

1 **A Reptilian Endogenous Foamy Virus Sheds Light on the Early Evolution** 2 **of Retroviruses**

3

4 Xiaoman Wei^{1,2†}, Yicong Chen^{1,2†}, Edward C. Holmes³, Jie Cui^{1*}

5

6

7 ¹Key Laboratory of Special Pathogens and Biosafety, Center for Emerging Infectious
8 Diseases, Wuhan Institute of Virology, Chinese Academy of Sciences, Wuhan 430071,
9 China.

10 ²University of Chinese Academy of Sciences, Beijing 100049, China.

11 ³Marie Bashir Institute for Infectious Diseases and Biosecurity, School of Life and
12 Environmental Sciences and Faculty of Medicine and Health, University of Sydney, Sydney,
13 NSW 2006, Australia.

14

15 [†]These authors contributed equally to this work.

16 ^{*}Corresponding author: E-mail: jiecui@wh.iov.cn

17 **Abstract**

18 Endogenous retroviruses (ERVs) represent host genomic fossils of ancient viruses. Foamy
19 viruses, including those that form endogenous copies, provide strong evidence for virus-host
20 co-divergence across the vertebrate phylogeny. Endogenous foamy viruses (EFV) have
21 previously been discovered in mammals, amphibians and fish. Here we report a novel
22 endogenous foamy virus, named SpuEFV, in genome of the tuatara (*Sphenodon punctatus*), a
23 reptile species endemic to New Zealand. Surprisingly, SpuEFV robustly grouped with the
24 coelacanth EFV on virus phylogenies, rather than with the mammalian foamy viruses as
25 expected with virus-host co-divergence, and indicative of a major cross-species transmission
26 event in the early evolution of the foamy viruses. In sum, the discovery of SpuEFV fills a
27 major gap in the fossil record of foamy viruses and provides important insights into the early
28 evolution of retroviruses.

29

30 **Key words:** endogenous retroviruses; foamy virus; reptiles; evolution; cross-species
31 transmission

32

33 Retroviruses (family *Retroviridae*) are viruses of major medical significance as some are
34 associated with severe infectious disease or are oncogenic (Hayward, et al. 2015; Aiewsakun
35 and Katzourakis 2017; Xu, et al. 2018). Retroviruses are also of note because of their ability
36 to integrate into the host germ-line, generating endogenous retroviruses (ERVs) that then
37 exhibit Mendelian inheritance (Stoye 2012; Johnson 2015). ERVs are widely distributed in
38 vertebrates and provide important molecular “fossils” for the study of retrovirus evolution.
39 ERVs related to all seven major retroviral genera have been described, although some of the
40 more complex retroviruses, such as lenti-, delta- and foamy viruses, rarely appear as
41 endogenous copies.

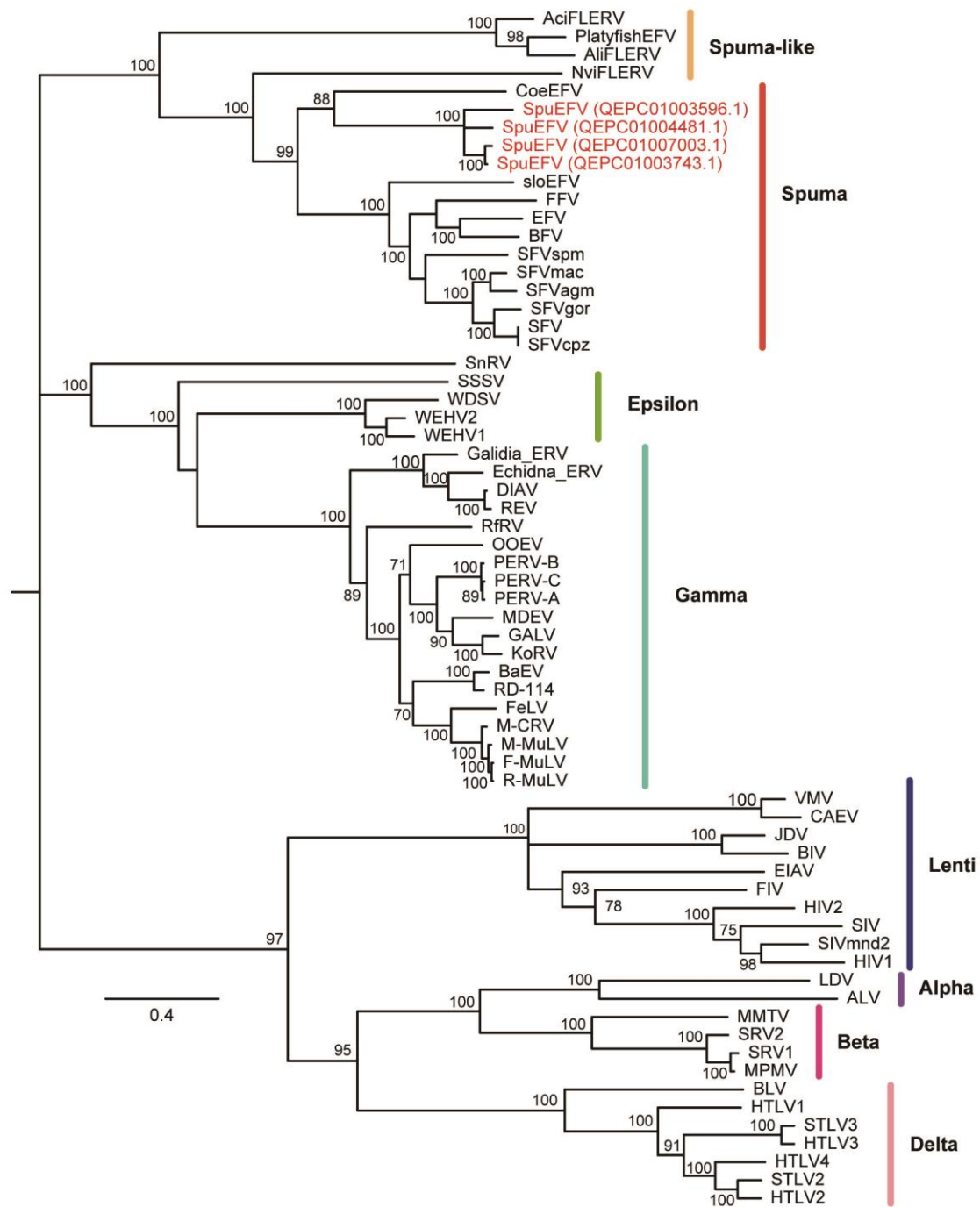
42
43 As well as being agents of disease, foamy viruses are of importance because of their long-
44 term virus-host co-divergence (Switzer, et al. 2005). Endogenous foamy viruses (EFVs), first
45 discovered in sloths (class Mammalia) (Katzourakis, et al. 2009) also exhibit co-divergence.
46 The later discovery of a fish EFV in the coelacanth genome indicated that foamy viruses have
47 an ancient evolutionary history (Han and Worobey 2012) and hence have likely co-diverged
48 with their vertebrate hosts for hundreds of million years (Aiewsakun and Katzourakis 2017).
49 However, although EFVs or foamy-like elements have been reported in fish, amphibians and
50 mammals, they have currently not been reported in genomes of reptiles.

51
52 To search for potential foamy (-like) viral elements in reptiles, we collated 28 reptilian
53 genomes (Supplementary Table S1) and performed *in silico* TBLASTN with full-length Pol
54 sequences of various foamy viruses, including EFVs, as screening probes (Supplementary
55 Table S2). We only considered viral hits within long genomic scaffold (>20 kilobases in
56 length) to be *bona fide* ERVs. This genomic mining identified 175 ERV hits in three species:
57 tuatara (*Sphenodon punctatus*), Schlegel's Japanese Gecko (*Gekko japonicas*) and

58 Madagascar ground gecko (*Paroedura picta*). However, because only one viral hit of each
59 was found in the Schlegel's Japanese Gecko and Madagascar ground gecko (accession
60 number: LNDG01066615.1 and BDOT01000314.1), which could represent false-positives,
61 they were excluded. Hence, a total of 173 ERV hits in the tuatara genome were extracted and
62 subjected to evolutionary analysis (Supplementary Table S3).

63

64 The long Pol (>700 amino acids) and Env (>350) sequences of these ERVs were then
65 selected for phylogenetic analysis. Our maximum likelihood (ML) phylogenetic tree revealed
66 that the ERVs discovered in tuatara genome formed a close monophyletic group within the
67 foamy clade, indicative of a single origin, and with high bootstrap supports in both
68 phylogenies (Fig. 1; Fig. S1). We named this new ERV as SpuEFV (*Sphenodon punctatus*
69 endogenous foamy virus). To our surprise, SpuEFV was consistently and robustly related to
70 the fish EFV – CoeEFV – derived from the coelacanth genome (Han and Worobey 2012),
71 and hence in conflict with the known host phylogeny. Although this phylogenetic pattern is
72 compatible with cross-class virus transmission from fish to reptiles, it is possible that this
73 pattern will change with a larger sampling of taxa such that the EFV phylogeny expands.
74 Failure to detect any SpuEFV related elements in the remaining reptilian genome screening
75 suggests that the virus was not vertically transmitted among reptiles, although this will clearly
76 need to be reassessed with a larger sample size.



77

78

79 We successfully retrieved two full-length SpuEFV viral genomes and annotated one in detail

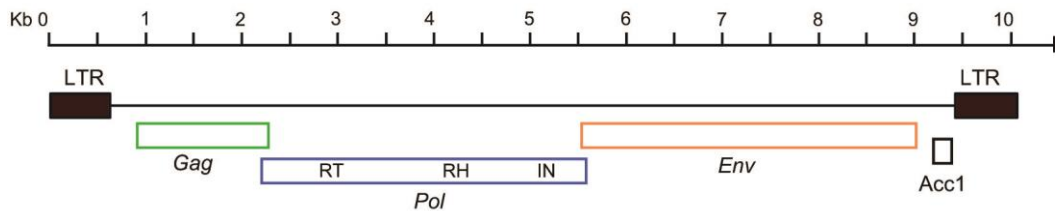
80 (Fig. S2). The annotated sequence exhibits a typical spuma virus structure, encoding three

81 mainly open reading frames (ORF) – *gag*, *pol* and *env* – and one additional accessory genes,

82 Acc1 (Fig. 2). Interestingly, this accessory ORF (Acc1) exhibit no sequence similarity to

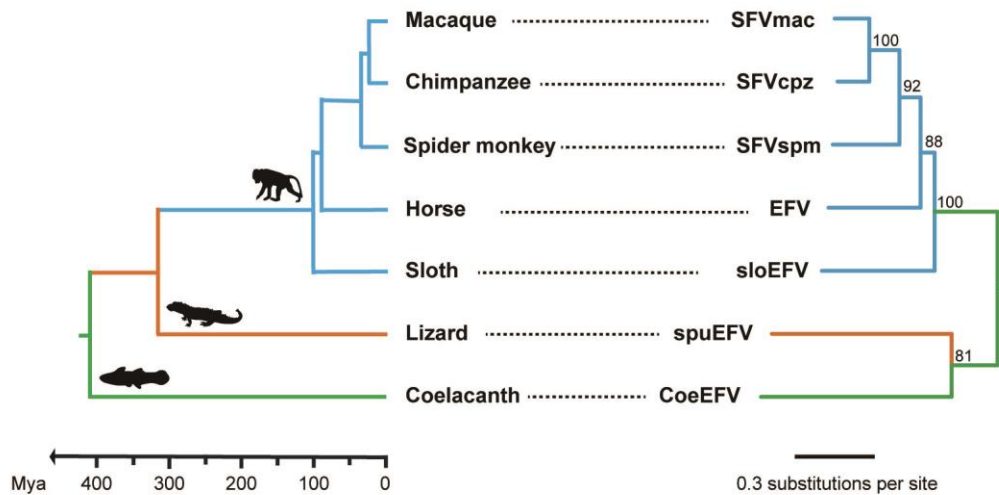
83 known genes. Notably, by searching the Conserved Domains Database

84 (www.ncbi.nlm.nih.gov/Structure/cdd), we identified a typical conserved foamy virus
85 envelope protein domain (pfam034308) (Han and Worobey 2012), they further confirming
86 that SpuEFV is of foamy virus origin.



87
88 To broadly estimate the integration time of SpuEFVs, we employed the LTR (long terminal
89 repeat)-divergence method, which analyzes the degree of divergence between 5' and 3'LTRs
90 assuming a known rate of nucleotide substitution (Johnson and Coffin 1999). In total, five
91 pairwise LTRs flanking SpuEFV elements were used for date estimation (Supplementary
92 Table S4), from which we estimated an integration time of SpuEFV ranging from 1.3 to
93 35.47 MYA (million years ago). Although these dates are young relative to the age of
94 reptiles, LTR dating may severely underestimate ERV ages (Kijima and Innan 2010;
95 Aiewsakun and Katzourakis 2017), such that all estimates of integration time should be
96 treated with caution.

97
98 Previous studies provided strong evidence for the co-divergence of foamy viruses and their
99 vertebrate hosts over extended time-periods (Katzourakis, et al. 2009). That the reptilian
100 SpuEFV newly described here was most closely related to fish EFVs than those found in
101 mammalian genomes (Fig. 3) indicates that cross-species virus transmission on a back-bone
102 of long-term virus-host co-divergence may also play a major role in shaping the early
103 evolution of retroviruses.



104

105 **Materials and Methods**

106 **Genomic mining**

107 To identify foamy viruses in reptiles, the TBLASTN program (Altschul, et al. 1990) was used

108 to screen relevant taxa from 28 reptile genomes downloaded from GenBank

109 (www.ncbi.nlm.nih.gov/genbank) (Supplementary Table S1). In each case amino acid

110 sequences of the Pol and Env genes of representative EFVs (endogenous foamy viruses),

111 foamy-like sequences, and foamy viruses were chosen as queries. As filters to identify

112 significant and meaningful hits, we chose sequences with more than 30% amino acid identity

113 over a 30% genomic region, with an e-value set to 0.00001. We extended viral flanking

114 sequences of the hits to identify the 5'- and 3'-LTRs using LTR finder (Xu and Wang 2007)

115 and LTR harvest (Ellinghaus, et al. 2008).

116

117 **Phylogenetic analysis**

118 To determine the evolutionary relationship of EFVs and retroviruses, Pol and Env protein

119 sequences were aligned in MAFFT 7.222 (Kato and Standley 2013) and confirmed

120 manually in MEGA7 (Kumar, et al. 2016). The phylogenetic relationships among these

121 sequences were then determined using the maximum-likelihood (ML) method in PhyML 3.1

122 (Guindon, et al. 2010), incorporating 100 bootstrap replicates to determine node robustness.

123 The best-fit models of amino acid substitution were determined by ProtTest 3.4.2 (Abascal, et
124 al. 2005): RtREV+ Γ +I for Pol, and WAG+ Γ for Env. All alignments used in the phylogenetic
125 analyses can be found in Data set S1.

126

127 **Molecular dating**

128 The ERV integration time can be calculated using the following simple relation: $T = (D/R)/2$,
129 in which T is the integration time (million years, MY), D is the number of nucleotide
130 differences per site between the two LTRs, and R is the genomic substitution rate (i.e.
131 number of nucleotide substitutions per site, per year). We used the previously estimated
132 neutral substitution rate for squamate reptiles (7.6×10^{-10} nucleotide substitutions per site,
133 per year) (Perry, et al. 2018). LTRs less than 300 bp in length were not included in this
134 analysis.

135

136 **Acknowledgments**

137 J.C. is supported by National Natural Science Foundation of China (31671324) and CAS
138 Pioneer Hundred Talents Program. ECH is supported by an ARC Australian Laureate
139 Fellowship (FL170100022).

140

141 **References**

142

143 Abascal F, Zardoya R, Posada D. 2005. ProtTest: selection of best-fit models of protein
144 evolution. *Bioinformatics* 21:2104-2105.

145 Aiewsakun P, Katzourakis A. 2017. Marine origin of retroviruses in the early Palaeozoic Era.
146 *Nat Commun* 8:13954.

147 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search
148 tool. *J Mol Biol* 215:403-410.

- 149 Ellinghaus D, Kurtz S, Willhoeft U. 2008. LTRharvest, an efficient and flexible software for
150 de novo detection of LTR retrotransposons. *BMC Bioinformatics*.
- 151 Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. 2010. New
152 algorithms and methods to estimate maximum-likelihood phylogenies: assessing the
153 performance of PhyML 3.0. *Syst Biol* 59:307-321.
- 154 Han GZ, Worobey M. 2012. An endogenous foamy-like viral element in the coelacanth
155 genome. *PLoS Pathog* 8:e1002790.
- 156 Hayward A, Cornwallis CK, Jern P. 2015. Pan-vertebrate comparative genomics unmasks
157 retrovirus macroevolution. *Proc Natl Acad Sci U S A* 112:464-469.
- 158 Johnson WE. 2015. Endogenous Retroviruses in the Genomics Era. *Annu Rev Virol* 2:135-
159 159.
- 160 Johnson WE, Coffin JM. 1999. Constructing primate phylogenies from ancient retrovirus
161 sequences. *Proc Natl Acad Sci U S A* 96:10254-10260.
- 162 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
163 improvements in performance and usability. *Mol Biol Evol* 30:772-780.
- 164 Katzourakis A, Gifford RJ, Tristem M, Gilbert MT, Pybus OG. 2009. Macroevolution of
165 complex retroviruses. *Science* 325:1512.
- 166 Kijima TE, Innan H. 2010. On the estimation of the insertion time of LTR retrotransposable
167 elements. *Mol Biol Evol* 27:896-904.
- 168 Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis
169 Version 7.0 for Bigger Datasets. *Mol Biol Evol* 33:1870-1874.
- 170 Perry BW, Card DC, McGlothlin JW, Pasquesi GIM, Adams RH, Schield DR, Hales NR,
171 Corbin AB, Demuth JP, Hoffmann FG, et al. 2018. Molecular Adaptations for Sensing and

172 Securing Prey and Insight into Amniote Genome Diversity from the Garter Snake Genome.

173 Genome Biol Evol 10:2110-2129.

174 Stoye JP. 2012. Studies of endogenous retroviruses reveal a continuing evolutionary saga.

175 Nat Rev Microbiol 10:395-406.

176 Switzer WM, Salemi M, Shanmugam V, Gao F, Cong ME, Kuiken C, Bhullar V, Beer BE,

177 Vallet D, Gautier-Hion A, et al. 2005. Ancient co-speciation of simian foamy viruses and

178 primates. Nature 434:376-380.

179 Xu X, Zhao H, Gong Z, Han GZ. 2018. Endogenous retroviruses of non-avian/mammalian

180 vertebrates illuminate diversity and deep history of retroviruses. PLoS Pathog 14:e1007072.

181 Xu Z, Wang H. 2007. LTR_FINDER: an efficient tool for the prediction of full-length LTR

182 retrotransposons. Nucleic Acids Res 35:W265-268.

183

184 **Figure Legends**

185

186 **Figure 1.** Phylogenetic tree of retroviruses, including SpuEFV, using amino acid sequences

187 of the Pol gene. The phylogenetic tree was rooted using *Caenorhabditis elegans*

188 retrotransposon Cer1 (GenBank accession no. U15406). The newly identified SpuEFVs are

189 labelled in red along with their accession numbers. The scale bar indicates the number of

190 amino acid changes per site. Bootstrap values <70% are not shown.

191

192 **Figure 2.** Genomic organizations of SpuEFV. LTR, long-terminal repeat; RT, reverse

193 transcriptase; RH, ribonuclease H; IN, integrase.

194

195 **Figure 3.** A simplified evolutionary relationship between foamy viruses and their hosts.

196 Phylogenies representing mammals, reptile and fish and their associated viruses are shown.

- 197 The scale bar indicates host speciation time (million years ago, MYA) or the number of
198 amino acid changes per site in the viral genomes.