1    **Discovery of flavivirus-derived endogenous viral elements in two *Anopheles***

2    **mosquito genomes supports the existence of *Anopheles*-associated insect-**

3    **specific flaviviruses**

4

5    *Sebastian Lequime[1,2,3*], Louis Lambrechts[1,2]*

6

7    [1]Insect-Virus Interactions Group, Department of Genomes and Genetics, Institut Pasteur,

8    Paris, France

9    [2]Centre National de la Recherche Scientifique, Unité de Recherche Associée 3012, Paris,

10   France

11   [3]University Pierre et Marie Curie, Cellule Pasteur UPMC, Paris, France

12

13

14   [*]Corresponding author. Address: Insect-Virus Interactions, Institut Pasteur, 28 rue du Docteur

15   Roux, 75724 Paris Cedex 15, France. E-mail: sebastian.lequime@gmail.com

16

17 **Abstract**

18    The *Flavivirus* genus encompasses several arboviruses of public health significance

19 such as dengue, yellow fever, and Zika viruses. It also includes insect-specific flaviviruses

20 (ISFs) that are only capable of infecting insect hosts. The vast majority of mosquito-infecting

21 flaviviruses have been associated with mosquito species of the *Aedes* and *Culex* genera in the

22 Culicinae subfamily, which also includes most arbovirus vectors. Mosquitoes of the

23 Anophelinae subfamily are not considered significant arbovirus vectors, however flaviviruses

24 have occasionally been detected in field-caught *Anopheles* specimens. Whether such

25 observations reflect occasional spillover or laboratory contamination or whether *Anopheles*

26 mosquitoes are natural hosts of flaviviruses is unknown. Here, we provide *in silico* and *in vivo*

27 evidence of transcriptionally active, flavivirus-derived endogenous viral elements (EVEs) in

28 the genome of *Anopheles minimus* and *Anopheles sinensis*. Such non-retroviral

29 endogenization of RNA viruses is consistent with a shared evolutionary history between

30 flaviviruses and *Anopheles* mosquitoes. Phylogenetic analyses of the two newly described

31 EVEs support the existence of a distinct clade of *Anopheles*-associated ISFs.

32

**Introduction**

33

34      Flaviviruses are positive-sense, single-stranded RNA viruses that infect a broad range of

35      hosts including vertebrates (e.g., birds, primates) and arthropods (e.g., ticks, mosquitoes). In

36      addition to major arboviruses of public health significance such as dengue, Zika, West Nile

37      and yellow fever viruses, the *Flavivirus* genus also includes vertebrate-specific (not known

38      vector; NKV) and insect-specific (insect-specific flaviviruses; ISFs) members (Moureau et al.

39      2015). The majority of mosquito-infecting flaviviruses have been associated with members of

40      the Culicinae subfamily, mainly from the *Culex* and *Aedes* genera. *Anopheles* mosquitoes in

41      the Anophelinae subfamily are well known for their role in the transmission of human malaria

42      parasites, but they are not considered important vectors of arboviruses in general, and of

43      flaviviruses in particular. Nevertheless, several studies have detected flaviviruses in field-

44      caught *Anopheles* mosquitoes from different parts of the world. In North America, West Nile

45      virus (WNV) was detected in *Anopheles punctipennis* (Bernard et al. 2001; Kulasekera et al.

46      2001). In Asia, Japanese encephalitis virus was detected in *Anopheles sinensis* (Feng et al.

47      2012; Liu et al. 2013). In Europe, *Anopheles maculipennis* was found positive for WNV

48      (Filipe 1972; Kemenesi et al. 2014), Usutu virus (Calzolari et al. 2013) and Batai virus

49      (Calzolari et al. 2010). Interestingly, some ISFs were also detected in *An. sinensis* (Zuo et al.

50      2014; Liang et al. 2015) and *Anopheles atroparvus* (Aranda et al. 2009). It is unknown

51      whether these detections reflect occasional spillover or laboratory contamination, or whether

52      *Anopheles* mosquitoes are in fact natural hosts of flaviviruses.

53      Endogenous viral elements (EVEs) are chromosomal integrations of partial or full viral

54      genetic material into the genome of a host species. Not only retroviruses, whose replication

55      cycle includes incorporation of a DNA form of the RNA viral genome into the host cell

56      genome, but virtually all types of eukaryotic viruses can become endogenous (Feschotte and

57      Gilbert 2012). Non-retroviral EVEs have been documented in the genome of a wide variety of

58    host species, including vertebrates and arthropods (Feschotte and Gilbert 2012). Unlike

59    detection of exogenous viruses, subject to possible laboratory or environmental

60    contamination, EVEs are likely to reflect a long-lasting evolutionary relationship between an

61    RNA virus and its natural host. This is because endogenization, for a single-stranded RNA

62    virus, requires (1) reverse transcription, (2) integration of virus-derived DNA into the genome

63    of germinal host cells and (3) fixation of the integrated sequence in the host population

64    (Holmes 2011; Aiewsakun and Katzourakis 2015). The low probability of this combination of

65    events makes endogenization exceedingly unlikely to occur unless the viral infection is

66    common in the host population over long evolutionary times. For example, flavivirus-derived

67    EVEs have been reported in the genome of *Aedes aegypti* (Crochu et al. 2004; Katzourakis

68    and Gifford 2010) and *Aedes albopictus* (Roiz et al. 2009; Chen et al. 2015). These EVEs are

69    phylogenetically related to the clade of *Aedes*-associated ISFs (Crochu et al. 2004; Roiz et al.

70    2009; Katzourakis and Gifford 2010; Chen et al. 2015), which is consistent with the ancient

71    evolutionary relationship between *Aedes* mosquitoes and ISFs.

72    Here, we report the discovery of two flavivirus-derived EVEs in the genomes of *An.*

73    *minimus* and *An. sinensis* mosquitoes. We screened 24 publicly available *Anopheles* genomes

74    (Holt et al. 2002; Zhou et al. 2014; Logue et al. 2015; Neafsey et al. 2015) for flaviviral

75    sequences, and validated *in silico* hits both at the DNA and RNA levels *in vivo*. The two

76    newly described flavivirus-derived EVEs are phylogenetically related to ISFs, and support the

77    existence of a previously unsuspected *Anopheles*-associated clade of ISFs.

78

79      **Material and Methods**

80      1.  *In silico* survey

81      1.1. Genome screen

82      Twenty-four *Anopheles* genomes (full list and accession numbers are provided in

83      Table 1) were screened by tBLASTn (BLAST+ 2.2.28) (Camacho et al. 2009) for the

84      presence of flavivirus-derived EVEs using a collection of 50 full flavivirus polyprotein

85      queries (full list and accession numbers are provided in Table S1) representing the currently

86      known diversity of the *Flavivirus* genus. The sequences of hits whose E-value was $> E^{-4}$ were

87      extracted using an in-house bash shell script. In order to reconstruct putative flavivirus-

88      derived EVEs, BLAST hits were clustered using CD-HIT v.4.6.1 (Li and Godzik 2006) and

89      aligned using MAFFT v7.017 (Katoh et al. 2002). Putative EVEs were extracted and used as

90      query for a reciprocal tBLASTn (BLAST+ 2.2.30) against the National Center for

91      Biotechnology Information (NCBI) nucleotide database (E-value $> 10^{-4}$). Genetic features of

92      identified EVE were analyzed using the NCBI Conserved Domain Database (Marchler-Bauer

93      et al. 2015). Nucleotide sequence data reported are available in the Third Party Annotation

94      Section of the DDBJ/ENA/GenBank databases under the accession numbers TPA:

95      BK009978-BK009980.

96      1.2. Phylogenetic analyses

97      Translated EVE sequences were aligned to the corresponding sections of several

98      flavivirus polyproteins (Table S1) with MAFFT v7.017 and phylogenetically uninformative

99      positions were trimmed using TrimAI v.1.3 (Capella-Gutiérrez et al. 2009) accessed through

100     the webserver Phylemon 2 (Sánchez et al. 2011). The trimmed alignments were used to

101     construct phylogenetic trees with PhyML Best AIC Tree (Sánchez et al. 2011). Best

102     substitution models were Blosum62+I+G+F for *An. minimus* and Blosum62+I+G for *An.*

103     *sinensis*.

104         1.3. Transcriptome screen

105         Published RNA sequencing (RNA-seq) data were retrieved from NCBI Sequence

106 Read Archive (Leinonen et al. 2011) and explored for the presence of previously identified

107 EVE sequences. Only one *An. minimus* transcriptome sequence read archive was found under

108 accession number SRX265162. Six *An. sinensis* transcriptome sequence read archives were

109 found under accession numbers SRX448985, SRX449003, SRX449006 and SRX277584 for

110 experiments using Illumina sequencing technology, and SRX218691 and SRX218721 for

111 experiments using Roche 454 sequencing technology. RNA-seq reads were mapped to the

112 EVE nucleotide sequence using Bowtie2 v2.1.0 (Langmead and Salzberg 2012). The

113 alignment file was converted, sorted and indexed with Samtools v0.1.19 (Li et al. 2009).

114 Coverage was assessed using bedtools v2.17.0 (Quinlan and Hall 2010).

115         2. *In vivo* validation

116         2.1. Mosquitoes

117         *Anopheles minimus* and *An. sinensis* mosquitoes were obtained through BEI Resources

118 (www.beiresources.org), National Institute of Allergy and Infectious Diseases, National

119 Institutes of Health (*An. minimus* MINIMUS1, MRA-729; *An. sinensis* SINENSIS, MRA-

120 1154). *Anopheles minimus* and *An. sinensis* specimens came from the $132^{nd}$ and $65^{th}$

121 generations of laboratory colonization, respectively. Eggs were hatched in filtered tap water,

122 reared in 24×34×9 cm plastic trays and fed with fish food (TetraMin, Tetra, Melle, Germany).

123 Adults were maintained in 30×30×30 cm screened cages under controlled insectary conditions

124 (28±1°C, 75±5% relative humidity, 12:12 hour light-dark cycle). They were provided with

125 cotton soaked in a 10% (m/v) sucrose solution *ad libitum. Anopheles stephensi* nucleic acids,

126 used as a reaction control, were kindly provided by the Genetics and Genomics of Insect

127 Vectors unit, Institut Pasteur, Paris.

128         2.2. EVE genomic integration

129    Mosquitoes were homogenized in pools of 10 separated by sex in 300 µL of

130    Dulbecco's phosphate-buffered saline (DPBS) during two rounds of 30 sec at 5,000 rpm in a

131    mixer mill (Precellys 24, Bertin Technologies, Montigny le Bretonneux, France). DNA was

132    extracted using All Prep DNA/RNA Mini Kit (Qiagen, Hilden, Germany) following the

133    manufacturer's instructions. EVE presence in genomic DNA was assessed by 35 cycles of

134    PCR using Taq Polymerase (Invitrogen, Thermo Fisher Scientific, Waltham, MA, USA)

135    (Table S2). PCR primers were designed to generate an amplicon spanning part of the EVE

136    sequence and a section of the flanking host sequence. Identity of the EVE sequence was

137    confirmed by Sanger sequencing of the PCR product.

138    2.3. EVE transcription level

139    Mosquitoes were homogenized in pools of 5 separated by sex or development stage in

140    300 µL of DPBS during two rounds of 30 sec at 5,000 rpm in a mixer mill (Precellys 24).

141    RNA was extracted from mosquito homogenates separated by sex using TRIzol Reagent (Life

142    Technologies, Thermo Fisher Scientific, Waltham, MA, USA) following the manufacturer's

143    instructions. Samples were treated with Turbo DNA-free kit (Life Technologies) and reverse

144    transcribed using random hexamers and M-MLV reverse transcriptase (Invitrogen).

145    Complementary DNA was amplified with 35 cycles of PCR for *An. minimus* and 40 cycles of

146    PCR for *An. sinensis,* respectively*,* using DreamTaq polymerase (Thermo Fisher Scientific)

147    and primers located within the EVE sequence (Table S2). To verify that RNA samples were

148    free of DNA contamination, two sets of primers spanning exons 3 and 4 of the RPS7 gene of

149    both *Anopheles* species (under VectorBase annotation number AMIN008193 and

150    ASIC017918 for *An. minimus* and *An. sinensis*, respectively) were designed (Table S2).

151    Because the corresponding DNA sequence includes intron 3, DNA contamination is expected

152    to result in a larger PCR product. The length of intron 3 is 252 base pairs (bp) for *An. minimus*

153    and 295 bp for *An. sinensis.*

154

155    **Results**

156    The *in silico* screen of 24 *Anopheles* genomes identified two flavivirus-derived EVEs, one

157    in the reference genome sequence of *An. minimus* and one in both genome sequences

158    available for *An. sinensis* (Table 1).

159    ***An. minimus* EVE**

160    The *An.  minimus* EVE is 1,881 bp long (627 amino acid residues) with Nienokoue virus

161    as the closest BLAST hit (44% amino-acid identity). The integrated sequence spans non-

162    structural protein 4A (NS4A), NS4B and NS5 (Figure 1A) without any stop codon.

163    Conserved domain search identified the NS5-methyltransferase domain involved in RNA

164    capping and part of the RNA-directed RNA-polymerase domain. Phylogenetically, the *An.*

165    *minimus* EVE is sister to the ISF clade (Figure 1B). The *An. minimus* genome supercontig

166    where the EVE was detected is 2,043 bp long and consists almost exclusively of the EVE

167    sequence.

168    Presence of the EVE was verified *in vivo* by PCR on genomic DNA (Figure 2A), followed

169    by amplicon sequencing to confirm identity. The *An. minimus* EVE was found in both male

170    and female genomic DNA, and was transcriptionally expressed for all combinations of sex

171    and development stages tested (Figure 2B). Evidence for transcriptional activity of the *An.*

172    *minimus* EVE was confirmed in published RNA-seq data (Figure S1A).

173    ***An. sinensis* EVE**

174    The *An. sinensis* EVE was detected in two distinct genome sequences that are available

175    for this mosquito species, derived from a Korean and a Chinese strain, respectively. The EVE

176    is 792 bp long (264 amino acid residues) and 799 bp long (266 amino acid residues) for the

177    Korean and Chinese strains, respectively (Table 1). The closest BLAST hit is *Culex* flavivirus

178    (45% amino-acid identity for the Korean strain, 44% amino-acid identity for the Chinese

179    strain). The integrated sequence corresponds to the middle part of NS3 (Figure 3A) and

180    contains six and eight stop codons in the Korean and Chinese strains, respectively. Conserved

181    domain search identified the presence of a P-loop containing the nucleoside triphosphate

182    hydrolase domain found in the NS3 protein of exogenous flaviviruses. Phylogenetically, the

183    *An. sinensis* EVE is sister to the ISF clade (Figure 3B).

184    The *An. sinensis* supercontig containing the EVE is 2,797 bp long for the Korean strain

185    and 7,380 bp long for the Chinese strain. Analyses of flanking regions revealed the presence

186    of another EVE in the same orientation, upstream of the flavivirus-derived EVE in both the

187    Korean and the Chinese strains. The closest BLAST hit of this additional EVE is Xincheng

188    mosquito virus (43% and 42% amino-acid identity for the Korean and Chinese strains,

189    respectively), an unclassified, negative-sense, single-stranded RNA virus detected in a pool of

190    Chinese mosquitoes including *An. sinensis* specimens (Table S3). BLAST and conserved

191    domain search identified a class II Mariner-like transposase close to a mariner mos1 element,

192    approximately 1,000 bp downstream of the flavivirus-derived EVE (Table S3). This was only

193    the case for the Chinese strain because the supercontig of the Korean strain was not long

194    enough.

195    Presence of the *An. sinensis* EVE was verified *in vivo* by PCR on genomic DNA from the

196    Korean strain (Figure 4A), followed by amplicon sequencing to confirm identity. The *An.*

197    *sinensis* EVE was found in both male and female genomic DNA, and was transcriptionally

198    expressed, although less abundantly than the *An. sinensis* EVE, for all combinations of sex

199    and development stages tested, especially in $L_4$ larvae (Figure 4B). Low expression observed

200    for the *An. sinensis* EVE is consistent with barely detectable transcriptional activity in

201    published RNA-seq data (Figure S1B).

202

**Discussion**

ISFs have attracted substantial interest in recent years after some of them were shown to enhance or suppress the replication of medically important flaviviruses in co-infected mosquito cells (Blitvich and Firth 2015). Over a dozen of ISFs have been formally identified to date, mainly in *Aedes* and *Culex* genera of the Culicinae subfamily (Blitvich and Firth 2015). ISFs were also reported in *Anopheles* mosquitoes of the Anophelinae subfamily (Aranda et al. 2009; Zuo et al. 2014; Liang et al. 2015). However, these *Anopheles*-associated ISFs are thought to infect a broad range of hosts including several mosquito species, mainly in the *Culex* genus, and are phylogenetically related to *Culex*-associated ISFs. Therefore, it is unclear whether *Anopheles* mosquitoes are true natural hosts of flaviviruses. Detection of ISFs in field-caught mosquitoes could result from incidental infection, or from a laboratory artifact. In this study, we discovered flavivirus-derived EVEs in the genomes of two *Anopheles* species. Phylogenetic analyses indicated that both EVEs are related to ISFs but belong to a clade that is distinct from *Aedes*-associated and *Culex*-associated ISFs.

Presence of flavivirus-derived EVEs in *Anopheles* genomes supports the hypothesis that *Anopheles* mosquitoes are natural hosts of flaviviruses. Endogenization of non-retroviral RNA viruses is unlikely to occur in the absence of recurrent host-virus interactions over a long evolutionary time scale. Endogenization requires reverse transcription, germ line integration and fixation in the host population, three steps whose combined frequency is exceedingly rare (Holmes 2011; Aiewsakun and Katzourakis 2015). The species-wide frequency of the *An. minimus* EVE is unknown because our *in silico* and *in vivo* analyses were based on the same mosquito strain. Presence of the *An. sinensis* EVE in two mosquito strains from different geographical locations, however, suggests that it could be fixed at the species level. Thus, our discovery of flavivirus-derived EVEs in *Anopheles* genomes is consistent with a long-lasting host-virus interaction between flaviviruses and mosquitoes of the Anophelinae subfamily.

228        ISFs are thought to be mainly maintained through vertical transmission from an infected

229        female to its offspring (Blitvich and Firth 2015). Vertical transmission is likely to favor co-

230        divergence of pathogens and hosts (Jackson and Charleston 2004), as illustrated by the

231        existence of *Aedes*-associated and *Culex*-associated clades of ISFs (Moureau et al. 2015).

232        Although extrapolation is limited by the scarcity of data on ISF host range and diversity,

233        phylogenetic position of *Anopheles*-associated ISFs as sister to all other known ISFs is

234        consistent with the co-divergence hypothesis. During the evolutionary history of mosquitoes,

235        the Anophelinae diverged from the Culicinae prior to the separation of *Culex* and *Aedes*

236        genera (Reidenbach et al. 2009). Further investigations will be necessary to determine

237        whether an *Anopheles*-associated clade of exogenous ISFs exists, or existed.

238        Non-retroviral EVEs are thought to be generated by interaction of exogenous viruses with

239        endogenous retro-elements, either with or without long terminal repeats (LTR) (Holmes

240        2011). The short size of the supercontigs containing the *An. minimus* and *An. sinensis* EVEs

241        limited our ability to investigate the integration mechanism(s). Another EVE sequence that we

242        identified close to the flavivirus-derived EVE in *An. sinensis* may point to an EVE hotspot.

243        Sequence conservation of the *An. minimus* EVE (i.e., absence of stop codon across 1,881

244        bp) is consistent with a recent integration or a more ancient integration followed by non-

245        neutral evolution. Our observation that the *An. minimus* EVE is abundantly transcribed may

246        reflect a selective advantage for the host (Holmes 2011). Transcriptionally active EVEs have

247        been suggested to confer protection or tolerance against related exogenous viruses (Flegel

248        2009; Holmes 2011; Aswad and Katzourakis 2012; Bell-Sakyi and Attoui 2013; Fujino et al.

249        2014). Despite the lack of empirical evidence so far, flavivirus-derived EVEs could contribute

250        to antiviral immunity and arbovirus vector competence in mosquitoes.

251

252     **Acknowledgements**

264

265 **References**

266 Aiewsakun P, Katzourakis A. 2015. Endogenous viruses: Connecting recent and ancient viral
267     evolution. Virology 479-480:26–37.

268 Aranda C, Sánchez-Seco MP, Cáceres F, Escosa R, Gálvez JC, Masià M, Marqués E, Ruíz S,
269     Alba A, Busquets N, et al. 2009. Detection and monitoring of mosquito flaviviruses in
270     Spain between 2001 and 2005. Vector Borne Zoonotic Dis. 9:171–178.

271 Aswad A, Katzourakis A. 2012. Paleovirology and virally derived immunity. Trends Ecol.
272     Evol. (Amst.) 27:627–636.

273 Bell-Sakyi L, Attoui H. 2013. Endogenous tick viruses and modulation of tick-borne pathogen
274     growth. Front Cell Infect Microbiol 3:25.

275 Bernard KA, Maffei JG, Jones SA, Kauffman EB, Ebel G, Dupuis AP, Ngo KA, Nicholas DC,
276     Young DM, Shi PY, et al. 2001. West Nile virus infection in birds and mosquitoes, New
277     York State, 2000. Emerging Infect. Dis. 7:679–685.

278 Blitvich BJ, Firth AE. 2015. Insect-specific flaviviruses: a systematic review of their
279     discovery, host range, mode of transmission, superinfection exclusion potential and
280     genomic organization. Viruses 7:1927–1959.

281 Calzolari M, Bonilauri P, Bellini R, Albieri A, Defilippo F, Tamba M, Tassinari M, Gelati A,
282     Cordioli P, Angelini P, et al. 2013. Usutu virus persistence and West Nile virus inactivity
283     in the Emilia-Romagna region (Italy) in 2011. PLoS ONE 8:e63978.

284 Calzolari M, Bonilauri P, Bellini R, Caimi M, Defilippo F, Maioli G, Albieri A, Medici A,
285     Veronesi R, Pilani R, et al. 2010. Arboviral survey of mosquitoes in two northern Italian
286     regions in 2007 and 2008. Vector Borne Zoonotic Dis. 10:875–884.

287 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
288     BLAST+: architecture and applications. BMC Bioinformatics 10:421.

289 Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated
290     alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25:1972–1973.

291 Chen X-G, Jiang X, Gu J, Xu M, Wu Y, Deng Y, Zhang C, Bonizzoni M, Dermauw W,
292     Vontas J, et al. 2015. Genome sequence of the Asian Tiger mosquito, Aedes albopictus,
293     reveals insights into its biology, genetics, and evolution. Proc. Natl. Acad. Sci. U.S.A.
294     112:E5907–E5915.

295 Crochu S, Cook S, Attoui H, Charrel RN, De Chesse R, Belhouchet M, Lemasson J-J, de
296     Micco P, De Lamballerie X. 2004. Sequences of flavivirus-related RNA viruses persist in
297     DNA form integrated in the genome of Aedes spp. mosquitoes. J. Gen. Virol. 85:1971–
298     1980.

299 Feng Y, Fu S, Zhang H, Li M, Zhou T, Wang J, Zhang Y, Wang H, Tang Q, Liang G. 2012.
300     Distribution of mosquitoes and mosquito-borne viruses along the China-Myanmar border
301     in Yunnan Province. Jpn. J. Infect. Dis. 65:215–221.

302 Feschotte C, Gilbert C. 2012. Endogenous viruses: insights into viral evolution and impact on

303     host biology. Nat. Rev. Genet. 13:283–296.

304     Filipe AR. 1972. Isolation in Portugal of West Nile virus from Anopheles maculipennis
305        mosquitoes. Acta Virol. 16:361.

306     Flegel TW. 2009. Hypothesis for heritable, anti-viral immunity in crustaceans and insects.
307        Biol. Direct 4:32.

308     Fujino K, Horie M, Honda T, Merriman DK, Tomonaga K. 2014. Inhibition of Borna disease
309        virus replication by an endogenous bornavirus-like element in the ground squirrel genome.
310        Proc. Natl. Acad. Sci. U.S.A. 111:13175–13180.

311     Holmes EC. 2011. The evolution of endogenous viral elements. Cell Host Microbe 10:368–
312        377.

313     Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, Nusskern DR, Wincker P,
314        Clark AG, Ribeiro JMC, Wides R, et al. 2002. The genome sequence of the malaria
315        mosquito Anopheles gambiae. Science 298:129–149.

316     Jackson AP, Charleston MA. 2004. A cophylogenetic perspective of RNA-virus evolution.
317        Mol. Biol. Evol. 21:45–57.

318     Katoh K, Misawa K, Kuma K-I, Miyata T. 2002. MAFFT: a novel method for rapid multiple
319        sequence alignment based on fast Fourier transform. Nucleic Acids Res. 30:3059–3066.

320     Katzourakis A, Gifford RJ. 2010. Endogenous viral elements in animal genomes. PLoS Genet.
321        6:e1001191.

322     Kemenesi G, Krtinić B, Milankov V, Kutas A, Dallos B, Oldal M, Somogyi N, Nemeth V,
323        Banyai K, Jakab F. 2014. West Nile virus surveillance in mosquitoes, April to October
324        2013, Vojvodina province, Serbia: implications for the 2014 season. Euro Surveill.
325        19:20779.

326     Kulasekera VL, Kramer L, Nasci RS, Mostashari F, Cherry B, Trock SC, Glaser C, Miller JR.
327        2001. West Nile virus infection in mosquitoes, birds, horses, and humans, Staten Island,
328        New York, 2000. Emerging Infect. Dis. 7:722–725.

329     Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nat. Methods
330        9:357–U54.

331     Leinonen R, Sugawara H, Shumway M, International Nucleotide Sequence Database
332        Collaboration. 2011. The sequence read archive. Nucleic Acids Res. 39:D19–D21.

333     Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
334        1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map
335        format and SAMtools. Bioinformatics 25:2078–2079.

336     Li W, Godzik A. 2006. Cd-hit: a fast program for clustering and comparing large sets of
337        protein or nucleotide sequences. Bioinformatics 22:1658–1659.

338     Liang W, He X, Liu G, Zhang S, Fu S, Wang M, Chen W, He Y, Tao X, Jiang H, et al. 2015.
339        Distribution and phylogenetic analysis of Culex flavivirus in mosquitoes in China. Arch.

340      Virol. 160:2259–2268.

341    Liu H, Lu H-J, Liu Z-J, Jing J, Ren J-Q, Liu Y-Y, Lu F, Jin N-Y. 2013. Japanese encephalitis
342        virus in mosquitoes and swine in Yunnan province, China 2009-2010. Vector Borne
343        Zoonotic Dis. 13:41–49.

344    Logue K, Small ST, Chan ER, Reimer L, Siba PM, Zimmerman PA, Serre D. 2015. Whole-
345        genome sequencing reveals absence of recent gene flow and separate demographic
346        histories for Anopheles punctulatus mosquitoes in Papua New Guinea. Mol. Ecol.
347        24:1263–1274.

348    Marchler-Bauer A, Derbyshire MK, Gonzales NR, Lu S, Chitsaz F, Geer LY, Geer RC, He J,
349        Gwadz M, Hurwitz DI, et al. 2015. CDD: NCBI's conserved domain database. Nucleic
350        Acids Res. 43:D222–D226.

351    Moureau G, Cook S, Lemey P, Nougairede A, Forrester NL, Khasnatinov M, Charrel RN,
352        Firth AE, Gould EA, De Lamballerie X. 2015. New insights into flavivirus evolution,
353        taxonomy and biogeographic history, extended by analysis of canonical and alternative
354        coding sequences. PLoS ONE 10:e0117849.

355    Neafsey DE, Waterhouse RM, Abai MR, Aganezov SS, Alekseyev MA, Allen JE, Amon J,
356        Arcà B, Arensburger P, Artemov G, et al. 2015. Mosquito genomics. Highly evolvable
357        malaria vectors: the genomes of 16 Anopheles mosquitoes. Science 347:1258522–
358        1258522.

359    Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
360        features. Bioinformatics 26:841–842.

361    Reidenbach KR, Cook S, Bertone MA, Harbach RE, Wiegmann BM, Besansky NJ. 2009.
362        Phylogenetic analysis and temporal diversification of mosquitoes (Diptera: Culicidae)
363        based on nuclear genes and morphology. BMC Evol. Biol. 9:298.

364    Roiz D, Vázquez A, Seco MPS, Tenorio A, Rizzoli A. 2009. Detection of novel insect
365        flavivirus sequences integrated in Aedes albopictus (Diptera: Culicidae) in Northern Italy.
366        Virol. J. 6:93.

367    Sánchez R, Serra F, Tárraga J, Medina I, Carbonell J, Pulido L, de María A, Capella-
368        Gutiérrez S, Huerta-Cepas J, Gabaldón T, et al. 2011. Phylemon 2.0: a suite of web-tools
369        for molecular evolution, phylogenetics, phylogenomics and hypotheses testing. Nucleic
370        Acids Res. 39:W470–W474.

371    Zhou D, Zhang D, Ding G, Shi L, Hou Q, Ye Y, Xu Y, Zhou H, Xiong C, Li S, et al. 2014.
372        Genome sequence of Anopheles sinensis provides insight into genetics basis of mosquito
373        competence for malaria parasites. BMC Genomics 15.

374    Zuo S, Zhao Q, Guo X, Zhou H, Cao W, Zhang J. 2014. Detection of Quang Binh virus from
375        mosquitoes in China. Virus Res. 180:31–38.

376

377

378   **Figures**

379   **Figure 1: Discovery of a flavivirus-derived EVE in *An. minimus*.** (A) EVE location in a

380   generic *Flavivirus* genome. Positioning is based on the genome sequence of Nienokoue virus

381   (GenBank accession no. JQ957875). C=capsid protein, E=envelope glycoprotein,

382   M=membrane glycoprotein, NS1=non-structural glycoprotein 1; NS2A=non-structural protein

383   2A; NS2B= non-structural protein 2B; NS3=non-structural protein 3 (protease/helicase);

384   NS4A=non-structural protein 4A; NS4B=non-structural protein 4B; NS5=non-structural

385   protein 5 (RNA-dependent RNA polymerase). (B) Phylogenetic relationships of the newly

386   discovered *An. minimus* flavivirus-derived EVE with exogenous flaviviruses. Maximum

387   likelihood trees were constructed based on the translated EVE sequence. Clades are color-

388   coded according to known host specificity: green, ISFs; purple, tick-borne arboviruses; black,

389   'not known vector' (vertebrate specific); blue, mosquito-borne arboviruses; red: EVEs. Scale

390   bar indicates the number of substitutions. Node values represent Shimodaira-Hasegawa (SH)-

391   like branch support (only values > 0.8 are shown).

392   **Figure 2: *In vivo* detection of the *An. minimus* flavivirus-derived EVE.** (A) EVE detection

393   in genomic DNA. Lane 1: size marker; lane 2: amplified genomic DNA from a pool of 10 *An.*

394   *minimus* adult females; lane 3: amplified genomic DNA from a pool of 10 *An. minimus* adult

395   males; lane 4: amplified genomic DNA from a pool of 10 *An. sinensis* adult females; lane 5:

396   amplified genomic DNA from a pool of 10 *An. sinensis* adult males; lane 6: amplified DNA

397   from a pool of 10 *An. stephensi* females; lane 7: no template control (NTC). (B) EVE

398   detection in total RNA. Lane 1: size marker; lanes 2 and 3: amplified cDNA from pools of 5

399   adult females; lanes 4 and 5: amplified cDNA from pools of 5 adult males; lanes 6 and 7:

400   amplified cDNA from pools of 5 $L_4$ larvae; lane 8: amplified DNA from a pool of 10 females;

401   lane 9: amplified cDNA from a pool of 5 *An. stephensi* females; lane 10: DNA contamination

402    control (no reverse transcription) using the same pool of 5 adult females as lane 2; lane 11: no

403    template control (NTC). First row: EVE; second row: RPS7 (control gene). The RPS7 target

404    DNA sequence includes an intron, so that DNA contamination is expected to result in a larger

405    PCR product.

406    **Figure 3: Discovery of a flavivirus-derived EVE in *An. sinensis*.** (A) EVE location in a

407    generic *Flavivirus* genome. Positioning is based on the genome sequence of *Culex* flavivirus

408    (GenBank accession no. JQ308188). C=capsid protein, E=envelope glycoprotein,

409    M=membrane glycoprotein, NS1=non-structural glycoprotein 1; NS2A=non-structural protein

410    2A; NS2B= non-structural protein 2B; NS3=non-structural protein 3 (protease/helicase);

411    NS4A=non-structural protein 4A; NS4B=non-structural protein 4B; NS5=non-structural

412    protein 5 (RNA-dependent-RNA polymerase). (B) Phylogenetic relationships of the newly

413    discovered *An. sinensis* flavivirus-derived EVEs with exogenous flaviviruses. Maximum

414    likelihood trees were constructed based on the translated EVE sequence. Clades are color-

415    coded according to known host specificity: green, ISFs; purple, tick-borne arboviruses; black,

416    'not known vector' (vertebrate specific); blue, mosquito-borne arboviruses; red: EVEs. Scale

417    bar indicates the number of substitutions. Node values represent Shimodaira-Hasegawa (SH)-

418    like branch support (only values > 0.8 are shown).

419    **Figure 4: *In vivo* detection of the *An. sinensis* flavivirus-derived EVE.** (A) EVE detection

420    in genomic DNA from the Korean strain of *An. sinensis*. Lane 1: size marker; lane 2:

421    amplified genomic DNA from a pool of 10 *An. minimus* adult females; lane 3: amplified

422    genomic DNA from a pool of 10 *An. minimus* adult males; lane 4: amplified genomic DNA

423    from a pool of 10 *An. sinensis* adult females; lane 5: amplified genomic DNA from a pool of

424    10 *An. sinensis* adult males; lane 6: amplified DNA from a pool of 10 *An. stephensi* females;

425    lane 7: no template control (NTC). (B) EVE detection in total RNA from the Korean strain of

426    *An. sinensis*. Lane 1: size marker; lanes 2 and 3: amplified cDNA from pools of 5 adult

427    females; lanes 4 and 5: amplified cDNA from pools of 5 adult males; lanes 6 and 7: amplified

428    cDNA from pools of 5 $L_4$ larvae; lane 8: amplified DNA from a pool of 10 females; lane 9:

429    amplified cDNA from a pool of 5 *An. stephensi* females; lane 10: DNA contamination control

430    (no reverse transcription) using the same pool of 5 adult females as lane 2; lane 11: no

431    template control (NTC). First row: EVE; second row: RPS7 (control gene). The RPS7 target

432    DNA sequence includes an intron, so that DNA contamination is expected to result in a larger

433    PCR product.

434     **Tables**

435     **Table 1:** *Anopheles* genomes screened in this study.

| Species | GenBank WGS Project | Assembly | GenBank Assembly ID |
|---|---|---|---|
| *Anopheles albimanus* | APCK01 | AalbS1 | GCA_000349125.1 |
| *Anopheles arabiensis* | APCN01 | AaraD1 | GCA_000349185.1 |
| *Anopheles atroparvus* | AXCP01 | AatrE1 | GCA_000473505.1 |
| *Anopheles christyi* | APCM01 | AchrA1 | GCA_000349165.1 |
| *Anopheles coluzzii* | ABKP02 | AcolM1 | GCA_000150765.1 |
| *Anopheles culicifacies A* | AXCM01 | AculA1 | GCA_000473375.1 |
| *Anopheles darlingi* | ADMH02 | AdarC3 | GCA_000211455.3 |
| *Anopheles dirus A* | APCL01 | AdirW1 | GCA_000349145.1 |
| *Anopheles epiroticus* | APCJ01 | AepiE1 | GCA_000349105.1 |
| *Anopheles farauti* | AXCN01 | AfarF1 | GCA_000473445.1 |
| *Anopheles funestus* | APCI01 | AfunF1 | GCA_000349085.1 |
| *Anopheles gambiae* | AAAB01 | AgamP4 | GCA_000005575.2 |
| *Anopheles koliensis* | JXXB01 | AKwgs3 | GCA_000956275.1 |
| *Anopheles maculatus B* | AXCL01 | AmacM1 | GCA_000473185.1 |
| *Anopheles melas* | ACXO01 | AmelC1 | GCA_000473525.1 |
| *Anopheles merus* | AXCQ01 | AmerM1 | GCA_000473845.1 |
| *Anopheles minimus A* | APHL01 | AminM1 | GCA_000349025.1 |
| *Anopheles nili* | ATLZ01 | Anili1 | GCA_000439205.1 |
| *Anopheles punctulatus* | JXXA01 | APwgs2 | GCA_000956255.1 |
| *Anopheles quadriannulatus A* | APCH01 | AquaS1 | GCA_000349065.1 |
| *Anopheles sinensis: Sinensis* (Korean strain) | AXCK02 | AsinS2 | GCA_000472065.2 |
| *Anopheles sinensis: China* (Chinese strain) | ATLV01 | AsinC2 | GCA_000441895.2 |

| | | | |
|---|---|---|---|
| *Anopheles stephensi*: Indian | ALPR002 | AsteI2 | GCA_000300775.2 |
| *Anopheles stephensi*: SDA-500 | APCG01 | AsteS1 | GCA_000349045.1 |

436    WGS: Whole-genome shotgun sequencing

437

438    **Table 2:** Newly described flavivirus-derived EVEs in *An. minimus* and *An. sinensis* genomes.

| Element | GenBank accession no. | Protein identity | Supercontig accession no. | Super contig length (bp) | Coordinates | | | |
|---|---|---|---|---|---|---|---|---|
| | | | | | Supercontig | | Viral genome* | |
| | | | | | Start | End | Start | End |
| *An. minimus* EVE | BK009978 | NS4-NS5 | KB664005.1 | 2,043 | 107 | 1,987 | 6,567 | 8,432 |
| *An. sinensis* EVE (Korean strain) | BK009979 | NS3 | AXCK02024744 | 2,797 | 822 | 1,613 | 5,214 | 5,822 |
| *An. sinensis* EVE (Chinese strain) | BK009980 | NS3 | ATLV01019207.1 | 7,380 | 3,578 | 4,376 | 5,214 | 5,822 |

439    * Viral genomes coordinates are based on the closest tBLASTn hits: Nienokoue virus

440    (GenBank accession no. JQ957875) for the *An. minimus* EVE and *Culex* flavivirus (GenBank

441    accession no. JQ308188) for the *An. sinensis* EVE.

442

443    **Supporting information**
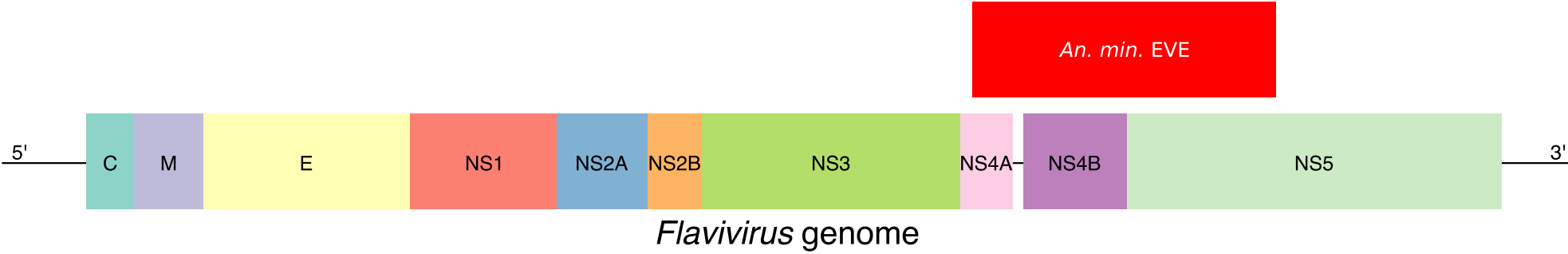
444

445    **Table S1:** Names and accession numbers of flavivirus polyprotein sequences used as queries

446    for the *Anopheles* genome screen and included in the EVE phylogenetic analyses.

447

448    **Table S2:** PCR primers and reaction conditions used to amplify EVE and RPS7 targets from

449    DNA or cDNA templates.

450

451    **Table S3:** Other genetic integrations detected in the *An. sinensis* genome supercontigs

452    containing the flavivirus-derived EVE.

453

454    **Figure S1:** Sequencing coverage of (A) *An. minimus* and (B) *An. sinensis* EVEs in published

455    RNA-seq experiments. For each experiment, the Sequence Read Archive (NCBI) accession
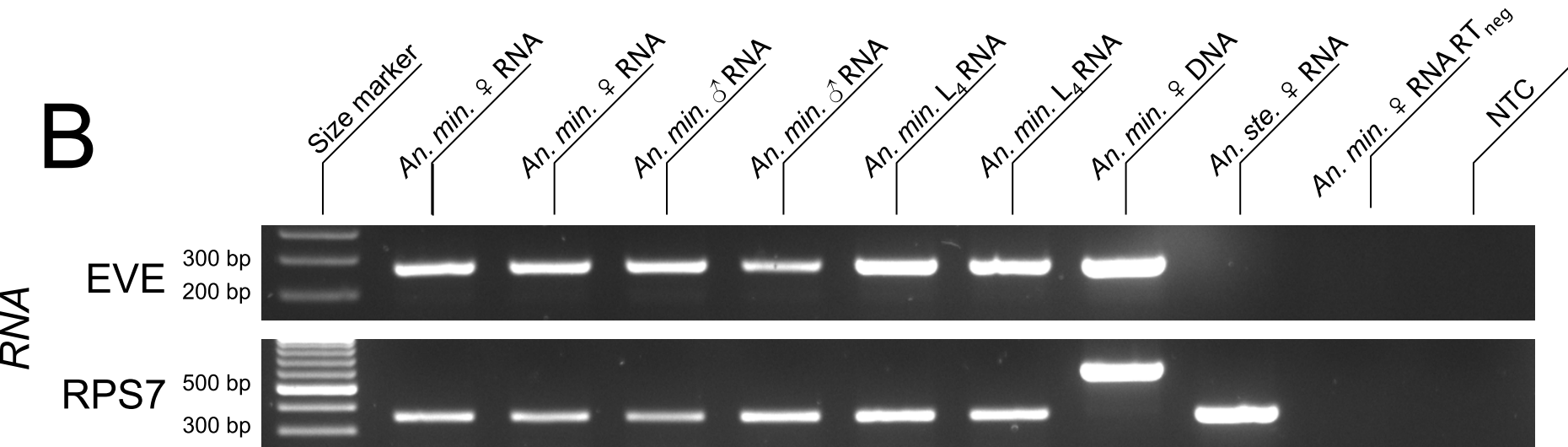
456    number is indicated above the graph.

**A**

An. min. EVE

5'  C  M  E  NS1  NS2A  NS2B  NS3  NS4A  NS4B  NS5  3'

*Flavivirus* genome

**B**

St. Louis encephalitis virus
West Nile virus
Kunjin virus
Usutu virus
Murray Valley encephalitis virus
Japanese encephalitis virus
Rocio virus
Ilheus virus
Sitiawan virus
Duck Tembusu virus
Ntaya virus
Israel turkey meningoencephalomyelitis virus
Bagaza virus
New Mapoon virus
Stratford virus
Kokobera virus
Zika virus
DENV4
DENV2
DENV3
DENV1
Nounane virus
Nhumirim virus
Barkedji virus
Ilomantsi virus
Donggang virus
Lammi virus
Chaoyang virus
Yellow fever virus
Wesselsbron virus
Yokose virus
Entebbe bat virus
Modoc virus
Tyuleniy virus
Tick-borne encephalitis virus
Langat virus
Alkhumra hemorrhagic fever virus
Powassan virus
Deer tick virus
*An. minimus* EVE
Quang Binh virus
*Culex theileri* flavivirus
*Culex* flavivirus
Palm Creek virus
Nakiwogo virus
Nienokoue virus
Hanko virus
*Aedes* flavivirus
Kamiti River virus
Cell fusing agent virus
Tamana bat virus

0.8

**A**

An. sin. EVE

5'    C    M    E    NS1    NS2A    NS2B    NS3    NS4A    NS4B    NS5    3'

*Flavivirus* genome

**B**

West Nile virus
0.99
Kunjin virus
0.98
Usutu virus
0.97
Murray Valley encephalitis virus
0.91
Japanese encephalitis virus
Rocio virus
0.97
Ilheus virus
0.97
Sitiawan virus
0.92
Duck Tembusu virus
0.97
Ntaya virus
0.95
Bagaza virus
0.85
Israel turkey meningoencephalomyelitis virus
0.97
St. Louis encephalitis virus
New Mapoon virus
1
Stratford virus
0.99
Kokobera virus
Zika virus
DENV4
DENV1
0.9
DENV3
0.97
DENV2
Chaoyang virus
0.97
Lammi virus
0.99
Donggang virus
0.96
Ilomantsi virus
Barkedji virus
1
Nhumirim virus
Nounane virus
0.92
0.84
Yellow fever virus
1
Wesselsbron virus
0.92
Entebbe bat virus
0.98
Yokose virus
Modoc virus
0.98
Tyuleniy virus
0.96
Langat virus
0.99
Tick-borne encephalitis virus
0.84
Alkhumra hemorrhagic fever virus
0.88
Powassan virus
0.98
Deer tick virus
1
*An. sinensis* EVE Korea
*An. sinensis* EVE China
Palm Creek virus
Nakiwogo virus
0.96
Quang Binh virus
0.92
*Culex theileri* flavivirus
0.97
*Culex* flavivirus
Nienokoue virus
0.85
Hanko virus
Kamiti River virus
0.92
*Aedes* flavivirus
1
Cell fusing agent virus
Tamana_bat_virus

0.6

**A**

DNA

| Size marker | An. min. ♀ DNA | An. min. ♂ DNA | An. sin. ♀ DNA | An. sin. ♂ DNA | An. ste. ♀ DNA | NTC |

EVE

1000 bp
900 bp

**B**

RNA

| Size marker | An. sin. ♀ RNA | An. sin. ♀ RNA | An. sin. ♂ RNA | An. sin. ♂ RNA | An. sin. L4 RNA | An. sin. L4 RNA | An. sin. ♀ DNA | An. ste. ♀ RNA | An. sin. ♀ RNA RT neg | NTC |

EVE

100 bp

RPS7

500 bp
300 bp