [5] Institute of Molecular Genetics; The Czech Academy of Sciences; Prague; Czech

15  Republic

16

17

18  *Corresponding author: Mary Poss (maryposs@gmail.com)

19  Current address: Department of Hematology and Oncology, University of Virginia,

20  Charlottesville, VA 22903, USA

21

22

23

24 **Abstract**

25

26 **Background:** All vertebrate genomes have been colonized by retroviruses along their

27 evolutionary trajectory. Although it is clear that endogenous retroviruses (ERVs) can

28 contribute important physiological functions to contemporary hosts, such benefits are

29 attributed to long-term co-evolution of ERV and host. Newly colonized ERVs are thought

30 unlikely to contribute to host genome evolution because germline infections are rare and

31 because the host effectively silences them.  The genomes of several outbred species

32 including mule deer (*Odocoileus hemionus*) are currently being colonized by ERVs,

33 which provides an opportunity to study ERV dynamics at a time when few are fixed.

34 Here we investigate the history of cervid endogenous retrovirus (CrERV) acquisition and

35 expansion in the mule deer genome to determine the potential impact of endogenizing

36 retroviruses on host genomic diversity.

37

38 **Methods:** A mule deer genome was de novo assembled from short and long insert

39 mate pair reads. Scaffolds were further assembled using reference assisted

40 chromosome assembly (RACA) to provide spatial orientation of CrERV insertion sites

41 and to facilitate assembly of CrERV sequences. We applied phylogenetic and

42 coalescent approaches to non-recombinant genomes to determine CrERV evolutionary

43 history, augmenting ancestral divergence estimates with the prevalence of each CrERV

44 locus in a population of mule deer. Recombination history was investigated on partial

45 genome alignments.

46

47    **Results:** The CrERV composition and diversity in the mule deer genome has recently

48    measurably increased by horizontal acquisition of a new retroviruses lineage and

49    because of recombination with existing CrERV. Resulting interlineage recombinants

50    also endogenized and subsequently retrotransposed. CrERV loci are significantly closer

51    to genes than expected if integration were random and gene proximity might explain the

52    recent expansion by retrotransposition of one recombinant CrERV lineage.

53

54    **Conclusions:**  There has been a burst of CrERV integrations during a recent retrovirus

55    epizootic that increased genomic CrERV burden and has resulted in extensive

56    insertional polymorphism in contemporary mule deer genomes. Recombination is a

57    defining feature of CrERV evolutionary dynamics driven by this colonization, increasing

58    CrERV burden and CrERV genetic diversity. These data support that retroviral

59    colonization during an epizootic provides a burst of genomic diversity to the host

60    population.

61

62    **Keywords:** endogenous retrovirus, ERV, recombination, genome diversity, mule deer,

63    Odocoileus hemionus, CrERV

64

65    **<u>Background</u>**

66

67    Retroviruses are unique among viruses in adopting life history strategies that enable

68    them to exist independently as an infectious RNA virus (exogenous retrovirus, XRV) [1]

69    or as an integral component of their host germline (endogenous retrovirus, ERV) [2,3].

3

70 An ERV is the result of a rare infection of a germ cell by an XRV and is maintained in

71 the population by vertical transmission. Germline colonization has been a successful

72 strategy for retroviruses as they comprise up to 10% of most contemporary vertebrate

73 genomes [4]. Over the evolutionary history of the species, ERV composition increases

74 by acquisition of new germ line XRV infections, and through retrotransposition or

75 reinfection of existing ERVs [5–8], which results in clusters of related ERVs. The ERV

76 profile in extant species therefore reflects both the history of retrovirus epizootics and

77 the fate of individual ERVs. Because the acquisition of retroviral DNA in a host genome

78 has the potential to affect host phenotype [9–11], the dynamic interactions among ERVs

79 and host could shape both retrovirus and host biology. However, the evolutionary

80 processes in play near the time of colonization are difficult to discern based on an ERV

81 colonization event that occurred in an ancestral species. A better understanding of both

82 host and virus responses to recent germ line invasion might inform homeostatic

83 changes in ERV-host regulation that are relevant to the pathogenesis of diseases in

84 which ERV involvement has been implicated [12–17].  Fortunately, there is now

85 evidence that retrovirus colonization is occurring in contemporary, albeit often non-

86 model, species [18–20], allowing for investigation of ERV dynamics near the time of

87 colonization. Our goal in this research is to investigate the evolutionary dynamics of the

88 phylogenetically distinct ERV lineages that have sequentially colonized mule deer over

89 the approximate million-year history of this species using the complete genome

90 sequence of a majority of coding ERVs in the context of a draft assembly of a newly

91 sequenced mule deer genome.

92

93    The life history strategy adopted by retroviruses indicates why this virus family has been

94    so successful in colonizing host germline. Retroviral replication requires that the viral

95    RNA genome be converted to DNA and then integrated into the genome of an infected

96    cell [21]. As with many RNA viruses, the virus polymerase enzyme, reverse

97    transcriptase (RT), is error prone, which contributes to a high mutation rate and enables

98    rapid host adaptation. In addition, RT moves between the two RNA copies that

99    comprise a retroviral genome [22]; this process can repair small genomic defects and

100   increases evolutionary rates via recombination if the two strands are not identical.

101   Retroviral DNA is called a provirus and is transported to the nucleus where it integrates

102   into host genomic DNA using a viral integrase enzyme. The provirus represents a newly

103   acquired gene that persists for the life of the cell and is passed to daughter cells, which

104   for XRV are often hematopoietic cells. For a retrovirus infecting a germ cell, all cells in

105   an organism will contain the new retroviral DNA if reproduction of the infected host is

106   successful.

107

108   The retroviral life cycle also demonstrates how ERVs can affect host biology [10,23].

109   ERVs require host transcription factors and RNA polymerases to bind to the retrovirus

110   promoter, called long terminal repeats (LTRs), to produce viral transcripts and the RNA

111   genome. Thus, the viral LTRs compete with host genes for transcription factors and

112   polymerases [24]. A retrovirus encodes at a minimum, genes for the capsid, viral

113   enzymes, and an envelope gene needed for cell entry, which is produced by a sub-

114   genomic mRNA. Hence an ERV also utilizes host-splicing machinery and can alter host

115   gene expression pattern if the site of integration is intronic [25,26]. While XRVs are

116    expressed from small numbers of somatic cells, ERVs are present in all cells and ERV

117    transcripts and proteins can be expressed in any cell type at any stage of host

118    development.  Hosts actively silence the expression of full or partial ERV sequences by

119    epigenetic methods [27,28] and by genes called viral restriction factors [29–33].

120    Because there will be no record of an ERV that causes reproductive failure of the newly

121    colonized host, ERVs in contemporary vertebrates are either effectively controlled by

122    host actions, are nearly neutral in effects on host fitness, or potentially contribute to the

123    overall fitness of the host [34–37].

124

125    The coding portion of a new ERV can be eliminated from the genome through non-

126    allelic homologous recombination (NAHR) between the LTRs, which are identical

127    regions that flank the viral coding portion. A single LTR is left at the site of integration as

128    a consequence of the recombination event and serves as a marker of the original

129    retrovirus integration site [38]. Most ERV integration sites in humans are solo LTRs

130    [39,40]. Because the efficiency of NAHR is highest between identical sequences [41],

131    conversion of a full-length ERV to a solo LTR likely arises early during ERV residency in

132    the genome before sequence identity of the LTR is lost as mutations accrue [42].

133    Because mutations are reported to arise in ERVs at the neutral mutation rate of the host

134    [43], sequence differences between the 5' and 3' LTR of an ERV have been used to

135    approximate the date of integration [44,45].

136

137    Although in humans most ERV colonization events occurred in ancestral species,

138    acquisition of new retroviral elements is an ongoing [46,47] or contemporary [18] event

139    in several animal species. The consequences of a recent ERV acquisition are important

6

140    to the host species because it creates an insertionally polymorphic site; the site is

141    occupied in some individuals but not in others. All ERVs are insertionally polymorphic

142    during the trajectory from initial acquisition to fixation or loss in the genome. Indeed, the

143    HERV-K (human endogenous retrovirus type K) family is insertionally polymorphic in

144    humans [48–51] and HERV-K prevalence at polymorphic sites differ among global

145    populations [52]. Phylogenetic analyses of the ERV population in a genome can inform

146    on the origins of ERV lineages to determine which are actively expanding in the genome

147    and the mutational processes that drive evolution.  These data indicate if expansion is

148    related to the site of integration or a feature of the virus, or both and coupled with

149    information of ERV prevalence at insertionally polymorphic sites, can inform ERV

150    effects on host phenotype.

151

152    To this end, we explored the evolutionary history of the mule deer (*Odocoileus*

153    *hemionus*) ERV (Cervid endogenous retrovirus, CrERV) because we have extensive

154    data for prevalence of CrERV loci in northern US mule deer populations [53] and

155    preliminary data on CrERV sequence variation and colonization history [19,54]. A

156    majority of CrERV loci is insertionally polymorphic in mule deer; 90% of animals shared

157    fewer than 10 of approximately 250 CrERV integrations per genome in one study [53].

158    Further, mule deer appear to have experienced several recent retrovirus epizootics with

159    phylogenetically distinct CrERV and, because none of the CrERV loci occupied in mule

160    deer are found in the sister species, white-tailed deer (*Odocoileus virginianus*) [19], all

161    endogenization events have likely occurred since the split of these sister taxa. Based on

162    the phylogeny of several CrERV identified in the mule deer genome, at least four

7

163  distinct epizootics resulted in germ line colonization [54]. A full-length retrovirus

164  representing the youngest of the CrERV lineages was recovered by co-culture on

165  human cells, indicating that some of these CrERV are still capable of infection [55]. In

166  this study, we expand on these preliminary data by sequencing a mule deer genome

167  and conducting phylogenetic analyses on a majority of reconstructed CrERV genomes.

168  Our results demonstrate that expression and recombination of recently acquired CrERV

169  with older CrERV have increased CrERV burden and diversity and consequently have

170  increased contemporary mule deer genome diversity.

171

172  **Results**

173

174  Establishing a draft mule deer reference genome to study CrERV evolution and

175  integration site preference

176  We developed a draft assembly of a mule deer genome from animal MT273 in order to

177  determine the sequence at each CrERV locus for phylogenetic analyses and to

178  investigate the effect of CrERV lineage or age on integration site preference. ERV

179  sequences are available in any genome sequencing data because a retrovirus

180  integrates a DNA copy into the host genome. However, there is extensive homology

181  among the most recently integrated ERVs making them difficult to assemble and

182  causing scaffolds to break at the site of an ERV insertion [56].  We assembled scaffolds

183  using a combination of high coverage Illumina short read whole genome sequencing

184  (WGS) and long insert mate pair sequencing. Our *de novo* assembly yielded an ~3.31

185  Gbp draft genome with an N50 of 156 Kbp (Table S1), which is comparable to the 3.33

8

186 Gbp (c-value of 3.41 pg) experimentally-determined genome size of reindeer (*Rangifer*

187 *tarandus*) [57,58].

188

189 Approximately half of CrERV loci are located at the ends of scaffolds based on mapping

190 our previously published junction fragment sequences [53], which is consistent with the

191 fact that repetitive elements such as ERVs break scaffolds. To determine the sequence

192 of these CrERVs and the genome context in which they are found, we developed a

193 higher order assembly using reference assisted chromosome assembly (RACA) [59].

194 RACA further scaffolds our *de novo* mule deer assembly into 'chromosome fragments'

195 by identifying synteny blocks among the mule deer scaffolds, the reference species

196 genome (cow), and the outgroup genome (human) (Figure 1A). We created a series of

197 RACA assemblies based on scaffold length to make efficient use of all data (Table S1).

198 RACA150K takes all scaffolds greater than 150,000 bp as input and yielded 41

199 chromosome fragments, 35 of which are greater than 1.5 Mbp; this is consistent with

200 the known mule deer karyotype of 2n=70 [60]. However, RACA150K only incorporates

201 48% of the total assembled sequences (1.59 Gbp) because of the scaffold size

202 constraint. In contrast, RACA10K uses all scaffolds 10,000 bp or longer and increases

203 the assembly size to 2.37 Gbp (~72% of total assembly) but contains 658 chromosome

204 fragments (Table S1). The majority of scaffolds that cannot be incorporated into a

205 RACA assembly are close to the ends of alignment chains (File S1, section 1a). Most

206 sequences not represented in any assemblies were repeats based on *k-mer* analyses

207 (File S1, section 1a and Figure S1).

208

209 Some scaffolds were excluded from the RACA assemblies, presumably because there

210 is no synteny between cow and human for these sequences. We oriented these

211 scaffolds using the cow-mule deer and sheep-mule deer alignments (RACA+, Table

212 S2). Approximately 124 Mbp of sequence (~4% of total assembly) is in scaffolds larger

213 than 10kb but cannot be placed in RACA10K, nearly all of which can be found on the

214 mule deer-cow alignment chain and the mule deer-sheep (oviAri3) alignment chain (123

215 Mbp in each chain). Because there is overlap between these alignments, only ~1 Mbp is

216 specific to cow and ~1 Mbp is specific to sheep. Therefore, RACA+ incorporated all but

217 69 scaffolds that are greater than 10 kbp, which consisted of 1.17 Mbp of sequence

218 (~0.04% of total scaffold size of the assembly) and yields an assembly size of 2.49 Gbp

219 (Table S1).

220

221 To enable the investigation of CrERV integration site preference relative to host genes,

222 we annotated the mule deer scaffolds. We used Maker2 [61,62] for the annotation,

223 which detects candidate genes based on RNA sequencing data and protein homology

224 to any of the three reference genomes: human, cow and sheep. After four Maker

225 iterations, 21,598 genes with an AED (annotation edit distance) [61] of less than 0.8

226 were annotated (Table S3). Approximately 92% of genes are found on RACA150K

227 scaffolds and 95% of genes are represented in RACA10K scaffolds.

228

229 <u>Establishing the location and sequence at CrERV loci</u>

230 Several lines of evidence suggest that most CrERVs are missing from the assemblies.

231 Only three CrERVs with coding potential were assembled by the *de novo* assembly.

10

232    The *k-mer* based analysis shows that less than 9.62% of all LTR repeat elements are in

233    the assemblies (Table S4). The CrERV-host junction fragments previously sequenced

234    [53] support that CrERV loci are near scaffold ends or in long stretches of 'N's.

235    Therefore, we took advantage of the different chromosome fragments generated by

236    RACA10K, RACA150K and RACA+ and the long insert mate pair sequencing data to

237    reconstruct CrERVs at each locus (Figure 1B). We identified 252 CrERV loci in the

238    MT273 genome, which is consistent with our estimates of an average of 240 CrERV loci

239    per mule deer by quantitative PCR [19] and 262 CrERV loci in animal MT273 by

240    junction fragment analysis [53]. The majority of CrERV loci (206/252) contains CrERVs

241    with some coding capacity and 46 are solo LTRs. Of the 206 CrERVs containing genes,

242    164 were sufficiently complete to allow phylogenetic analysis on the entire genome or, if

243    a deletion was present, on a subset of viral genes; at 42 loci we were unable to obtain

244    sufficient lengths of high-quality data for further analyses.

245

246    Evolutionary history of CrERV

247    We previously showed that mule deer genomes have been colonized multiple times

248    since the ancestral split with white tailed deer approximately one million years ago

249    (MYA) [54] because none of the CrERV integration sites are found in white-tailed deer.

250    To better resolve the colonization history, we conducted a coalescent analysis based on

251    an alignment spanning position 1,477-8,633 bp (omitting a portion of *env*) of the CrERV

252    genome (GenBank: JN592050) using 34 reconstructed CrERV sequences with high

253    quality data that had no signature of recombination and that were representative of the

254    phylogeny in a larger data set (Figure 2). The majority of the *env* gene, which has

255    distinct variable and conserved region [63], was manually blocked because of alignment

256    difficulties (6,923-7,503 bp by JN592050 coordinates; see Figure 2, right panel for

257    diagram of *env* variable regions and Table S5, column C for *env* structure of each

258    CrERV). This tree shows four well-supported CrERV lineages, each diverged from a

259    common ancestor at several points since the split of mule deer and white-tailed deer.

260    Although *env* sequence is not included in the phylogenetic analysis, CrERV assigned to

261    each of the four identified lineages share the same distinct *env* variable region structure

262    of insertions and deletions, which define the receptor-binding domain of the envelope

263    protein (Figure 2, right panel).

264

265    Lineage A CrERVs are the youngest ERV family in mule deer. Our estimates indicate

266    that Lineage A colonization has occurred over the last 300 thousand years to the

267    present (Figure 2; Table S6, node f, 95% high posterior density (HPD) interval 110-470

268    thousand years ago (KYA)) and is represented by three well-supported CrERV

269    subgroups evolving over this time frame. All have a complete open reading frame (ORF)

270    in *env* and likely represent a recent retrovirus epizootic. An infectious virus recovered by

271    co-culture belongs to this lineage [19]. Lineage A represents 30% of all CrERV sampled

272    from MT273 (Table S5). Our age estimates for each subgroup of Lineage A CrERV are

273    consistent with their prevalence in populations of mule deer in the Northern Rocky

274    Mountain ecosystem (Figure 2); [64]. For example, S29996 and S10113 are estimated

275    to derive from an older Lineage A CrERV subgroup and occur in our sampled mule deer

276    at higher prevalence than those estimated to have entered the genome more recently

277    (see S22897 and S111665, Figure 2).

278

279    Lineage B CrERV shared a common ancestor with Lineage A approximately 300 KYA

280    (node i, Figure 2). Lineage B CrERVs have a short insertion in the 5' portion of *env*

281    followed by a deletion that removes most of the *env* surface unit (SU) relative to

282    Lineage A *env*. Because our coalescent analysis does not include *env* sequence, these

283    results suggest that two phylogenetically distinct XRV with different envelope proteins

284    were circulating about the same time in mule deer populations. Lineage B CrERV

285    represent 32% of sampled viruses from our sequenced genome (Table S5). Like

286    Lineage A, the prevalence of CrERV from Lineage B among mule deer in the northern

287    Rockies region is low, reflecting their more recent colonization of the mule deer

288    genome. Indeed, six Lineage B CrERVs were identified only in MT273, while only one

289    Lineage A CrERV is found only in MT273 (Table S5), which could be indicative of a

290    recent expansion of some Lineage B CrERV. Of note, there are two related groups of

291    CrERV affiliated with Lineage B (Lineage B1 and B2, Table S5). One shares the short 5'

292    insertion in *env* but has a full-length *env* with an additional short insertion relative to the

293    *env* of Lineage A CrERV (Lineage B1, Figure 2). CrERV with this *env* configuration

294    represent 9% of coding CrERV in MT273. Because the prevalence of Lineage B1 is

295    high in mule deer, this group could represent the ancestral state for Lineage B CrERVs.

296    The second group has a unique *env* not found in any other CrERV lineages (Lineage

297    B2, Figure 2, node k; S16113 and S6404). We are unable to estimate the prevalence of

298    this unusual *env* containing CrERV in mule deer because the host junction fragments

299    are not represented in our draft mule deer assembly. It is possible that these viruses

300    represent a cross-species infection and it would be interesting to determine if

13

301    representatives of Lineage B2 are found in the genomes of other species that occupied

302    the ecosystem in the past.

303

304    Our coalescent estimates indicate that Lineage C CrERV emerged about 500 KYA

305    (Table S6). Several members of this lineage are found in all mule deer sampled (Figure

306    2; Table S5), consistent with a longer residence in the genome. There is a 59 bp

307    insertion (C) and 362 bp deletion (E) in *env* (Figure 2; Table S5) compared to the full

308    length *env* of Lineage A; none have an intact *env* ORF. Despite evidence that Lineage

309    C is an older CrERV, the approximately 13% of identified CrERV in MT273 belonging to

310    this lineage share a common ancestor ~50 KYA (95% HPD: 16-116 KYA, Table S6).

311    These data are consistent with a recent expansion of a long-term resident CrERV.

312

313    The first representatives of the CrERV family still identifiable in mule deer colonized

314    shortly after their split from white-tailed deer, approximately one MYA [19].  Lineage D

315    CrERVs comprise 12% of reconstructed CrERV in MT273 and appear to be near

316    fixation. Indeed, all mule deer in a larger survey of over 250 deer had CrERV S26536,

317    which is not found in white-tailed deer [54].  This lineage shares an *env* insertion with

318    Lineage C but lacks the deletion, which removes the transmembrane region of *env*.

319

320    These data expand our previous findings that over the approximately one million year

321    history of mule deer, the mule deer genome has been colonized at least four times by

322    phylogenetically distinct CrERVs; this likely reflects several retroviral epizootics each

323    characterized by a unique *env* structure. The two lineages responsible for most recent

14

324    endogenization events comprise 62% of sampled CrERV. In addition, these data

325    capture the evolutionary processes acting on the *env* gene of exogenous retroviruses,

326    which are characterized by gain or loss of variable regions of this important viral protein.

327

328    <u>Recombination among CrERV lineages</u>

329    Our coalescent estimates (Figure 2) indicate that two phylogenetically distinct CrERV

330    lineages have been expanding in contemporary mule deer genomes over the last

331    100,000 years. Both lineages have been actively colonizing contemporary mule deer

332    genomes based on divergence estimates, which include zero. While CrERVs

333    represented by Lineage A are capable of infection [19], all Lineage B CrERVs have an

334    identical deletion of the SU portion of *env* and should not be able to spread by

335    reinfecting germ cells. However, the mule deer genome is comprised of approximately

336    equal percentages of Lineage B and Lineage A CrERVs so we considered two modes

337    by which defective Lineage B CrERVs could expand in the genome at a similar rate with

338    Lineage A. Firstly, ERVs that have lost *env* are proposed to preferentially expand by

339    retrotransposition [65] because a functional envelope is not necessary for intracellular

340    replication. Secondly, we consider that Lineage B CrERVs could increase in the

341    genome by infection if the co-circulating Lineage A group provided a functional

342    envelope protein, a process called complementation [5,66]. This latter mechanism

343    requires that a member of each CrERV lineage be transcriptionally active at the same

344    time in the same cell, and that intact proteins from the 'helper' genome be used to

345    assemble a particle with a functional envelope for reinfection. If two different CrERV loci

346    are expressed in the same cell, both genomes could be co-packaged in the particle.

15

347    Because the reverse transcriptase moves between the two RNA genomes as first

348    strand DNA synthesis proceeds, evidence of inter-lineage recombination would support

349    that the molecular components necessary for complementation were in place. We

350    assessed Lineage B CrERV for recombination with Lineage A to determine if coincident

351    expression of the RNA genomes of these two lineages could explain the expansion by

352    infection through complementation of the *env*-less Lineage B CrERV.

353

354    There is good support for recombination between Lineage A and B in a region spanning

355    a portion of *pol* to the beginning of the variable region in *env* (4,422-7,076 based on

356    coordinates of JN592050). In this region, several CrERV that we provisionally classified

357    as Lineage B because they carried the prototypical *env* deletion of SU form a

358    monophyletic group that is affiliated with Lineage A CrERV (Figure 3, upper collapsed

359    clade containing orange diamonds). These Lineage B recombinants all share the same

360    recombination breakpoint just 5' of the characteristic short insertion for these viruses

361    (Figure S2, indicated by "**"; Table S7). In addition, several other CrERVs with Lineage

362    B *env* branch between lineages A and B, indicating that the recombination breakpoints

363    fall within the region assessed (Figure S2). Indeed, the breakpoint in a group of three

364    CrERV is at position 6630 based on coordinates of JN592050, which is near the

365    predicted splice site for *env* at position 6591 [19]; this confers an additional 500 bp of

366    the Lineage B *env* on these viruses (Figure S2) resulting in their observed phylogenetic

367    placement. Because recombination between the two retroviral RNA genomes occurs

368    during reverse transcription, our data indicate that both Lineage A and B CrERVs were

369    expressed and assembled in a particle containing a copy of each genome.  A functional

16

370   envelope from Lineage A would therefore have been available for infection. These data

371   support our premise that complementation with a replication competent Lineage A

372   CrERV or CrXRV (cervid exogenous gammaretrovirus, an exogenous version of

373   CrERV) contributes to the 32% prevalence of *env*-deleted Lineage B CrERV in the

374   genome. It is notable that recent retrotransposition of the lineage A-B recombinant

375   CrERVs likely occurred because they are nearly identical and the branches supporting

376   them are short (Figure 3, orange diamonds in the Lineage A type *env* cluster).

377

378   There is additional data to support the transcriptional activity of a Lineage B CrERV,

379   which is requisite for recombination with an infectious Lineage A CrERV or for

380   retrotransposition. We identified a non-recombinant Lineage B CrERV (S24870 in Table

381   S5) with extensive G to A changes (184 changes) compared to other members of this

382   monophyletic group. These data are indicative of a cytidine deaminase acting on the

383   single stranded DNA produced during reverse transcription [67].

384

385   Lineage C CrERV are enigmatic because based on full length sequences lacking a

386   signature of recombination it diverged around 500KYA (Figure 2) but all extant

387   members of this group diverged recently. From Figure 3, it is evident that over the

388   region of *pol* assessed, CrERVs containing the Lineage C *env* cluster with an older

389   Lineage A subgroup. Given that the *env* of Lineage C CrERV shares sequence

390   homology and an insertion with that of the oldest Lineage D, it is likely that Lineage C is

391   in fact the result of recombination between an early member of Lineage A and a relative

392   of a Lineage D CrERV. Many, but not all, Lineage C CrERVs are found at high

393    prevalence in the mule deer population (Figure 2; Table S5), supporting that the initial

394    recombination event occurred early during the Lineage A colonization. Our identification

395    of Lineage C as derived from a non-recombinant CrXRV is therefore incorrect. Instead,

396    Lineage C CrERVs are derived from a CrERV or CrXRV that is not currently

397    represented in mule deer genomes either because it was lost or it never endogenized.

398    Fourteen of the twenty-two CrERV in Lineage C have multiple signatures of

399    recombination predominantly with Lineage A CrERV. The expansion of a subset of

400    Lineage C as a monophyletic group approximately 50 KYA (Figure 2; Table S6)

401    suggests that like some members of Lineage B, CrERVs generated by recombination

402    with Lineage A have recently retrotransposed.

403

404    <u>Genomic distribution of CrERV lineages</u>

405    Of the 164 CrERV that we reconstructed from MT273, only 12 can be detected in all

406    mule deer that we have sampled [53,54] (Table S5). This means that the majority of

407    CrERV loci in mule deer are insertionally polymorphic; not all animals will have a CrERV

408    occupying a given locus. ERVs can impact genome function in multiple ways but the

409    best documented is by altering host gene regulation, which occurs if the integration site

410    is near a host gene [68]. Thus, we investigated the spatial distribution of CrERV loci

411    relative to host genes to determine the potential of either fixed or polymorphic CrERV to

412    impact gene expression, which could affect host phenotype.

413

414    The actual distance between genes is likely to be unreliable in our assembly because

415    most high copy number repeats are missing in the mule deer assembly (Figure S1,

18

416    Table S4, section 1a of File S1). To investigate potential problems determining the

417    spatial distribution of CrERV insertions imposed by using a draft assembly, we

418    simulated the distribution of retrovirus insertions (File S1, section 2l) in mule deer

419    (scaffold N50=156 Kbp) and the genomes of cow (Btau7, scaffold N50=2.60 Mbp) and

420    human (hg19, scaffold N50=46.4 Mbp). The mean distance between insertion and the

421    closest gene for all simulation replicates (Figure S3) is significantly higher in the cow

422    and human (Mann-Whitney U test $p < 2.2 \times 10^{-16}$ for any pair of comparison among the

423    three species). Therefore, we determined if the number of CrERV loci observed to be

424    within 20Kbp of a gene differed from that expected if the distribution was random. There

425    are significantly more observed insertions that fall within 20 Kbp of the translation start

426    site of a gene than occur randomly (Figure 4A). In contrast, intronic CrERV insertions

427    are significantly less than expected based on our simulations (Figure 4B). Among

428    Lineage A CrERVs, only a single sub-lineage (CrERVs that are associated to node 'a' in

429    Figure 2) are found in closer proximity to genes (bold font in Column G of Table S5)

430    than expected if integrations are random (Fisher's exact test $p = 0.002891$). We also

431    investigated whether any of the recombinant CrERV with a signature of recent

432    expansion was integrated within 20 Kbp of a gene. Two of the three recombinant

433    clusters (Figure 3) contain members that are close to a gene (Table S5, bold font in

434    column G). In particular, Lineage A/B recombinant CrERV S10 is 494 bp from the start

435    of a gene. Remarkably, four Lineage C CrERVs with the typical *env* sequence are within

436    20 Kbp of a gene (Table S5, bold font in column G). Our data indicate that integration

437    site preference overall favors proximity to genes but that this is not reflected in all

438    lineages. In particular, the history of Lineage C CrERV suggests they could have

439    acquired a different integration site preference through recombination that facilitated

440    recent genome expansion.

441

**Discussion**

443

444    The wealth of data on human ERVs (HERVs) provides the contemporary status of

445    events that initiated early in hominid evolution. Potential impacts of an ERV near the

446    time of colonization on a host population is thought to be minimal because infection of

447    host germ line by an XRV is a rare event and ERVs that affect host fitness are quickly

448    lost. Potentially deleterious ERVs that are not lost due to reproductive failure can be

449    removed by recombination leaving a solo LTR at the integration site or can suffer

450    degradation presumably because there is no benefit to retain function at these loci;

451    most HERVs are represented by these two states. In addition, humans and other

452    vertebrate hosts have invested extensive genomic resources [4,9,69] to control the

453    expression of ERVs that are maintained. The dynamics between host and ERV are

454    described as an evolutionary arms race [70,71]. This narrative may underrepresent any

455    contributions of ERVs to fitness as they were establishing in a newly colonized host

456    population. Because there are now several species identified to be at different points

457    along the evolutionary scale initiated by the horizontal acquisition of retroviral DNA it is

458    possible to investigate dynamics of ERV that are not yet fixed in a contemporary

459    species. Considering the numerous mechanisms by which newly integrated retroviral

460    DNA affect host biology, such as by introducing new hotspots for recombination [72],

461    altering host gene regulation [68,73,74], and providing retroviral transcripts and proteins

462    for host exaptation [75–78], colonizing ERVs could make a substantive contribution to

463    species' evolution. Our research on the evolutionary dynamics of mule deer CrERV

464    demonstrates that genomic CrERV content and diversity increased significantly during a

465    recent retroviral epizootic due to acquisition of new XRV and from endogenization and

466    retrotransposition of recombinants generated between recent and older CrERVs. These

467    data suggest that CrERV provide a pulse of genetic diversity, which could impact this

468    species' evolutionary trajectory.

469

470    Our analyses of CrERV dynamics in mule deer are based on the sequence of a majority

471    of coding CrERVs in MT273. Of the 252 CrERV loci identified in the MT273 assembly,

472    we were able to reconstruct CrERV sequences from long insert mate pair and Sanger

473    sequencing to use for phylogenetic analysis at 164 sites; 46 sites were solo LTR and 42

474    were occupied by CrERV retaining some coding capacity. We complimented

475    phylogenetic analyses with our previous data on the frequency of each CrERV locus

476    identified in MT273 in a population of mule deer in the northern Rocky Mountain

477    ecosystem [53,64]. In addition, we incorporated information on the variable structure of

478    the retroviral envelope gene, *env*, which is characteristic of retrovirus lineages but was

479    excluded from phylogenetic analyses. The variable regions of retroviral *env* result from

480    balancing its role in receptor-mediated, cell specific infection while evading host

481    adaptive immune response [79,80]. Despite excluding most of *env* from our phylogenic

482    analysis because of alignment problems, each of the lineages we identified has a

483    similar distinct *env* structure, as is well known for infectious retroviruses. By integrating

484    population frequency, coalescent estimation, and the unique structural features of *env*

21

485 we provide an integrated approach to explore the evolutionary dynamics of an

486 endogenizing ERV.

487

488 The most recent CrERV epizootic recorded by germline infection was coincident with

489 the last glacial period, which ended about 12 KYA.  The retroviruses that endogenized

490 during this epizootic belong to Lineage A, have open reading frames for all genes and

491 have been recovered by co-culture as infectious viruses [55]. There are several sub-

492 lineages within Lineage A, which likely reflect the evolutionary history of CrXRV

493 contributing to germline infections over this time period. Lineage A retroviruses

494 constitute approximately one third of all retroviral integrations in the genome. Only four

495 of the fifty Lineage A CrERV that we were able to reconstruct did not have a full length

496 *env*. An important implication of this result is that over the most recent approximately

497 100,000 years of the evolution of this species, the mule deer genome acquired up to

498 half a megabase of new DNA, which introduced new regulatory elements with promoter

499 and enhancer capability, new splice sites, and sites for genome rearrangements. Thus,

500 there is a potential to impact host fitness through altered host gene regulation even if

501 host control mechanisms suppress retroviral gene expression. None of the Lineage A

502 CrERV is fixed in mule deer populations (Table S5, column F) so any effect of CrERV

503 on the host will not be experienced equally in all animals. However, none of the Lineage

504 A CrERV is found only in M273 indicating that the burst of new CrERV DNA acquired

505 during the most recent epizootic has not caused reproductive failure among mule deer.

506 These data demonstrate that in mule deer, a substantial accrual of retroviral DNA in the

507    genome can occur over short time spans in an epizootic and could impose differential

508    fitness in the newly colonized population.

509

510    Lineage A CrERV has an open reading frame for *env* but Lineages B-D do not. Lineage

511    B CrERVs are intriguing in this regard because they also constitute approximately a

512    third of the CrERV in the genome. Yet all have identical deletions of the extracellular

513    portion of *env*, which should render them incapable of genome expansion by reinfection.

514    ERV that have deleted *env* are reportedly better able to expand by retrotransposition

515    [65], which could account for the prevalence of Lineage B. However, because we have

516    evidence for recent expansion of Lineage A and B recombinants, we considered an

517    alternative explanation; that *env*-deficient Lineage B CrERV was complemented with an

518    intact Lineage A CrERV envelope glycoprotein allowing for germline infection.

519    Complementation is not uncommon between XRV and ERV [81,82], is well established

520    for murine Intracisternal A-type Particle (IAP) [83] and has been reported for ERV

521    expansion in canids [84]. Complementation requires that two different retroviruses are

522    co-expressed in the same cell [85]. During viral assembly functional genes supplied by

523    either virus are incorporated into the virus particle and either or both retroviral genomes

524    can be packaged. Because the retroviral polymerase uses both strands of RNA during

525    reverse transcription to yield proviral DNA, a recombinant can arise if the two co-

526    packaged RNA strands are not identical. We investigated the possibility of

527    complementation by searching for Lineage A-B recombinants. Our data show that

528    Linage A and B recombination has occurred several times. A group of CrERV that

529    encode a Lineage B *env* cluster with Lineage A CrERV in a phylogeny based on a

23

530    partial genome alignment (JN592050: 4422-7076bp). The recombinant breakpoint

531    within this monophyletic group is identical, suggesting that the inter-lineage recombinant

532    subsequently expanded by retrotransposition. Notably, two of the CrERV in this

533    recombinant cluster were only found in M273, indicating that retrotransposition was a

534    recent event. There are other clusters of CrERV with Lineage B *env* affiliated with

535    Lineage A CrERV that have different breakpoints in this partial phylogeny.

536    Recombination between an XRV and ERV is also a well-documented property of

537    retroviruses [86–88]. However, the recombinant retroviruses that result are typically

538    identified because they are XRV and often associated with disease or a host switch.

539    Our data indicate that multiple recombination events between Lineage A and B CrERV

540    have been recorded in germline; this in itself is remarkable given that endogenization is

541    a rare event. Thus both the burden of CrERV integrations and the sequence diversity of

542    CrERV in the mule deer genome increase concomitant with a retrovirus epizootic by

543    CrERV inter-lineage recombination.

544

545    Recombination is a dominant feature of CrERV dynamics and is also displayed in the

546    evolutionary history of Lineage C CrERV. Our phylogenetic analysis places the ancestor

547    of Lineage C CrERV at 500 KYA and indeed, Lineage C and Lineage D, which is

548    estimated to be the first CrERV to colonize mule deer after splitting from white-tailed

549    deer [19,54], share many features in *env* that are distinct from those of Lineage A and

550    B. Consistent with a long-term residency in the genome, many Lineage C CrERV are

551    found in most or all mule deer surveyed. A recent expansion of a CrERV that has been

552    quiescent in the genome since endogenizing could explain the estimated 50 KYA time

24

553    to most recent common ancestor of extant members of this lineage. Although this

554    scenario is consistent with the paradigm that a single XRV colonized the genome and

555    recently expanded by retrotransposition, our analysis shows that all Lineage C CrERV

556    are recombinants of a Lineage A CrERV and a CrERV not recorded in or lost from

557    contemporary mule deer genomes. Hence the resulting monophyletic lineage does not

558    arise from retrotransposition of an ancient colonizing XRV. Rather, as is the case with

559    Lineage B CrERV, recombination between an older CrERV and either a Lineage A

560    CrXRV or CrERV occurred, infected germline, and recently expanded by

561    retrotransposition. It is noteworthy that all retrotransposition events detectable in our

562    data involve recombinant CrERV. Further, recombination often leads to duplications and

563    deletions in the retroviral genome, therefore some of the deletions we document in

564    Lineages B-D are not a consequence of slow degradation in the genome but rather are

565    due to reverse transcription and as was recently reported for Koala retrovirus [88]**.**

566

567    These data highlight that expansion of CrERV diversity and genomic burden has

568    occurred in the recent evolutionary history of mule deer by new acquisitions,

569    complementation, and pulses of retrotransposition of inter-lineage recombinants.

570    Indeed, several of the recombinant Lineage C CrERVs that have expanded by

571    retrotransposition are within 20kbp of a gene raising the question as to whether there is

572    a fitness effect at these loci that is in balance with continued expression of the

573    retrovirus. It is remarkable that so many of the events marking the dynamics of

574    retrovirus endogenization are preserved in contemporary mule deer genomes. Given

575    that germline infection is a rare event, it is likely that the dynamics we describe here

576    also resulted in infection of somatic cells. It is worthwhile to consider the potential for

577    ERVs in other species, in particular in humans where several HERVs are expressed, to

578    generate novel antigens through recombination or disruptive somatic integrations that

579    could contribute to disease states.

580

581    **Methods**

582

583    <u>Sequencing</u>

584    Whole genome sequencing (WGS) was performed for a male mule deer, MT273, at

585    ~30x depth using the library of ~260 bp insert size, ~10x using the library of ~1,400-

586    5,000 bp insert size and ~30x using the library of ~6,600 bp insert size. 3' CrERV-host

587    junction fragment sequencing was performed as described by Bao *et al.* [53]. 5' CrERV-

588    host junction fragment sequencing was performed on the Roche 454 platform, with a

589    target size of ~500bp containing up to 380 bp of CrERV LTR.

590

591    <u>Assembly and mapping</u>

592    The draft assembly of MT273 was generated using SOAPdenovo2 [89] (File S1, section

593    2a). WGS data were then mapped back to the assembly using the default setting of bwa

594    mem [90] for further use in RACA and CrERV reconstruction. RNA-seq data was

595    mapped to the WGS scaffolds using the default setting of tophat [91,92]. 3' junction

596    fragments were clustered as described in Bao *et al.* [53]. 3' junction fragment clusters

597    and 5' junction fragment reads were mapped to the WGS assembly using the default

598    setting of blat [93]. A perl script was used to filter for the clusters or reads whose host

26

599     side of the fragment maps to the host at its full length and high identity. 5' junction

600     fragments were then clustered using the default setting of bedtools merge.

601

602     RACA

603     Synteny based scaffolding using RACA was performed based on the genome alignment

604     between the mule deer WGS assembly, a reference genome (cow, bosTau7 or Btau7),

605     and an outgroup genome (hg19). Genome alignments were performed with lastz [94]

606     under the setting of '--notransition --step=20', and then processed using the UCSC

607     axtChain and chainNet tools. The mule deer-cow-human phylogeny was derived from

608     Bininda-Emonds *et al.* [95] using the 'ape' package of R.

609

610     CrERV sequence reconstruction

611     CrERV locations and sequences were retrieved based on junction fragment and long

612     insert mate pair WGS data. The long insert mate pair WGS reads were mapped to the

613     reference CrERV (GenBank: JN592050) using bwa mem. Mates of reads that mapped

614     to the reference CrERV were extracted and then mapped to the WGS assembly using

615     bwa mem. Mates mapped to the WGS assembly were then clustered using the 'cluster'

616     function of bedtools. Anchoring mate pair clusters on both sides of the insertion site

617     were complemented by junction fragments to localize CrERVs. Based on the RACA

618     data, CrERVs that sit between scaffolds were also retrieved in this manner. CrERV

619     reads were then assigned to their corresponding cluster and were assembled using

620     SeqMan (DNASTAR). Sanger sequencing was performed to complement key regions

27

621 used in CrERV evolutionary analyses. All reconstructed CrERV sequences used in the

622 phylogenetic analyses are included in File S2 in fasta format.

623

624 <u>CrERV evolution analyses</u>

625 CrERV sequences of interest were initially aligned using the default setting of muscle

626 [96], manually trimmed for the region of interest, and then re-aligned using the default

627 setting of Prank [97]. Lineage-specific regions are manually curated to form lineage-

628 specific blocks. Models for phylogeny were selected by AICc (Akaike Information

629 Criterion with correction) using jModelTest [98]. Coalescent analysis and associated

630 phylogeny (Figure 2) was generated using BEAST2 [99]. In the coalescent analysis, we

631 used GTR substitution matrix, four Gamma categories, estimated among-site variation,

632 Calibrated Yule tree prior with ucldMean ucldStddev from exponential distribution,

633 relaxed lognormal molecular clock, shared common ancestor of all CrERVs 0.47-1 MYA

634 as a prior [19,54]. Maximum likelihood phylogeny in Figure 3 was generated using

635 PhyML [100] using the models selected by AICc and the setting of '-o tlr -s BEST'

636 according to the selected model.

637

638 <u>CrERV spatial distribution</u>

639 We simulated 274 insertions per genome to approximate the average number of

640 CrERVs in a mule deer [53]. The simulation was performed 10,000 times on three

641 genomes: the mule deer WGS scaffolds, cow (Btau7) and human (hg19). Distance

642 between simulated insertions and the closest start of the coding sequence of a gene

643 was calculated using the 'closest' function of bedtools, and the simulated insertions that

28

644　overlap with a gene were marked with the 'intersect' function of bedtools.  Number of

645　simulated simulations that are within 20 Kbp or intronic to a gene was counted for each

646　of the 10,000 replicates. Counts were then normalized by the total number of insertions

647　and plotted using the 'hist' function of R.

648

649　Supplementary methods

650　Methods with extended details are available in File S1.

651

652　**Availability of supporting data**

653

654　The raw sequencing data was deposited in SRR9121136. Other data generated are

655　included in supplementary file and figures.

656

657　**List of abbreviations**

658

659　ERV: endogenous retrovirus

660　XRV: exogenous retrovirus or infectious retrovirus

661　LTR: long terminal repeat

662　CrERV: cervid endogenous gammaretrovirus

663　CrXRV: cervid exogenous gammaretrovirus

664　HERV-K: human endogenous retrovirus type K

665　RACA: reference-assisted chromosome assembly

666　WGS: whole genome sequencing

667     NAHR: non-allelic homologous recombination

668     KYA: thousand years ago

669     MYA: million years ago

670     RT: reverse transcriptase

671     ORF: open reading frame

672     IAP: intracisternal A-type particle

673

674     **Author's contributions**

675     LY, RM, RC, TK, JR, PM, and MP conducted analyses; LY, RM, DE, MP interpreted

676     data; LY and MP wrote the manuscript.

677

678     **Acknowledgements**

679     The authors would like to thank David Chen for contributions to genome analysis.

680

681     **Declarations**

682

683     Funding

684     This work was funded in part by USGS 06HQAG0131.  The funders had no role in study

685     design, data collection and analysis, decision to publish, or preparation of the

686     manuscript.

687

688     Competing interests

689     The authors claim no competing interests.

690

Ethics approval and consent to participate

692 Not applicable.

693

Consent for publication

695 Not applicable.

696

**References**

698

699 1. Coffin JM. Retroviridae and their replication. Fields Virol. Lippincott-Raven; 1996. p.

700 1767–1848.

701 2. Weiss RA. The discovery of endogenous retroviruses. Retrovirology. 2006;3:67.

702 3. Löwer R, Löwer J, Kurth R. The viruses in all of us: characteristics and biological

703 significance of human endogenous retrovirus sequences. Proc Natl Acad Sci U S A.

704 1996;93:5177–84.

705 4. Stoye JP. Studies of endogenous retroviruses reveal a continuing evolutionary saga.

706 Nat Rev Microbiol. 2012;10:395–406.

707 5. Belshaw R, Katzourakis A, Pac□es J, Burt A, Tristem M. High Copy Number in

708 Human Endogenous Retrovirus Families is Associated with Copying Mechanisms in

709 Addition to Reinfection. Mol Biol Evol. 2005;22:814–7.

710 6. Johnson WE. Endogenous Retroviruses in the Genomics Era. Annu Rev Virol.

711 2015;2:135–59.

712 7. Belshaw R, Pereira V, Katzourakis A, Talbot G, Paces J, Burt A, et al. Long-term

713    reinfection of the human genome by endogenous retroviruses. Proc Natl Acad Sci U S

714    A. 2004;101:4894–9.

715    8. Boeke JD, Stoye JP. Retrotransposons, endogenous retroviruses, and the evolution

716    of retroelements. Cold Spring Harbor Laboratory Press, Cold Spring Harbor (NY);

717    1997;343–435.

718    9. Feschotte C, Gilbert C. Endogenous viruses: Insights into viral evolution and impact

719    on host biology. Nat Rev Genet. 2012;13:283–96.

720    10. Jern P, Coffin JM. Effects of retroviruses on host genome function. Annu Rev

721    Genet. 2008;42:709–32.

722    11. Kurth R, Bannert N. Beneficial and detrimental effects of human endogenous

723    retroviruses. Int. J. Cancer. 2010. p. 306–14.

724    12. Antony JM, DesLauriers AM, Bhat RK, Ellestad KK, Power C. Human endogenous

725    retroviruses and multiple sclerosis: Innocent bystanders or disease determinants?

726    Biochim Biophys Acta - Mol Basis Dis. 2011;1812:162–76.

727    13. Magiorkinis G, Belshaw R, Katzourakis A. "There and back again": revisiting the

728    pathophysiological roles of human endogenous retroviruses in the post-genomic era.

729    Philos Trans R Soc Lond B Biol Sci. 2013;368:20120504.

730    14. Wildschutte JH, Ram D, Subramanian R, Stevens VL, Coffin JM. The distribution of

731    insertionally polymorphic endogenous retroviruses in breast cancer patients and cancer-

732    free controls. Retrovirology. 2014;11:62.

733    15. Xue B, Sechi LA, Kelvin DJ. Human Endogenous Retrovirus K (HML-2) in Health

734    and Disease. Front Microbiol. 2020;11.

735    16. Li W, Lee M, Henderson L, Tyagi R, Bachani M, Steiner J, et al. Human

32

736    endogenous retrovirus-K contributes to motor neuron disease. Sci Transl Med.

737    2015;7:307ra153.

738    17. Li W, Yang L, Harris RS, Lin L, Olson TL, Hamele CE, et al. Retrovirus insertion site

739    analysis of LGL leukemia patient genomes. BMC Med Genomics. BMC Medical

740    Genomics; 2019;12:88.

741    18. Roca AL, O'Brien SP, Greenwood AD, Eiden M V., Ishida Y. Transmission,

742    Evolution, and Endogenization: Lessons Learned from Recent Retroviral Invasions.

743    Microbiol Mol Biol Rev. 2017;82:1–41.

744    19. Elleder D, Kim O, Padhi A, Bankert JG, Simeonov I, Schuster SC, et al. Polymorphic

745    integrations of an endogenous gammaretrovirus in the mule deer genome. J Virol.

746    2012;86:2787–96.

747    20. Arnaud F, Caporale M, Varela M, Biek R, Chessa B, Alberti A, et al. A Paradigm for

748    Virus–Host Coevolution: Sequential Counter-Adaptations between Endogenous and

749    Exogenous Retroviruses. PLoS Pathog. 2007;3:e170.

750    21. Coffin JM, Hughes SH, Varmus HE. Retroviral Virions and Genomes--Retroviruses.

751    Cold Spring Harbor Laboratory Press; 1997.

752    22. Luo GX, Taylor J. Template switching by reverse transcriptase during DNA

753    synthesis. J Virol. 1990;64:4321–8.

754    23. Bolinger C, Boris-Lawrie K. Mechanisms employed by retroviruses to exploit host

755    factors for translational control of a complicated proteome. Retrovirology. 2009;6:8.

756    24. Sofuku K, Honda T. Influence of Endogenous Viral Sequences on Gene Expression.

757    Gene Expr Regul Mamm Cells - Transcr From Gen Asp. InTech; 2018.

758    25. Kim H-S. Genomic impact, chromosomal distribution and transcriptional regulation

759    of HERV elements. Mol Cells. 2012;33:539–44.

760    26. Isbel L, Whitelaw E. Endogenous retroviruses in mammals: An emerging picture of

761    how ERVs modify expression of adjacent genes. BioEssays. 2012;34:734–8.

762    27. Yao S, Sukonnik T, Kean T, Bharadwaj RR, Pasceri P, Ellis J. Retrovirus Silencing,

763    Variegation, Extinction, and Memory Are Controlled by a Dynamic Interplay of Multiple

764    Epigenetic Modifications. Mol Ther. 2004;10:27–36.

765    28. Hurst TP, Magiorkinis G. Epigenetic control of human endogenous retrovirus

766    expression: Focus on regulation of long-terminal repeats (LTRs). Viruses. 2017;9:1–13.

767    29. Geis FK, Goff SP. Silencing and Transcriptional Regulation of Endogenous

768    Retroviruses: An Overview. Viruses. 2020;12:884.

769    30. Lavie L, Kitova M, Maldener E, Meese E, Mayer J. CpG Methylation Directly

770    Regulates Transcriptional Activity of the Human Endogenous Retrovirus Family HERV-

771    K(HML-2). J Virol. 2005;79:876–83.

772    31. Matsui T, Leung D, Miyashita H, Maksakova IA, Miyachi H, Kimura H, et al. Proviral

773    silencing in embryonic stem cells requires the histone methyltransferase ESET. Nature.

774    2010;464:927–31.

775    32. Bruno M, Mahgoub M, Macfarlan TS. The Arms Race Between KRAB–Zinc Finger

776    Proteins and Endogenous Retroelements and Its Impact on Mammals. Annu Rev

777    Genet. 2019;53:annurev-genet-112618-043717.

778    33. Sze A, Olagnier D, Lin R, van Grevenynghe J, Hiscott J. SAMHD1 Host Restriction

779    Factor: A Link with Innate Immune Sensing of Retrovirus Infection. J Mol Biol.

780    2013;425:4981–94.

781    34. Blanco-Melo D, Gifford RJ, Bieniasz PD. Co-option of an endogenous retrovirus

782 envelope for host defense in hominid ancestors. Elife. 2017;6.

783 35. Haig D. Retroviruses and the Placenta. Curr Biol. 2012;22:R609–13.

784 36. Fu B, Ma H, Liu D. Endogenous Retroviruses Function as Gene Expression

785 Regulatory Elements During Mammalian Pre-implantation Embryo Development. Int J

786 Mol Sci. 2019;20:790.

787 37. Göke J, Ng HH. CTRL + INSERT: retrotransposons and their contribution to

788 regulation and innovation of the transcriptome. EMBO Rep. 2016;17:1131–44.

789 38. Hughes JF, Coffin JM. Human endogenous retrovirus K solo-LTR formation and

790 insertional polymorphisms: implications for human and viral evolution. Proc Natl Acad

791 Sci U S A. 2004;101:1668–72.

792 39. Belshaw R, Dawson ALA, Woolven-Allen J, Redding J, Burt A, Tristem M.

793 Genomewide screening reveals high levels of insertional polymorphism in the human

794 endogenous retrovirus family HERV-K(HML2): implications for present-day activity. J

795 Virol. 2005;79:12507–14.

796 40. Subramanian RP, Wildschutte JH, Russo C, Coffin JM. Identification,

797 characterization, and comparative genomic distribution of the HERV-K (HML-2) group of

798 human endogenous retroviruses. Retrovirology. BioMed Central Ltd; 2011;8:90.

799 41. Hoang ML, Tan FJ, Lai DC, Celniker SE, Hoskins RA, Dunham MJ, et al.

800 Competitive repair by naturally dispersed repetitive DNA during non-allelic homologous

801 recombination. PLoS Genet. 2010;6:e1001228.

802 42. Belshaw R, Watson J, Katzourakis A, Howe A, Woolven-Allen J, Burt A, et al. Rate

803 of recombinational deletion among human endogenous retroviruses. J Virol.

804 2007;81:9437–42.

805    43. Kijima TE, Innan H. On the Estimation of the Insertion Time of LTR

806    Retrotransposable Elements. Mol Biol Evol. 2010;27:896–904.

807    44. Johnson WE, Coffin JM. Constructing primate phylogenies from ancient retrovirus

808    sequences. Proc Natl Acad Sci. 1999;96:10254–60.

809    45. Zhuo X, Rho M, Feschotte C. Genome-Wide Characterization of Endogenous

810    Retroviruses in the Bat Myotis lucifugus Reveals Recent and Diverse Infections. J Virol.

811    2013;87:8493–501.

812    46. Stocking C, Kozak CA. Murine endogenous retroviruses. Cell Mol Life Sci.

813    2008;65:3383–98.

814    47. Anai Y, Ochi H, Watanabe S, Nakagawa S, Kawamura M, Gojobori T, et al.

815    Infectious endogenous retroviruses in cats and emergence of recombinant viruses. J

816    Virol. 2012;86:8634–44.

817    48. Wildschutte JH, Williams ZH, Montesion M, Subramanian RP, Kidd JM, Coffin JM.

818    Discovery of unfixed endogenous retrovirus insertions in diverse human populations.

819    Proc Natl Acad Sci U S A. 2016;113:E2326-34.

820    49. Soriano P, Gridley T, Jaenisch R. Retroviruses and insertional mutagenesis in mice:

821    proviral integration at the Mov 34 locus leads to early embryonic death. Genes Dev.

822    1987;1:366–75.

823    50. Moyes D, Griffiths DJ, Venables PJ. Insertional polymorphisms: a new lease of life

824    for endogenous retroviruses in human disease. Trends Genet. 2007;23:326–33.

825    51. Turner G, Barbulescu M, Su M, Jensen-Seaman MI, Kidd KK, Lenz J. Insertional

826    polymorphisms of full-length endogenous retroviruses in humans. Curr Biol.

827    2001;11:1531–5.

36

828    52. Li W, Lin L, Malhotra R, Yang L, Acharya R, Poss M. A computational framework to

829    assess genome-wide distribution of polymorphic human endogenous retrovirus-K In

830    human populations. Wilke CO, editor. PLOS Comput Biol. 2019;15:e1006564.

831    53. Bao L, Elleder D, Malhotra R, DeGiorgio M, Maravegias T, Horvath L, et al.

832    Computational and Statistical Analyses of Insertional Polymorphic Endogenous

833    Retroviruses in a Non-Model Organism. Computation. 2014;2:221–45.

834    54. Kamath PL, Poss M, Elleder D, Powell JH, Bao L, Cross PC.  The Population

835    History of Endogenous Retroviruses in Mule Deer ( Odocoileus hemionus ) . J Hered.

836    2013;105:173–87.

837    55. Fábryová H, Hron T, Kabíčková H, Poss M, Elleder D. Induction and

838    characterization of a replication competent cervid endogenous gammaretrovirus

839    (CrERV) from mule deer cells. Virology. 2015;485:96–103.

840    56. Chaisson MJP, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, et

841    al. Resolving the complexity of the human genome using single-molecule sequencing.

842    Nature. Nature Publishing Group; 2015;517:608–11.

843    57. Gregory TR. Animal Genome Size Database. 2019.

844    58. Vinogradov AE. Genome size and GC-percent in vertebrates as determined by flow

845    cytometry: The triangular relationship. Cytometry. 1998;31:100–9.

846    59. Kim J, Larkin DM, Cai Q, Asan, Zhang Y, Ge R-L, et al. Reference-assisted

847    chromosome assembly. Proc Natl Acad Sci U S A. 2013;110:1785–90.

848    60. Gallagher DS, Derr JN, Womack JE. Chromosome conservation among the

849    advanced pecorans and determination of the primitive bovid karyotype. J Hered.

850    1994;85:204–10.

851    61. Cantarel BL, Korf I, Robb SMC, Parra G, Ross E, Moore B, et al. MAKER: an easy-

852    to-use annotation pipeline designed for emerging model organism genomes. Genome

853    Res. 2008;18:188–96.

854    62. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database

855    management tool for second-generation genome projects. BMC Bioinformatics. BioMed

856    Central Ltd; 2011;12:491.

857    63. Benit L, Dessen P, Heidmann T. Identification, Phylogeny, and Evolution of

858    Retroviral Elements Based on Their Envelope Genes. J Virol. 2001;75:11709–19.

859    64. Hunter DR, Bao L, Poss M. Assignment of endogenous retrovirus integration sites

860    using a mixture model. Ann Appl Stat. 2017;11:751–70.

861    65. Gifford RJ, Katzourakis A, De Ranter J, Magiorkinis G, Belshaw R. Env-less

862    endogenous retroviruses are genomic superspreaders. Proc Natl Acad Sci.

863    2012;109:7385–90.

864    66. MAGER DL, FREEMAN JD. HERV-H Endogenous Retroviruses: Presence in the

865    New World Branch but Amplification in the Old World Primate Lineage. Virology.

866    1995;213:395–404.

867    67. Suspène R, Sommer P, Henry M, Ferris S, Guétard D, Pochet S, et al. APOBEC3G

868    is a single-stranded DNA cytidine deaminase and functions independently of HIV

869    reverse transcriptase. Nucleic Acids Res. 2004;32:2421–9.

870    68. Rebollo R, Romanish MT, Mager DL. Transposable Elements: An Abundant and

871    Natural Source of Regulatory Sequences for Host Genes. Annu Rev Genet.

872    2012;46:21–42.

873    69. Zheng Y-H, Jeang K-T, Tokunaga K. Host restriction factors in retroviral infection:

874    promises in virus-host interaction. Retrovirology. 2012;9:112.

875    70. Duggal NK, Emerman M. Evolutionary conflicts between viruses and restriction

876    factors shape immunity. Nat Rev Immunol. 2012;12:687–95.

877    71. Daugherty MD, Malik HS. Rules of Engagement: Molecular Insights from Host-Virus

878    Arms Races. Annu Rev Genet. 2012;46:677–700.

879    72. Campbell IM, Gambin T, Dittwald P, Beck CR, Shuvarikov A, Hixson P, et al.

880    Human endogenous retroviral elements promote genome instability via non-allelic

881    homologous recombination. BMC Biol. 2014;12:74.

882    73. Cohen CJ, Lock WM, Mager DL. Endogenous retroviral LTRs as promoters for

883    human genes: A critical assessment. Gene. 2009;448:105–14.

884    74. Maksakova IA, Romanish MT, Gagnier L, Dunn CA, van de Lagemaat LN, Mager

885    DL. Retroviral elements and their hosts: insertional mutagenesis in the mouse germ

886    line. PLoS Genet. 2006;2:e2.

887    75. Lu X, Sachs F, Ramsay L, Jacques P-É, Göke J, Bourque G, et al. The retrovirus

888    HERVH is a long noncoding RNA required for human embryonic stem cell identity. Nat

889    Struct Mol Biol. Nature Publishing Group; 2014;21:423–5.

890    76. Kawasaki J, Nishigaki K. Tracking the Continuous Evolutionary Processes of an

891    Endogenous Retrovirus of the Domestic Cat: ERV-DC. Viruses. 2018;10:179.

892    77. Bénit L, De Parseval N, Casella JF, Callebaut I, Cordonnier A, Heidmann T. Cloning

893    of a new murine endogenous retrovirus, MuERV-L, with strong similarity to the human

894    HERV-L element and with a gag coding sequence closely related to the Fv1 restriction

895    gene. J Virol. 1997;71:5652–7.

896    78. Finnerty H, Mi S, Veldman GM, McCoy JM, LaVallie E, Edouard P, et al. Syncytin is

39

897    a captive retroviral envelope protein involved in human placental morphogenesis.

898    Nature. 2002;403:785–9.

899    79. Murin CD, Wilson IA, Ward AB. Antibody responses to viral infections: a structural

900    perspective across three different enveloped viruses. Nat Microbiol. 2019;4:734–47.

901    80. Stamatatos L, Morris L, Burton DR, Mascola JR. Neutralizing antibodies generated

902    during natural HIV-1 infection: good news for an HIV-1 vaccine? Nat Med. 2009;15:866–

903    70.

904    81. Evans LH, Alamgir ASM, Owens N, Weber N, Virtaneva K, Barbian K, et al.

905    Mobilization of Endogenous Retroviruses in Mice after Infection with an Exogenous

906    Retrovirus. J Virol. 2009;83:2429–35.

907    82. Hanafusa H. Analysis of the defectiveness of rous sarcoma virus III. Determining

908    influence of a new helper virus on the host range and susceptibility to interference of

909    RSV. Virology. 1965;25:248–55.

910    83. Dewannieux M, Dupressoir A, Harper F, Pierron G, Heidmann T. Identification of

911    autonomous IAP LTR retrotransposons mobile in mammalian cells. Nat Genet.

912    2004;36:534–9.

913    84. Halo J V., Pendleton AL, Jarosz AS, Gifford RJ, Day ML, Kidd JM. Origin and recent

914    expansion of an endogenous gammaretroviral lineage in domestic and wild canids.

915    Retrovirology. 2019;16:6.

916    85. Ali L, Rizvi T, Mustafa F. Cross- and Co-Packaging of Retroviral RNAs and Their

917    Consequences. Viruses. 2016;8:276.

918    86. Kozak C. Origins of the Endogenous and Infectious Laboratory Mouse

919    Gammaretroviruses. Viruses. 2014;7:1–26.

87. Bamunusinghe D, Naghashfar Z, Buckler-White A, Plishka R, Baliji S, Liu Q, et al. Sequence Diversity, Intersubgroup Relationships, and Origins of the Mouse Leukemia Gammaretroviruses of Laboratory and Wild Mice. Beemon KL, editor. J Virol. 2016;90:4186–98.

88. Löber U, Hobbs M, Dayaram A, Tsangaras K, Jones K, Alquezar-Planas DE, et al. Degradation and remobilization of endogenous retroviruses by recombination during the earliest stages of a germ-line invasion. Proc Natl Acad Sci. 2018;115:8609–14.

89. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. Gigascience. 2012;1:18.

90. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics. 2009;25:1754–60.

91. Trapnell C, Pachter L, Salzberg SL. TopHat: discovering splice junctions with RNA-Seq. Bioinformatics. 2009;25:1105–11.

92. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. Genome Biol. BioMed Central Ltd; 2013;14:R36.

93. Kent WJ. BLAT---The BLAST-Like Alignment Tool. Genome Res. 2002;12:656–64.

94. Harris RS. Improved pairwise alignment of genomic DNA. The Pennsylvania State University; 2007.

95. Bininda-Emonds ORP, Cardillo M, Jones KE, MacPhee RDE, Beck RMD, Grenyer R, et al. The delayed rise of present-day mammals. Nature. 2007;446:507–12.

96. Edgar RC. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Res. 2004;32:1792–7.

943    97. Löytynoja A, Goldman N. An algorithm for progressive multiple alignment of

944    sequences with insertions. Proc Natl Acad Sci U S A. 2005;102:10557–62.

945    98. Posada D. jModelTest: phylogenetic model averaging. Mol Biol Evol. 2008;25:1253–

946    6.

947    99. Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, et al. BEAST 2: a

948    software platform for Bayesian evolutionary analysis. PLoS Comput Biol.

949    2014;10:e1003537.

950    100. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large

951    phylogenies by maximum likelihood. Syst Biol. 2003;52:696–704.

952

953    **Figure legends**

954

955    **Figure 1.  Diagram of CrERV reconstruction and RACA.** (A) Mule deer chromosome

956    fragment reconstruction using syntenic fragments. Gray, green and blue boxes

957    correspond to aligned human, cow and mule deer scaffold respectively. Lighter shades

958    represent regions that can only be aligned between two species. Dashed boxes

959    highlight syntenic fragments where the region is conserved among all three species,

960    which yield a chromosome fragment that orients mule deer scaffolds. (B)

961    Reconstruction of CrERV sequences. CrERV and mule deer scaffolds are shown in bold

962    orange and blue boxes, respectively. Long insert mate pair reads are connected by

963    dotted lines and are colored to indicate whether they derive from the mule deer scaffold

964    or CrERV genome. CrERV genomes were assembled by gathering the broken mate

965    pairs surrounding each CrERV loci as described.

966

**Figure 2. Coalescent phylogeny, *env* structural variation and population frequency of representative full-length non-recombinant CrERVs.** Nodes with at least 95% posterior probability support are marked by black dots. The high posterior density for each labeled node is shown in Table S6. Boxes next to CrERV names display the frequency of the CrERVs in the mule deer population with a gray scale (annotated at the top-left corner). Diagrams on the right side depict the lineage-specific structural variations in the CrERV envelope gene. White triangles represent insertions (A, B, C), and white rectangles represent deletions (D and E).

975

**Figure 3. Recombination among CrERVs.** Shown is a maximum likelihood phylogeny based on a region spanning a portion of *pol* to 5'*env* (JN592050: 4422-7076). Taxa used are a subset of full-length non-recombinant CrERVs representing the four lineages shown in Figure 2 and CrERVs with a recombinant signature containing a Lineage B *env*. Supported nodes (aLRT >= 0.85) are represented by black dots on the backbone of the tree. Lineage designation is assigned to supported branches based on the non-recombinant CrERV. Over this interval, Lineage B CrERVs are found as a sister group to lineage A CrERV but some CrERV containing a prototypical Lineage B *env* are dispersed among Lineage A CrERV. Note that in this interval lineage C CrERVs cluster with Lineage A CrERVs.

986

**Figure 4. CrERV insertions are enriched within 20 kbp of genes and depleted in introns.** We simulated the expected number of CrERV insertions by randomly placing

43

989   them on the *de novo* assembled MT273 genome. The proportion of insertions expected

990   within 20kb of a gene from the 10,000 replicates is shown in Panel A. The proportion of

991   intronic insertions is in panel B. The distribution of insertions within 20kb of a gene or an

992   intron from the simulation is shown as a histogram. Blue dashed lines indicate the mean

993   of the simulated data. Red dashed lines indicate the observed data in MT273. Black

994   dashed lines indicate the 5th and 95th percentile of the simulated data, which are used

995   to call significant differences.

996

997   **Additional files**

998

999   **Figure S1. *K-mer* representation of missing data in assemblies.** *60-mers* were

1000  generated from the raw paired-end WGS, then ranked and classified into 20 bins

1001  containing equal number of *60-mers*. Sorted by the *60-mer* frequency range, bars

1002  represent genome repeats from high copy number (left) to low copy number (right).

1003  Color bars show percent of all *60-mers* present in the scaffolds/contigs (blue),

1004  RACA10K (green), and RACA150K (red) assemblies. *K-mer* counts were based on the

1005  total number of *k-mers* (*k-mer* of frequency *n* were counted *n* times). *K-mers* in the raw

1006  sequencing data were normalized based on sequencing depth, genome ploidy, read

1007  length and *k-mer* length, so that the *k-mer* fractions reflect the proportion of *k-mers* that

1008  are present in each assembly compared to the raw sequencing data. All bars start from

1009  0% instead of being stacked. Numbers beneath each bar indicates the range of

1010  frequency of *60-mers* in the raw paired-end genome sequencing data in that bin. M:

1011  million, K: thousand.

1012

1013    **Figure S2.  Diagram of CrERV recombination breakpoints.**  Gray lines point at the

1014    key recombination breakpoints on the CrERV.  Text box connected to the gray lines

1015    indicate the coordinate and adjusted p-value of the breakpoint.  Solid gray lines indicate

1016    breakpoints of recombinant lineages; dashed gray lines indicate additional breakpoints

1017    detected by testing the alignment of reference non-recombinant and candidate

1018    recombinant CrERVs.  All coordinates are relative to GenBank entry JN592050. Double

1019    star (**) indicates breakpoints used in the Lineage B recombinant analysis.

1020

1021    **Figure S3.  Distribution of simulated mean distance to gene per replicate in mule**

1022    **deer, cow and human genome.**  Distribution of mule deer, cow and human are colored

1023    in red, green and blue respectively.  Mann-Whitney U test p-values in all three

1024    comparisons are less than $2.2\text{x}10^{-16}$.

1025

1026    **Table S1. de novo and RACA assembly statistics.**  As the resolution increases,

1027    scaffolds can be placed into less chromosome fragments at the expense of less

1028    scaffolds incorporated. PE: paired-end sequencing.  MP: long insert mate pair

1029    sequencing. N50: length of the shortest contig at 50% of the total genome length, a

1030    measurement of assembly contiguity. RACA: reference-assisted chromosome

1031    assembly. RACA assembly size: total size of RACA chromosome fragments at given

1032    RACA resolution.

1033

45

1034  **Table S2. Scaffold assignment to RACA chromosome fragments.** The tab

1035  "RACA10K" corresponds to the RACA at 10 Kbp resolution, and "RACA150K"

1036  corresponds to 150 Kbp resolution. "R150K.R10K.CowChain10KSort" contains the

1037  chromosome fragment assignment information of both RACA150K (column A-G) and

1038  RACA10K (column H-N), sorted by the scaffolds' alignment chain (column O-V) to the

1039  cow genome, with column W indicating the genes that are present on the scaffold.

1040  Similarly, "R150K.R10K.SheepChain10KSort" represents scaffold assignments to

1041  RACA150K and RACA10K chromosome fragments along the sheep genome.

1042

1043  **Table S3. Number of gene structures annotated after each maker annotation.** Each

1044  column represents a Maker iteration by their order.

1045

1046  **Table S4. Summary of *k-mers* with > 50 frequency.** Numbers outside of parentheses

1047  show the cumulative frequency of *k-mers* with >50 frequency belonging to each repeat

1048  family.  Given the read length of 100 bp and k-mer size of 60, frequency of >50 in the

1049  whole genome sequencing library corresponds to >4 copies in the genome.  Percentage

1050  in parentheses show the fraction of *k-mers* of the repeat family that are present in each

1051  assembly denominated by the total of that family in the raw paired-end sequencing

1052  library.

1053

1054  **Table S5. Inventory of CrERVs used in phylogenetic analysis.** Column A: CrERV

1055  names, by scaffold number. Multiples CrERVs on the same scaffold are discerned by

1056  additional character after scaffold numbers. Column B: Lineage designation, consistent

46

1057    with assignments mentioned in text and Figure 2. B1 and B2 are the two lineages that

1058    are closely related to lineage B but have a different type of env. Bold font of lineage A

1059    indicates the sub-lineage that are significantly closer to genes. Column C: Env status,

1060    as illustrated in Figure 2. A, B: insertion; D, E: deletion; C: insertion and deletion; M:

1061    missing data for the corresponding insertion or deletion. Column D: Inter-lineage

1062    recombination status. Text indicates the recombination partner and section. Empty cell

1063    means no recombination was detected. Column E: Assignment to RACA10K

1064    chromosome fragment (also refer to Table S2). Parentheses mean that they are not on

1065    RACA10K and provisionally assigned by being the closest to sheep or cow aligmnt

1066    chain. Column F: Frequency in all 63 animals used by Bao et al. 2014, using >=0.95

1067    probability cutoff. The repetitive junction fragments (multiple mapping) are designated

1068    'uncertain' and their frequency in the probability table was listed in the parenthesis. Bold

1069    font indicates the singletons. Column G: Distance to the closest gene.  'NA' means that

1070    the virus-containing scaffold cannot be assigned to RACAs and no genes can be found

1071    on the CrERV-containing scaffold. Negative value means intronic insertion. Bold font

1072    indicates recombinant CrERVs (Figure 3) that are close to genes. Column H: Presence

1073    in Figure 2, with letter in parenthesis representing the node it corresponds to.

1074

1075    **Table S6. Intervals for the 95% highest probability density (HPD) intervals for key**

1076    **nodes in Figure 2.**

1077

1078    **Table S7. Sequences used for recombination breakpoint tests and results of the**

1079    **tests.** (A) List of reference non-recombinant CrERVs used for each lineage. (B)

47

1080    Summary of recombination breakpoint tests. Statistically significant breakpoints

1081    (adjusted $p < 0.01$) are shown. Breakpoint coordinates are based on GenBank
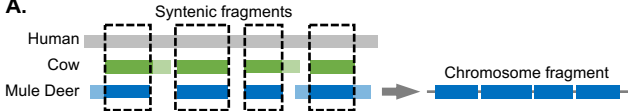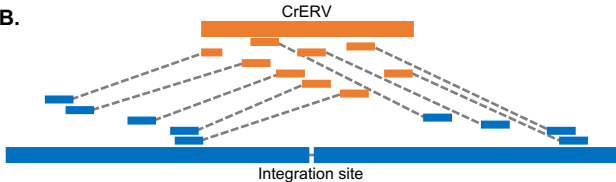
1082    JN592050.

1083

1084    **File S1. Supplementary analyses and methods.**
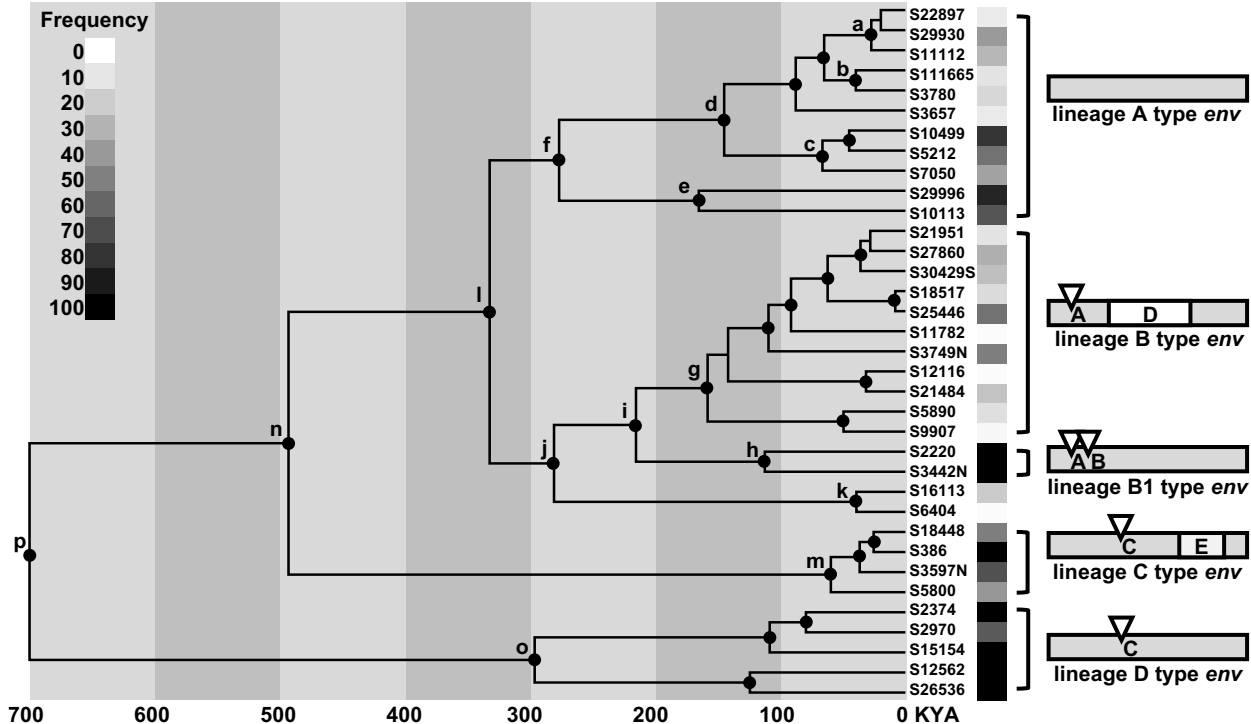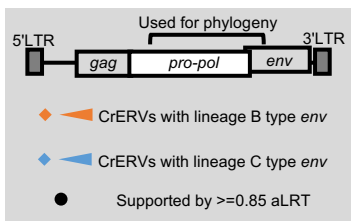
1085

1086    **File S2. Reconstructed CrERV sequences used in phylogenetic analyses.** The file

1087    is in fasta format. CrERVs are named after the scaffold they came from, as listed in

1088    Table S5.

1089

**A.** Syntenic fragments

Human
Cow
Mule Deer

Chromosome fragment

**B.** CrERV

Integration site

Used for phylogeny

5'LTR | gag | pro-pol | env | 3'LTR

◆ ◀ CrERVs with lineage B type *env*

◆ ◀ CrERVs with lineage C type *env*

● Supported by >=0.85 aLRT

lineage A

lineage B

lineage D

**A. <=20kb**

**B. in 100kb**

Observed
Simulated mean

Frequency

Frequency

Frequency

Percent of insertions that fall in category