1    # Long-read assembly of the *Aedes aegypti* genome reveals the
2    # nature of heritable adaptive immunity sequences
3

4    Zachary J. Whitfield[1*], Patrick T. Dolan[1*], Mark Kunitomi[1*#], Michel Tassetto[1],
5    Matthew G. Seetin[2], Steve Oh[2], Cheryl Heiner[2], Ellen Paxinos[2], and Raul
6    Andino[1]
7
8    [1] Department of Microbiology and Immunology, University of California, 600 16[th]

9    Street, GH-S572, UCSF Box 2280, San Francisco, California 94143-2280, USA.

10    [2] Pacific Biosciences, 1305 O'Brien Drive, Menlo Park, California, 94025, USA.

11    *These authors contributed equally to this work

12    # Current address: IBM Almaden Research Center, 650 Harry Road, San Jose,

13    California, 95120-6099, USA.

14    To whom correspondence should be addressed. E-mail: raul.andino@ucsf.edu
15

16    Abstract

17         The *Aedes aegypti* mosquito is a major vector for arboviruses including

18    dengue, chikungunya and zika. Combating the spread of these viruses requires a

19    more complete understanding of mosquito-virus interactions. Recent studies

20    have implicated DNA derived from non-retroviral RNA viruses in insect immunity.

21    To better define the role and origin of these elements, we generated a high-

22    quality assembly of the *Ae. aegypti*-derived Aag2 cell line genome using single-

23    molecule, real-time sequencing technology. The new assembly improves

24    contiguity by one to two orders of magnitude with respect to previously released

25    assemblies. This improved quality enables characterization of the collection of

26    Endogenous Viral Elements (EVEs) in the mosquito genome, providing insight

27    into their integration and role in mosquito immunity. Additionally, we find a distinct

28    repertoire of EVEs present in the genomes of *Ae. aegypti* and *Ae. albopictus*,

29    suggesting the intriguing possibility that differences in EVE composition may play

30    a role in establishing vector competence.

31

32

33

**Introduction**

35

36      Mosquito transmission of arboviruses such as Dengue virus (DENV),

37   Chikungunya virus (CHIKV), and the newly emerging Zika virus causes

38   widespread and debilitating disease across the globe (Bhatt et al., 2013). The

39   primary vector of these viruses, *Aedes aegypti,* has a global tropical/subtropical

40   distribution (Kraemer et al., 2015) creating geographically isolated populations of

41   *Ae. aegypti* that have diversified over time. This genetic diversity has resulted in

42   differential competence for vectoring virus(Bennett et al., 2002). Comparative

43   genomics may explain these differences in vector competence between *Ae.*

44   *aegypti* populations. However, this comparison is limited by the repetitive nature

45   of the *Ae. aegypti* genome and the absence of an assembly based on long-read

46   sequencing capable of spanning these regions. The underlying genetic

47   contributions to vector competence are of critical importance to understanding

48   the epidemiology of both acute epidemics as well as endemic occurrences of

49   arbovirus infection.

50

51      One critical factor underlying vector competence is the insect immune

52   system (Kramer, 2016; Kramer & Ciota, 2015). However, the connection between

53   the mosquito immune system, tolerance to viral infection, and vector competence

54   at a genomic level is still unclear. At their core, insects make use of an RNAi-

55   based immune defense, where foreign viral dsRNA intermediates are recognized

56   and processed through a dicer and argonaute mediated pathway, leading to

57   cleavage of viral RNA and protection from infection (Mongelli & Saleh, 2016).  In

58   addition, mosquitoes utilize an additional RNAi pathway mediated by Piwi

59   proteins and piRNAs as an antiviral defense system (Miesen, Joosten, & van Rij,

60   2016). Previous reports have indicated piRNAs can be produced from virus

61   derived DNA sequences, and identified a set of proteins responsible for their

62   processing and maturation (Goic et al., 2016; Miesen, Girardi, & van Rij, 2015;

63   Miesen, Ivens, Buck, & van Rij, 2016; Miesen, Joosten, et al., 2016).  Given that

64    endogenous viral elements (EVEs) (Katzourakis & Gifford, 2010)  are capable of

65    producing piRNAs (Arensburger, Hice, Wright, Craig, & Atkinson, 2011), we

66    propose that rigorous EVE identification and examination in *Ae. aegypti* is central

67    to understand arbovirus/vector interaction. However, due to the highly repetitive

68    nature of the *Ae. aegypti* genome and EVEs' tendency to cluster within such

69    repetitive regions, many EVEs are likely to be missing from the current *Ae.*

70    *aegypti* assemblies, which are based on relatively short read lengths (Nene et al.,

71    2007; Vicoso & Bachtrog, 2015).

72          Accordingly, we sought to generate a *de novo* genome assembly using

73    long read deep sequencing, which is more suited to characterize the many

74    repetitive regions that make up the majority of the *Ae. aegypti* genome. The

75    resulting assembly greatly improves our understanding of the *Ae. aegypti*

76    genome, a foundational tool for studying arboviral disease, spread, and

77    prevention. As a demonstration of its utility, we use this improved assembly to

78    characterize the structure and composition of EVE-containing loci across the

79    entire *Ae. aegypti* genome.

80

81    **Sequencing and Assembly of the Aag2 Genome**

82          Current assemblies of the *Ae. aegypti* genome are based on two

83    sequencing strategies: one produced with the Illumina sequencing platform (to be

84    referred to as 'UCB') (Vicoso & Bachtrog, 2015) and one based on conventional

85    Sanger sequencing (to be referred to as 'LVP') (Nene, et al., 2007). In all

86    instances, the Liverpool strain of *Ae. aegypti* was sequenced (Table 1).  A more

87    recent study used Hi-C to further organize the sanger-based *Ae. aegypti*

88    assembly into chromosome level scaffolds (Dudchenko et al., 2017). Long-read

89    assembly will further aid these efforts by reading through and resolving large

90    repetitive regions unable to be identified by previous methods.  To this end, we

91    employed single-molecule, real-time sequence technology (Pacific Biosystems)

92    to generate long read sequences of the genome of the cell line Aag2.  Of note,

93    the *Ae. aegypti*-derived Aag2 cell line is an important, widely used tool in the

94    study of *Ae. aegypti* biology and its associated arboviruses. Here, we present the

95    draft assembly of the Aag2 genome and use it to map endogenous viral elements

96    (EVEs) genome-wide.

97

98    Approximately 76-fold coverage of the *Ae. aegypti*-based Aag2 genome

99    was achieved using the Single Molecule Real Time (SMRT) sequencing platform

100   (P6/C4 chemistry) to shotgun sequence 116 SMRT cells generating 92.7 GB of

101   sequencing data with an average read length of 15.5 kb. We used Falcon and

102   Quiver to generate a *de novo* 1.7 Gbp assembly with a contig N50 of ~1.4 Mbp.

103

104   Our assembly improves upon previous *Aedes* assemblies as measured by

105   N50, L50, and simply by the number of contigs (Table 1 and Figure 1). The long-

106   read sequencing approach enabled a substantial improvement in assembly

107   contiguity (Figure 1a).  A majority of the Aag2 assembly sequence is found on

108   contigs 10-100x longer than previous assemblies.  This increased contiguity

109   allows the mapping of numerous contigs from the initial LVP assembly to single

110   Aag2 contigs (Figure 1b and c), and makes for an overall more complete and

111   much more ordered genome assembly for use in downstream processes.  In this

112   report, we focused on examining repetitive regions of the genome, and as

113   expected, uncovered a plethora of Endogenous Viral Elements (EVEs). These

114   EVEs were found to produce *bona fide* small RNAs, and were organized into loci

115   akin to a CRISPR-like system(de Vanssay et al., 2012).

116

117   **Repetitive nature of the Aag2 genome**

118   The genome of *Ae. aegypti* was previously shown to contain high

119   proportions of repeat-DNA (Nene, et al., 2007). The *Ae. aegypti*-derived Aag2

120   genome is no different, and is comprised of almost 55% repeat sequence (Table

121   1 and Table 2). Our sequencing strategy allows more repeats to be sequenced

122   within a single read, and therefore better reflects the structure and organization

123   of these repetitive elements. Direct alignment of contigs in the Aag2 assembly

124   and those of previous *Ae. aegypti* assemblies reveal resolved rearrangements

125   and distinct repeated regions that were collapsed into single sequences in the

126    previous assemblies (Figure 1c and d).  These regions can span 10-20kb

127    (uncollapsed), illustrating the need for long read lengths to properly order the

128    vast amount of repetitive regions in the *Ae. aegypti* genome.  Of these repetitive

129    regions, over 75% is made up of transposon-derived sequence (Table 2).

130

131    **Transposon-derived sequences in the *Ae. aegypti* genome**

132        Transposable elements (TEs) play important roles in gene regulation and

133    genome evolution, providing a source of genomic variation that is a driver of

134    evolution (Gifford, Pfaff, & Macfarlan, 2013; Thompson, Macfarlan, & Lorincz,

135    2016).  TEs account for a large portion of the *Aedes aegypti* genome and our

136    long-read assembly allows us to explore the large-scale structure of these

137    prominent features (Figure 2a and b). Individual contigs of the Aag2 assembly

138    contained on average a significantly higher number of transposons than

139    previously observed in the other two assemblies (Figure 2c).

140

141        The TEs in the Aag2 genome are derived from a number of different

142    families and distributed throughout the genome (Figure 2a, d , Figure 2-Figure

143    Supplement 1). However, the distribution of TEs is not completely uniform across

144    the genome (Figure 2e).  Local density plots reveal regions where particular TE

145    classes are overrepresented, likely reflecting integration site bias and/or piRNA

146    cluster formation (discussed below).  Kimura distribution analysis (Kimura, 1980)

147    of TEs in the Aag2 genome shows a relatively recent expansion of TEs,

148    particularly LINE, LTR, and MITEs elements (Figure 2f; low Kimura scores

149    indicate TEs that are closer to the element's consensus sequence, while higher

150    scores indicated more diverged TE sequences). With our assembly a resource

151    particularly suited for repeat identification and analysis has been produced.

152

153    **Identification of EVEs**

154        Given their propensity to integrate into repetitive TE clusters (Figure

155    3b)(Parrish et al., 2015), our understanding of the EVE composition and structure

156    has been limited. We thus used our improved long-read assembly genome to

157    better define the complete set of EVEs contained within the Aag2 genome,

158    hereby called the "EVEome". Using a BLASTx-based approach (see Methods),

159    we searched for all EVE sequences within the newly assembled Aag2

160    genome.  A total of 368 EVEs were identified in the Aag2 assembly (Figure 3a).

161    We were able to detect EVEs derived from at least 8 viral families, dominated by

162    sequences derived from Rhabdoviridae, Flaviviridae, and Chuviridae (Figure 3c).

163

164         Given that TE clusters have been shown to produce piRNAs, we next

165    examined the populations of these small RNAs mapping to EVEs across the

166    entire Aag2 genome (Figure 3a, b, and d).  While only a small proportion of *Ae.*

167    *aegypti* TE elements produce piRNAs(Arensburger, et al., 2011), we observed

168    281 EVEs (73.2% of all EVEs)  each producing at least 10 (Figure 3a, red lines).

169    Characterization of these EVE piRNAs suggest they are *bona fide* piRNAs

170    because they are methylated (resistant to $\beta$–elimination), present a 5' U bias, and

171    range in size from 24 to 30 nucleotides) (data not shown). Furthermore, when

172    compared to the previous *Ae. aegypti* assemblies, ~$2x10^6$ more piRNAs could be

173    mapped to the more contiguous Aag2 assembly (Figure 3d), indicating our

174    assembly gives a more complete tool for the study of the EVE composition within

175    repetitive regions of the genome. With their propensity to produce mature

176    piRNAs, and their integration in the *Ae. aegypti* genome, we speculate that EVEs

177    comprising the "EVEome" represent a collection of adaptive immune "cassettes"

178    in *Ae. aegypti,* some of which could function anti-virally.

179

180    **Insights into the mechanism of EVE integration**

181         DNA derived from RNA viruses is produced in persistently infected

182    Drosophila cell lines (Goic et al., 2013) and in infected *Aedes albopictus*

183    mosquitoes (and multiple mosquito-derived cell culture lines) (Goic, et al., 2016).

184    This viral DNA (vDNA) synthesis depends on the activity of endogenous reverse

185    transcriptases (Goic, et al., 2016; Goic, et al., 2013). Further, sequencing of viral

186    DNA isolated from Drosophila cell lines (Goic, et al., 2013) has demonstrated

187    formation of DNA hybrids between viral DNA sequences and transposable

188    elements. EVEs are typically found nearby transposons and can be found within

189    piClusters (Figure 3a and b) (Feschotte & Gilbert, 2012; Honda & Tomonaga,

190    2016; Miesen, Joosten, et al., 2016). Together this information suggests that

191    EVEs are generated and integrated into the host genome in the context of the

192    replication cycle of endogenous transposable elements. To determine whether a

193    particular transposable element type is responsible for EVE integration, we

194    identified TEs whose coordinates in the genome directly overlap EVE sequences

195    (as called by RepeatMasker and BLASTX respectively). This approach identifies

196    mobile elements most likely responsible for genomic integration of non-retroviral

197    virus sequence.  In line with observations in *Ae. albopictus (X. G. Chen et al.,*

198    *2015)*, LTR retrotransposons were found to be greatly enriched near EVEs

199    (Figure 4-Figure Supplement 1 (i)). A similar pattern was observed when

200    classifying the nearest upstream and downstream (non-overlapping) TE

201    sequences around each EVE (Figure 4a(i)). These results further implicate LTR

202    retrotransposons in the acquisition of EVEs and indicate that the typical

203    integration sites are composed of clusters of similar LTR retrotransposons.

204

205        Strikingly, a vast majority of these LTR TEs shared the same

206    strandedness as their nearest EVE. Of 571 (non-overlapping) TEs with the same

207    strandedness as their nearest EVE, 419 were LTRs (Figure 4a(i); p-value =

208    $3.42 \times 10^{-188}$ by one-sided binomial test).  This bias is consistent with a copy-

209    choice mechanism of recombination between LTR retrotransposon sequence

210    and viral RNA leading to EVE integration, as previously proposed (Cotton,

211    Steinbiss, Yokoi, Tsai, & Kikuchi, 2016; Geuking et al., 2009). Our analysis of

212    transposons in the Aag2 genome shows LTR-retrotransposons display less

213    diversity (by Kimura Divergence score; Figure 2f), indicating that they are

214    currently (or were recently) actively replicating in the Aag2 cell line. Consistent

215    with this idea, LTR-retrotransposon transcripts and proteins are readily detected

216    in Aag2 cells (Maringer et al., 2017). We thus conclude that LTR-

217    retrotransposons are responsible for the acquisition of the majority of EVEs

218    observed in the *Ae. aegypti* genome.

219

220        Within the LTR retrotransposon family, both Ty3/gypsy and Pao Bel TEs

221    are enriched surrounding EVEs (Figure 4a iv,v). Again, this enrichment for

222    Ty3/gypsy and Pao Bel elements near EVE loci is strongest when the EVE and

223    TE are in the same orientation (p-value = $7.99 \times 10^{-23}$ and $2.00 \times 10^{-3}$

224    respectively).  The drastic reduction in associated transposons based on

225    directionality is not observed for other TE categories (Figure 4a, iii, vi). These

226    data support Ty3/gypsy (and to a lesser extent Pao Bel) as the primary

227    transposon type facilitating EVE genomic integration in *Ae.*

228    *aegypti*.   Interestingly, an association between LTR Ty3/gypsy elements and

229    integrated viral sequence has also been observed previously in plants (Lee,

230    Nolan, Watson, & Tristem, 2013; Staginnus et al., 2007), suggesting a conserved

231    mechanism for the acquisition of invading virus sequences and generation of

232    EVEs.

233

234        EVE-proximal TEs of the Ty3/gypsy and Pao Bel families can be further

235    partitioned into individual elements. Of these, many specific elements were

236    enriched for being the nearest TE to an EVE (Figure 4a(iv, v)). Interestingly,

237    EVEs derived from different virus families show different patterns of enrichment

238    for nearby TEs (Figure 4b).  Both Flaviviridae and Rhabdoviridae-derived EVEs

239    show strong enrichment for Ty3/gypsy transposable elements, while Chuviridae-

240    derived EVEs are most typically adjacent to Pao Bel elements. Of the Ty3

241    elements near Rhabdoviridae and Flaviviridae-derived EVEs, Ele152 is most

242    enriched near Flaviviridae sequences, while Ty3/gypsy Ele134 and Ele135 are

243    most closely associated with Rhabdoviridae EVE sequences.  We hypothesize

244    these data reflect the particular mobile elements whose replication overlapped

245    with a given RNA virus' replication (in both space and time).

246

247        The strong enrichment for multiple LTRs around EVEs (Figure 3b, 4a,

248    Figure 4-Figure Supplement 1), as well as the presence of EVE-derived piRNAs

249    (Figure 3a), are consistent with previous observations that EVEs integrate into

250    piRNA clusters (Parrish, et al., 2015).  We identified 469 piRNA-encoding loci

251    (piClusters) using proTRAC (Rosenkranz, Rudloff, Bastuck, Ketting, & Zischler,

252    2015; Rosenkranz & Zischler, 2012), accounting for 5,774,304 bp (0.335%) of

253    the genome. Depending on the mapping algorithm used, between 63% (bowtie)

254    and 77% (sRNAmapper.pl, see Methods) of beta-eliminated small RNAs from

255    Aag2 cells mapped to these loci. Of the identified piClusters, 66 (14.1%) have

256    EVE sequences associated with them and 65 of these piCluster-resident EVE

257    sequences act as the template for piRNAs. Of the 384 EVEs identified, 256

258    (66.7%) or 280,475 bp of the 411,239 EVE bp mapped to piClusters (68.2%,

259    Fisher's test $p < 2.2e-16$, OR=203.42). Furthermore, a vast majority of piRNAs

260    which map to EVEs are anti-sense to the EVE itself (544,429/547,014; 99.5%),

261    consistent with the idea that EVEs produce functional piRNAs and are selected

262    through evolution for their antiviral potential.

263

264          EVE/TE piRNA clusters are typically found in unidirectional orientations

265    (Figure 3b, Figure 3-Figure Supplement 1), and strikingly, some LTR/EVE

266    clusters have 'crowded out' almost all other transposons in the same region

267    (Figure 3b). We hypothesize these TE/EVE clusters are the result of canonical

268    piRNA cluster formation (Yamanaka, Siomi, & Siomi, 2014), with the rare

269    occurrence of LTR-facilitated EVE integration, resulting in localized genomic

270    expansion (due to non-random LTR integrase directed integration)(Lesbats,

271    Engelman, & Cherepanov, 2016). These data together support the notion that

272    EVEs are located at specific genomic loci associated with high piRNA production.

273

274    **EVEs in *Ae. aegypti* and *Aedes albopictus***

275          If EVEs serve as a representative record of viral infection over time, and

276    the majority are a 'recent' acquisition (based on Ty3/gypsy's low Kimura

277    Divergence scores (Figure 2-Figure Supplement 2)), we hypothesized that EVEs

278    present in two different species of mosquito would also differ (particularly given

279    the relatively rare occurrence of genome fixation) (Holmes, 2011; Katzourakis &

280    Gifford, 2010).  The *Ae. aegypti* and *Ae. albopictus* species of mosquito occupy

281  distinct (yet overlapping) regions around the globe (Kraemer, et al., 2015) and so

282  have not faced the exact same viral challenges over time.  While the EVEs

283  present in the Aag2 and LVP *Ae. aegypti*-based genomes correspond well, *Ae.*

284  *aegypti* and *Ae. albopictus* do not share any specific EVEs. However, exploring

285  the Flaviviridae family of viruses in greater detail, the viral species from which

286  these EVEs are derived do partially overlap (Figure 5a). However, the relative

287  abundance of EVEs derived from various viral species in *Ae. aegypti* and *Ae.*

288  *albopictus* differs. The lack of specific EVEs in common between the mosquito

289  genomes indicates EVE acquisition by *Ae. aegypti and Ae. albopictus* occurred

290  post-speciation, an important factor when considering any differences in vector

291  competence between these two species.

292

293      We next determined whether any particular region of the virus genome is

294  more frequently acquired and converted into EVEs. EVEs were mapped to the

295  locations of the viral ORFs from which they originate (Figure 5b).  Interestingly,

296  EVEs deriving from Flaviviridae primarily map toward the 5' end of the single

297  Flaviviral ORF, leaving a relative dearth of EVEs at the 3' end. EVEs deriving

298  from Rhabdoviridae primarily originate from the Nucleoprotein (N) and

299  Glycoprotein (G) coding sequences, with only a few originating from the RNA-

300  dependent RNA polymerase (L). The lack of EVEs mapping to the polymerase

301  may be the result of RNA expression levels, with L being the least expressed

302  gene (Conzelmann, 1998). This suggests that the template for cDNA synthesis

303  and recombination with a TE genome are viral mRNAs. The lack of EVEs

304  originating from the Phosphoprotein (P) or Matrix protein (M) is more difficult to

305  explain, potentially reflecting the localization and/or availability of the RNA

306  species for recombination. Interestingly, EVEs derived from Chuviridae primarily

307  map to the ORF of the Glycoprotein. Given the complex and diverse nature of

308  Chuviridae genomes (which include unsegmented, bi-segmented, and possibly

309  circularized negative-sense genomes)(Li et al., 2015), this pattern could also be

310  the result of RNA abundance.  Within each ORF however, there is no obvious

311    selection for EVEs from a particular location. Rather, EVEs map to regions

312    evenly distributed across their respective ORFs.

313

314        We next analyzed the Kimura divergence scores for TEs closest to (but

315    not overlapping) EVEs (both upstream and downstream; Figure 5c). When TEs

316    were grouped both by family and the family of the virus from which its closest

317    EVE was derived from, stark differences in Kimura score distribution emerge

318    (Figure 5c, Figure 5-Figure Supplement 1). Pao Bel elements near EVEs

319    originating from Chuviridae and Flaviviridae had much higher Kimura divergence

320    scores than Ty3/gypsy elements matching the same criteria.  If Kimura

321    divergence score is taken as a proxy for 'genomic age', these Pao Bel elements

322    (and potentially the EVEs they are most closely associated with) are a more

323    ancient acquisition than their Ty3/gypsy counterparts.  On the other hand, both

324    Pao Bel and Ty3/gypsy elements closest to EVEs derived from Rhabdoviridae

325    show a more uniform distribution of scores, suggesting there may have been

326    continued periods of Rhabodviral infection (and subsequent EVE acquisition by

327    transposable elements) in *Ae. aegypti*.

328

329        Recent publications have highlighted the integration of genetic material

330    from non-retroviral RNA viruses into the genome of the host during infection that

331    relies upon endogenous retro-transcriptase activity from transposons. A subset of

332    EVEs found in mammalian systems seem to be under purifying selection, which

333    suggests that they are beneficial to the host (Horie et al., 2010). A species'

334    "EVEome" could represent a built-in, yet adaptable, viral defense system. All

335    mosquito species share the same basic RNAi-based immune system. We

336    propose differences in a given specie's (or subpopulation's) "EVEome", such as

337    those between *Ae. aegypti* and *Ae. albopictus* (Figure 5a), may represent an

338    important factor contributing to inherent differences in vector competence. In

339    support of this hypothesis, the *Ae. albopictus* assembly contains an EVE that

340    clusters phylogenetically with extant DENV (Figure 5a) and has been shown to

341    be less competent at disseminating DENV in some cases (W. J. Chen, Wei, Hsu,

342    & Chen, 1993; Whitehorn et al., 2015).

343

344

345    **Discussion**

346        *In vivo* studies of mosquito immunity are a valuable, but challenging

347    approach to understanding arboviral life cycles.  The *Ae. aegypti*-based cell line

348    Aag2 provides a tool with which molecular characterization of the arboviral

349    replication in mosquitos can be accomplished on a much wider scale.  An

350    improved, fully assembled *Ae. aegypti*-derived Aag2 genome is a significant step

351    in best utilizing this cell line and furthering our understanding of *Ae. aegypti*

352    biology.

353

354        The presence of piRNA producing EVEs in the *Ae. aegypti* genome is

355    reminiscent of the CRISPR system in bacteria.  Both take advantage of the

356    invading pathogen's genetic material to create small RNAs capable of restricting

357    an invading virus' replication.  Furthermore, both end up integrated into the host's

358    genome, likely providing some level of protection against future infections.  EVEs

359    themselves can present a unique opportunity to track viral evolution and

360    historical interactions between host and virus (Holmes, 2011; Katzourakis,

361    Tristem, Pybus, & Gifford, 2007; Keckesova, Ylinen, Towers, Gifford, &

362    Katzourakis, 2009). Our Aag2 genome assembly refines our understanding of

363    EVEs in the *Ae. aegypti* genome, and the breadth of coverage by which they may

364    protect hosts against viral infections.

365

366        The finding that LTR retrotransposon elements (specifically Ty3/gypsy and

367    Pao Bel) are closest in proximity to EVEs is intriguing. A number of (non-mutually

368    exclusive) possibilities could explain this observation. It may be that Ty3/gypsy

369    and Pao Bel elements were in the 'right place at the right time' to participate in

370    acquisition of viral sequences as EVEs.  LTR retrotransposon replication and

371    virus replication must have overlapped both physically within the cell and

372    temporally within the natural history of *Ae. aegypti* evolution so as to provide the

373    opportunity for LTR-mediated EVE integration.  Template switching during

374    reverse transcription has previously been proposed to play a role in creating the

375    transposon-virus hybrids which integrate into the host genome to form EVEs

376    (Cotton, et al., 2016; Geuking, et al., 2009).  The enrichment of Pao Bel TEs near

377    Chuviridae-derived sequences and Ty3/gypsy TEs near Flaviviridae and

378    Rhabdoviridae sequences (Figure 4b) could have occurred by chance, or may

379    hint at an even deeper level of specificity directing capture of viral sequences by

380    LTRs. Possibly these TE/EVE pairs have increased sequence homology leading

381    to more frequent template switching of the reverse transcriptase (Delviks-

382    Frankenberry et al., 2011), or their subcellular localization of replication better

383    coincide. One key distinction between LTR and non-LTR retrotransposons is the

384    cellular location in which reverse transcription (RT) occurs.  LTR

385    retrotransposons undergo reverse transcription in the cytoplasm, while non-LTR

386    retrotransposons undergo RT within the nucleus (Servant & Deininger, 2015).

387    Given most RNA viruses also replicate within the cytoplasm, the opportunity for

388    template switching onto an RNA virus genome would be more readily available to

389    an LTR retrotransposon.

390

391         It is also possible that only LTR/EVE pairs were selected for after

392    integrating into the mosquito genome. The role of selection on the organization

393    and placement of EVEs in the *Ae. aegypti* genome is an intriguing one, and many

394    EVEs have integrated into piRNA clusters from which *bona fide* piRNAs are

395    produced (Figure 3a,b).  However, a lack of synteny between EVEs in *Ae.*

396    *aegypti* and *Ae. albopictus* suggests their acquisition likely took place post-

397    speciation.  Given this, along with LTRs being most active most recently (among

398    TEs) in the Aag2 cell line (Figure 2f) (Maringer, et al., 2017), many EVEs in the

399    *Aedes aegypti* genome likely represent relatively recent acquisitions. While an

400    informative observation, this also limits the ability to date these EVEs with

401    molecular-clock based techniques*,* and thus properly study selection. However,

402  the evidence so far suggests that some EVE sequences (and their organization

403  in the genome) are being maintained evolutionarily.

404

405      As mentioned above, only the *Ae. albopictus* genome contains EVEs

406  which cluster with extant DENV, and it also happens to be a less suitable vector

407  for the virus (as compared to *Ae. aegypti*).  We hypothesize these particular EVE

408  insertions may provide a buffer against DENV infection. *Ae. albopictus* is still

409  susceptible to DENV infection, thus the DENV-derived EVE does not provide

410  complete protection against DENV infection.  However, as piRNAs are found at

411  various levels within different tissues of the mosquito (Akbari et al., 2013), EVE-

412  derived piRNAs likely confer differing levels of viral resistance depending on the

413  tissue in question. Thus, EVE-derived piRNAs could play a role in keeping viral

414  infections in check at an organismal level, maybe indefinitely in the case of

415  persistent infections. The hypothesis that many EVE-derived piRNAs in *Ae.*

416  *aegypti* are functional (and their organization in the genome has been preserved)

417  is further supported by the shared orientation of LTR retrotransposons and their

418  nearby EVEs.  In this arrangement, precursor transcripts originating from

419  EVE/TE-containing piRNA clusters encode anti-viral piRNA sequence with the

420  same directionality they contain anti-TE piRNA sequence.  The resulting anti-

421  sense piRNAs can then go on to silence their complementary (viral or TE) RNA.

422

423      A solid foundation with which to study the genetic factors contributing to

424  vector competence is of utmost importance. With this in mind, we generated a

425  highly contiguous assembly of the *Aedes aegypti* cell line, Aag2. With this long-

426  read assembly, we then identified nearly the entire set of endogenous viral

427  elements and their surrounding genomic context in the Aag2 cell line at a

428  genome-wide scale. Uncovering the genomic context of this EVE-derived piRNA

429  system in mosquitos provides the foundation for future studies on the role of

430  EVEs in vector competence.  The potential to manipulate a heritable, anti-viral

431  system opens up new avenues to understand the complexities of the insect

432  immune system and work to prevent spread of viral disease dependent on such

433    insect vectors.

434

## Acknowledgements

437

**MATERIALS AND METHODS**

**Cells culture**

*Aedes aegypti* Aag2 *(Lan & Fallon, 1990; Peleg, 1968)* cells were cultured at 28 °C without $CO_2$ in Schneider's *Drosophila* medium (GIBCO-Invitrogen), supplemented with 10% heat-inactivated fetal bovine serum (FBS), 1X non-essential amino acids (NEAA, UCSF Cell Culture Facility, 100X stock is 0.1 µM filtered, 10 mM each of Glycine, L-Alanine, L-Asparagine, L-Aspartic acid, L-Glutamic Acid, L-Proline, and L-Serine in de-ionized water), and 1X Penicillin-Streptomycin-Glutamine (Pen/Strep/G, 100X = 10,000 units of penicillin, 10,000 µg of streptomycin, and 29.2 mg/ml of L-glutamine, Gibco).

**DNA sequencing**

Aag2 cells were grown in T-150 Flasks until ~80% confluent. Cells were then washed with dPBS twice and scrapped off in dPBS + 10 µg/ml RNase A (ThermoFisher). Genomic DNA (gDNA) was extracted from ~10^8 Aag2 cells using the QIAamp DNA Mini Kit according to the manufacturer's instructions with the optional RNase A treatment.  Aag2 gDNA was re-suspended in 10mM Tris pH8, and the quality and quantity of the sample was assessed using the Agilent DNA12000 kit and 2100 Bioanalyzer system (Agilent Technologies), as well as the Qubit dsDNA Broad Range assay kit and Qubit Fluorometer (Thermo Fisher) and visualized by gel electrophoresis (1% TBE gel).   After purification and quality control, a total of 130 ug of DNA was available for library preparation and sequencing.

SMRTbell libraries were prepared using Pacific Biosciences' Template Prep Kit 1.0 (PacBio) and a slightly modified version of the Pacific Biosciences' protocol, "Procedure & Checklist - 20-kb Template Preparation Using BluePippin Size-Selection System (15-kb Size Cutoff)".   Specifically, 52.5ug of gDNA were hydrodynamically sheared to target sizes of 30kb (26 µg) and 35 kb (26 µg) using the Megaruptor® (Diogenode) with long hydropores according to the manufacturer's protocols.  Size distributions of the final sheared gDNA were verified by pulse field electrophoresis of a 100ng sub-aliquot through 0.75%

469 agarose using the Pippin Pulse (Sage Science), run according to the

470 manufacturer's "10-48 kb protocol" for 16 hrs. The two sheared samples were

471 then pooled, for a total of 37ug sheared DNA to be used as input into SMRTbell

472 preparation. Sheared DNA was subjected to DNA damage repair and ligated to

473 SMRTbell adapters. Following ligation, extraneous DNA was digested with exo-

474 nucleases and the resulting SMRTbell library was cleaned and concentrated with

475 AMPure PB beads (Pacific Biosciences). A total of 20.5ug of library was available

476 for size selection.

477 Approximately half (10ug) of the SMRTbell pooled SMRTbell library was

478 size-selected using the BluePippin System (Sage Science) using a 15 kb cutoff

479 and 0.75% agarose cassettes. To obtain longer read lengths, an additional 5ug

480 of the library was selected using a 17kb cutoff.

481 Library quality and quantity were assessed using the Agilent 12000 DNA

482 Kit and 2100 Bioanalyzer System (Agilent Technologies), as well as the Qubit

483 dsDNA Broad Range Assay kit and Qubit Fluorometer (Thermo Fisher). An

484 additional DNA Damage Repair step and AMPure bead cleanup were included

485 after size-selection of the libraries.

486 Annealed libraries were then bound to DNA polymerases using 3nM of the

487 SMRTbell library and 3X excess DNA polymerase at a concentration of 9nM

488 using Pacific Biosciences DNA/Polymerase Binding Kit 1.0, (Pacific Biosciences).

489 Bound libraries were sequenced on the Pacific Biosciences RSII using P6/C4

490 chemistry (PacBio), magnetic bead loading (PacBio) and 6 hour collection times.

491 84 SMRTcells of the > 15 kb library were loaded at concentrations of 75-100 pM

492 on-plate. 32 SMRTcells of the > 17 kb library was prepared separately and

493 loaded at on-plate concentrations of 40 pM and 60 pM. These 116 SMRTcells

494 generated 92.7 GB of sequencing data, which resulted in approximately 76X

495 coverage of the Aag2 genome. Average raw read length of 15.5KB, with average

496 sub-reads length of 13.2kb. Assembly was performed using Quiver/FALCON

497

498 **Genome assembly statistics**

499 Basic statistics (e.g. Size, Gaps, N50, L50, # contigs) for each genome

500 analyzed was produced using Quast (Gurevich, Saveliev, Vyahhi, & Tesler,

501 2013).

502 As a complementary approach Benchmarking sets of Universal Single-

503 Copy Orthologs (BUSCO) was also run using the Arthropod dataset in order to

504 assess the completeness of genome assembly. Of the 2675 BUSCO groups

505 searched only 81 were missing from the Aag2 assembly, indicating good

506 assembly completeness. Of the 2315 BUSCOs found only 279 of them were

507 annotated as fragmented, emphasizing the continuity of the assembly.

508

509 **Repeat Identification and Kimura Divergence**

510 In order to *de novo* identify and classify novel repetitive elements from the

511 Aag2 genome, RepeatModeler was run on the assembled genome using

512 standard parameters.  Outputs from RepeatModeler were cross-referenced with

513 annotated entries for Aedes aegypti from TEfam. All entries from RepeatModeler

514 that were >80% identical to TEfam entries were discarded as redundant. This

515 combined annotated and de novo identified list of repeat elements was used to

516 identify the genome wide occurrences of repeats using RepeatMasker using

517 standard parameters.

518 Kimura scores and corresponding alignment information were extracted

519 from the ".align" file as output by RepeatMasker. This information was used to

520 make the stacked plot in figure 2 using R (version 3.30) and the ggplot2 package.

521 Information from the header lines was then used to match Kimura

522 divergence score with the appropriate EVE-proximal TEs based on contig, TE

523 name, start point, and end point. TEs whose coordinates did not exactly match

524 output of the align file were not used in the Kimura analysis. The violin plot of

525 Kimura divergence scores was plotted using R (version 3.3.0) and the ggplot2

526 plugin.

527

528 Citations:

529  Smit, AFA, Hubley, R. *RepeatModeler Open-1.0*.

530  2008-2015 <http://www.repeatmasker.org>

531  Smit, AFA, Hubley, R & Green, P. *RepeatMasker Open-4.0*.

532  2013-2015 <http://www.repeatmasker.org>.

533  Dr. Zhijian Jake Tu. TEfam. <http://tefam.biochem.vt.edu/tefam/index.php>

534

**EVE identification**

536  Identification of EVEs was achieved using standalone Blast+ (Altschul,

537  Gish, Miller, Myers, & Lipman, 1990). Blast Searches were run using the Blastx

538  command specifying the genome as the query and a refseq library composed of

539  the ssRNA and dsRNA viral protein-coding sequences from the NCBI genomes

540  as the database. The E-value threshold was set at 10-6.

541  The EVE with the lower E-value was chosen for further analysis to predict

542  EVEs that overlapped. Several Blast hits to viral protein genes were identified as

543  artifacts because of their homology to eukaryotic genes (e.g. closteroviruses

544  encode an Hsp70 homologue). These artifacts were filtered by hand.

545

**Identification of LTR enrichment near EVEs**

547  Separate BED files containing all TEs in the Aag2 assembly and all EVEs

548  in the Aag2 assembly were used as input to Bedtools (*bedtools closest* command

549  using the *–io* flag, and *–id* or *-iu*) to find the single closest non-overlapping TE to

550  each EVE (both upstream and downstream).

551  An in-house script compiled these two output files together and filtered

552  them for the TE content of interest. TE categories (subclass, family, element)

553  were assigned by RepeatMasker. Enrichment was compared to the prevalence

554  of the TE element genome wide based on a one-sided binomial test. Stacked

555  histograms were produced based on TE categories as found in Figure 3. The

556  legend lists (up to) the 10 most prevalent TE elements of TE/EVE pairs in the

557  same orientation. Plots were produced using Python (version 2.7.6) with the

558  pandas and matplotlib plugins.

559

**Classification of nearest TE to EVEs by virus taxonomy**

Taxonomy categories for viruses from which each EVEs derived were assigned using an in-house script. Assignments were made based on NCBI's taxonomy database (ftp://ftp.ncbi.nih.gov/pub/taxonomy/), with the following additional annotations by hand.

*Virus species; assigned family*:

*Wuhan Mosquito Virus 8; Chuviridae*

*Wuchang Cockraoch Virus 3; Chuviridae*

*Lishi Spider Virus 1; Chuviridae*

*Shayang Fly Virus 1; Chuviridae*

*Wenzhou Crab Virus 2; Chuviridae*

*Bole Tick Virus 2; Rhabdoviridae*

*Shayang Fly Virus 2; Rhabdoviridae*

*Wuhan Ant Virus; Rhabdoviridae*

*Wuhan Fly Virus 2; Rhabdoviridae*

*Wuhan House Fly Virus 1; Rhabdoviridae*

*Wuhan Mosquito Virus 9; Rhabdoviridae*

*Yongjia Tick Virus 2; Rhabdoviridae*

*Cilv-C; Virgaviridae*

*Citrus leprosis virus C; Virgaviridae*

*Blueberry necrotic ring blotch virus; Virgaviridae*

*Wutai Mosquito Virus; Bunyaviridae*

Heat maps were produced using the Seaborn plugin for python. Only TEs with >=10% proportion in at least one sample (Flaviviridae, Chuviridae, or Rhabdoviridae) are shown. Color was assigned based on proportion of TE element/family in each viral category.

591    Enrichment was scored as above using a one-sided binomial test

592    (significant is p-value < 0.0001).

593

594    **Small RNA bioinformatics**

595    Adaptors were trimmed using Cutadapt

596    (http://dx.doi.org/10.14806/ej.17.1.200) using the --discard-untrimmed and -m 19

597    flags to discard reads without adaptors and below 19 nt in length. Reads were

598    mapped using bowtie (Langmead, Trapnell, Pop, & Salzberg, 2009)using the –v

599    1 flag. Read distance overlaps were generated by viROME (Watson, Schnettler,

600    & Kohl, 2013). Sequence biases were determined by Weblogo (Crooks, Hon,

601    Chandonia, & Brenner, 2004).

602    **piCluster Analysis**

603    piClusters were identified using PROtrac (Rosenkranz & Zischler, 2012)

604    based on mapping with positions for beta-eliminated small RNAs libraries from

605    Aag2 cells from sRNAmapper.pl. Based on these predictions, visualizations of

606    clusters were produced using EasyFig (Sullivan, Petty, & Beatson, 2011) for

607    visualization of TEs and R for comparison of TEs, piRNA abundance and EVE

608    positions.

609    **Sequence alignment and phylogenetic analysis**

610    For phylogenetic analysis of Flaviviridae, polyprotein sequences from 61

611    members of the Flaviviridae family were aligned with MUSCLE (Edgar, 2004) and

612    a maximum likelihood tree was generated with FastTree (Price, Dehal, & Arkin,

613    2009) using the generalised time reversible substitution model ("-gtr").  Trees

614    were visualized and annotated with ggtree (DOI: 10.1111/2041-210X.12628).

615

616    **EVE coverage**

617    Base R (version 3.3.0) was used to show regions individual EVEs span on

618    the indicated viral family (and protein). EVE length is a function of the percentage

619    of the respective ORF from which it derives.

620

621 **Statistics summary**

622     Enrichment for TE elements near EVEs (Figure 3B, 3F) was determined

623 with a one-sided binomial test (alternative hypothesis 'greater'). Enrichment of

624 EVEs in piClusters was determined by Fisher's Test. Difference in Kimura

625 Divergence distributions (not necessarily normally distributed) of TEs near EVEs

626 vs total EVE populations (Figure 4C) was determined by Kolmogorov-Smirnov

627 tests.

628

629 **Code availability**

630     The code used to generate the datasets used for visualization have been

631 provided.

632 **Data availability**

633     The Aag2 genome (v 1.00) is available through VectorBase

634 (https://www.vectorbase.org/organisms/aedes-aegypti/aag2/aag2).

635 Main datasets produced during this work have been provided in excel format.

636

## REFERENCES

Akbari, O. S., Antoshechkin, I., Amrhein, H., Williams, B., Diloreto, R., Sandler, J., & Hay, B. A. (2013). The developmental transcriptome of the mosquito Aedes aegypti, an invasive species and major arbovirus vector. *G3 (Bethesda), 3*(9), 1493-1509. doi: 10.1534/g3.113.006742

g3.113.006742 [pii]

Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *J Mol Biol, 215*(3), 403-410. doi: 10.1016/S0022-2836(05)80360-2

S0022-2836(05)80360-2 [pii]

Arensburger, P., Hice, R. H., Wright, J. A., Craig, N. L., & Atkinson, P. W. (2011). The mosquito Aedes aegypti has a large genome size and high transposable element load but contains a low proportion of transposon-specific piRNAs. *BMC Genomics, 12*, 606. doi: 10.1186/1471-2164-12-606

1471-2164-12-606 [pii]

Bennett, K. E., Olson, K. E., Munoz Mde, L., Fernandez-Salas, I., Farfan-Ale, J. A., Higgs, S., . . . Beaty, B. J. (2002). Variation in vector competence for dengue 2 virus among 24 collections of Aedes aegypti from Mexico and the United States. *Am J Trop Med Hyg, 67*(1), 85-92.

Bhatt, S., Gething, P. W., Brady, O. J., Messina, J. P., Farlow, A. W., Moyes, C. L., . . . Hay, S. I. (2013). The global distribution and burden of dengue. *Nature, 496*(7446), 504-507. doi: 10.1038/nature12060

nature12060 [pii]

Chen, W. J., Wei, H. L., Hsu, E. L., & Chen, E. R. (1993). Vector competence of Aedes albopictus and Ae. aegypti (Diptera: Culicidae) to dengue 1 virus on Taiwan: development of the virus in orally and parenterally infected mosquitoes. *J Med Entomol, 30*(3), 524-530.

Chen, X. G., Jiang, X., Gu, J., Xu, M., Wu, Y., Deng, Y., . . . James, A. A. (2015). Genome sequence of the Asian Tiger mosquito, Aedes albopictus, reveals insights into its biology, genetics, and evolution. *Proc Natl Acad Sci U S A, 112*(44), E5907-5915. doi: 10.1073/pnas.1516410112

1516410112 [pii]

Conzelmann, K. K. (1998). Nonsegmented negative-strand RNA viruses: genetics and manipulation of viral genomes. *Annu Rev Genet, 32*, 123-162. doi: 10.1146/annurev.genet.32.1.123

Cotton, J. A., Steinbiss, S., Yokoi, T., Tsai, I. J., & Kikuchi, T. (2016). An expressed, endogenous Nodavirus-like element captured by a retrotransposon in the genome of the plant parasitic nematode Bursaphelenchus xylophilus. *Sci Rep, 6*, 39749. doi: 10.1038/srep39749

srep39749 [pii]

Crooks, G. E., Hon, G., Chandonia, J. M., & Brenner, S. E. (2004). WebLogo: a sequence logo generator. *Genome Res, 14*(6), 1188-1190. doi: 10.1101/gr.849004

14/6/1188 [pii]

de Vanssay, A., Bouge, A. L., Boivin, A., Hermant, C., Teysset, L., Delmarre, V., . . . Ronsseray, S. (2012). Paramutation in Drosophila linked to emergence of a piRNA-producing locus. *Nature, 490*(7418), 112-115. doi: 10.1038/nature11416

nature11416 [pii]

686  Delviks-Frankenberry, K., Galli, A., Nikolaitchik, O., Mens, H., Pathak, V. K., & Hu, W. S.
687      (2011). Mechanisms and factors that influence high frequency retroviral
688      recombination. *Viruses, 3*(9), 1650-1680. doi: 10.3390/v3091650
689  viruses-03-01650 [pii]
690  Dudchenko, O., Batra, S. S., Omer, A. D., Nyquist, S. K., Hoeger, M., Durand, N. C., . . .
691      Aiden, E. L. (2017). De novo assembly of the Aedes aegypti genome using Hi-C
692      yields chromosome-length scaffolds. *Science*. doi: eaal3327 [pii]
693  10.1126/science.aal3327
694  science.aal3327 [pii]
695  Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high
696      throughput. *Nucleic Acids Res, 32*(5), 1792-1797. doi: 10.1093/nar/gkh340
697  32/5/1792 [pii]
698  Feschotte, C., & Gilbert, C. (2012). Endogenous viruses: insights into viral evolution and
699      impact on host biology. *Nat Rev Genet, 13*(4), 283-296. doi: 10.1038/nrg3199
700  nrg3199 [pii]
701  Geuking, M. B., Weber, J., Dewannieux, M., Gorelik, E., Heidmann, T., Hengartner, H., .
702      . . Hangartner, L. (2009). Recombination of retrotransposon and exogenous RNA
703      virus results in nonretroviral cDNA integration. *Science, 323*(5912), 393-396. doi:
704      10.1126/science.1167375
705  323/5912/393 [pii]
706  Gifford, W. D., Pfaff, S. L., & Macfarlan, T. S. (2013). Transposable elements as genetic
707      regulatory substrates in early development. *Trends Cell Biol, 23*(5), 218-226. doi:
708      10.1016/j.tcb.2013.01.001
709  S0962-8924(13)00003-2 [pii]
710  Goic, B., Stapleford, K. A., Frangeul, L., Doucet, A. J., Gausson, V., Blanc, H., . . .
711      Saleh, M. C. (2016). Virus-derived DNA drives mosquito vector tolerance to
712      arboviral infection. *Nat Commun, 7*, 12410. doi: 10.1038/ncomms12410
713  ncomms12410 [pii]
714  Goic, B., Vodovar, N., Mondotte, J. A., Monot, C., Frangeul, L., Blanc, H., . . . Saleh, M.
715      C. (2013). RNA-mediated interference and reverse transcription control the
716      persistence of RNA viruses in the insect model Drosophila. *Nat Immunol, 14*(4),
717      396-403. doi: 10.1038/ni.2542
718  ni.2542 [pii]
719  Gurevich, A., Saveliev, V., Vyahhi, N., & Tesler, G. (2013). QUAST: quality assessment
720      tool for genome assemblies. *Bioinformatics, 29*(8), 1072-1075. doi:
721      10.1093/bioinformatics/btt086
722  btt086 [pii]
723  Holmes, E. C. (2011). The evolution of endogenous viral elements. *Cell Host Microbe,*
724      *10*(4), 368-377. doi: 10.1016/j.chom.2011.09.002
725  S1931-3128(11)00285-X [pii]
726  Honda, T., & Tomonaga, K. (2016). Endogenous non-retroviral RNA virus elements
727      evidence a novel type of antiviral immunity. *Mob Genet Elements, 6*(3),
728      e1165785. doi: 10.1080/2159256X.2016.1165785
729  1165785 [pii]
730  Horie, M., Honda, T., Suzuki, Y., Kobayashi, Y., Daito, T., Oshida, T., . . . Tomonaga, K.
731      (2010). Endogenous non-retroviral RNA virus elements in mammalian genomes.
732      *Nature, 463*(7277), 84-87. doi: 10.1038/nature08695
733  nature08695 [pii]
734  Katzourakis, A., & Gifford, R. J. (2010). Endogenous viral elements in animal genomes.
735      *PLoS Genet, 6*(11), e1001191. doi: 10.1371/journal.pgen.1001191

736    Katzourakis, A., Tristem, M., Pybus, O. G., & Gifford, R. J. (2007). Discovery and
737        analysis of the first endogenous lentivirus. *Proc Natl Acad Sci U S A, 104*(15),
738        6261-6265. doi: 0700471104 [pii]
739    10.1073/pnas.0700471104
740    Keckesova, Z., Ylinen, L. M., Towers, G. J., Gifford, R. J., & Katzourakis, A. (2009).
741        Identification of a RELIK orthologue in the European hare (Lepus europaeus)
742        reveals a minimum age of 12 million years for the lagomorph lentiviruses.
743        *Virology, 384*(1), 7-11. doi: 10.1016/j.virol.2008.10.045
744    S0042-6822(08)00718-6 [pii]
745    Kimura, M. (1980). A simple method for estimating evolutionary rates of base
746        substitutions through comparative studies of nucleotide sequences. *J Mol Evol,*
747        *16*(2), 111-120.
748    Kraemer, M. U., Sinka, M. E., Duda, K. A., Mylne, A. Q., Shearer, F. M., Barker, C. M., . .
749        . Hay, S. I. (2015). The global distribution of the arbovirus vectors Aedes aegypti
750        and Ae. albopictus. *Elife, 4*, e08347. doi: 10.7554/eLife.08347
751    Kramer, L. D. (2016). Complexity of virus-vector interactions. *Curr Opin Virol, 21*, 81-86.
752        doi: S1879-6257(16)30104-3 [pii]
753    10.1016/j.coviro.2016.08.008
754    Kramer, L. D., & Ciota, A. T. (2015). Dissecting vectorial capacity for mosquito-borne
755        viruses. *Curr Opin Virol, 15*, 112-118. doi: 10.1016/j.coviro.2015.10.003
756    S1879-6257(15)00153-4 [pii]
757    Lan, Q., & Fallon, A. M. (1990). Small heat shock proteins distinguish between two
758        mosquito species and confirm identity of their cell lines. *Am J Trop Med Hyg,*
759        *43*(6), 669-676.
760    Langmead, B., Trapnell, C., Pop, M., & Salzberg, S. L. (2009). Ultrafast and memory-
761        efficient alignment of short DNA sequences to the human genome. *Genome Biol,*
762        *10*(3), R25. doi: 10.1186/gb-2009-10-3-r25
763    gb-2009-10-3-r25 [pii]
764    Lee, A., Nolan, A., Watson, J., & Tristem, M. (2013). Identification of an ancient
765        endogenous retrovirus, predating the divergence of the placental mammals.
766        *Philos Trans R Soc Lond B Biol Sci, 368*(1626), 20120503. doi:
767        10.1098/rstb.2012.0503
768    rstb.2012.0503 [pii]
769    Lesbats, P., Engelman, A. N., & Cherepanov, P. (2016). Retroviral DNA Integration.
770        *Chem Rev, 116*(20), 12730-12757. doi: 10.1021/acs.chemrev.6b00125
771    Li, C. X., Shi, M., Tian, J. H., Lin, X. D., Kang, Y. J., Chen, L. J., . . . Zhang, Y. Z. (2015).
772        Unprecedented genomic diversity of RNA viruses in arthropods reveals the
773        ancestry of negative-sense RNA viruses. *Elife, 4*. doi: 10.7554/eLife.05378
774    Maringer, K., Yousuf, A., Heesom, K. J., Fan, J., Lee, D., Fernandez-Sesma, A., . . .
775        Davidson, A. D. (2017). Proteomics informed by transcriptomics for
776        characterising active transposable elements and genome annotation in Aedes
777        aegypti. *BMC Genomics, 18*(1), 101. doi: 10.1186/s12864-016-3432-5
778    10.1186/s12864-016-3432-5 [pii]
779    Miesen, P., Girardi, E., & van Rij, R. P. (2015). Distinct sets of PIWI proteins produce
780        arbovirus and transposon-derived piRNAs in Aedes aegypti mosquito cells.
781        *Nucleic Acids Res, 43*(13), 6545-6556. doi: 10.1093/nar/gkv590
782    gkv590 [pii]
783    Miesen, P., Ivens, A., Buck, A. H., & van Rij, R. P. (2016). Small RNA Profiling in
784        Dengue Virus 2-Infected Aedes Mosquito Cells Reveals Viral piRNAs and Novel
785        Host miRNAs. *PLoS Negl Trop Dis, 10*(2), e0004452. doi:
786        10.1371/journal.pntd.0004452

787    PNTD-D-15-01748 [pii]

788    Miesen, P., Joosten, J., & van Rij, R. P. (2016). PIWIs Go Viral: Arbovirus-Derived

789           piRNAs in Vector Mosquitoes. *PLoS Pathog, 12*(12), e1006017. doi:

790           10.1371/journal.ppat.1006017

791    PPATHOGENS-D-16-02114 [pii]

792    Mongelli, V., & Saleh, M. C. (2016). Bugs Are Not to Be Silenced: Small RNA Pathways

793           and Antiviral Responses in Insects. *Annu Rev Virol, 3*(1), 573-589. doi:

794           10.1146/annurev-virology-110615-042447

795    Nene, V., Wortman, J. R., Lawson, D., Haas, B., Kodira, C., Tu, Z. J., . . . Severson, D.

796           W. (2007). Genome sequence of Aedes aegypti, a major arbovirus vector.

797           *Science, 316*(5832), 1718-1723. doi: 1138878 [pii]

798    10.1126/science.1138878

799    Parrish, N. F., Fujino, K., Shiromoto, Y., Iwasaki, Y. W., Ha, H., Xing, J., . . . Tomonaga,

800           K. (2015). piRNAs derived from ancient viral processed pseudogenes as

801           transgenerational sequence-specific immune memory in mammals. *RNA, 21*(10),

802           1691-1703. doi: 10.1261/rna.052092.115

803    rna.052092.115 [pii]

804    Peleg, J. (1968). Growth of arboviruses in monolayers from subcultured mosquito

805           embryo cells. *Virology, 35*(4), 617-619.

806    Price, M. N., Dehal, P. S., & Arkin, A. P. (2009). FastTree: computing large minimum

807           evolution trees with profiles instead of a distance matrix. *Mol Biol Evol, 26*(7),

808           1641-1650. doi: 10.1093/molbev/msp077

809    msp077 [pii]

810    Rosenkranz, D., Rudloff, S., Bastuck, K., Ketting, R. F., & Zischler, H. (2015). Tupaia

811           small RNAs provide insights into function and evolution of RNAi-based

812           transposon defense in mammals. *RNA, 21*(5), 911-922. doi:

813           10.1261/rna.048603.114

814    rna.048603.114 [pii]

815    Rosenkranz, D., & Zischler, H. (2012). proTRAC--a software for probabilistic piRNA

816           cluster detection, visualization and analysis. *BMC Bioinformatics, 13*, 5. doi:

817           10.1186/1471-2105-13-5

818    1471-2105-13-5 [pii]

819    Servant, G., & Deininger, P. L. (2015). Insertion of Retrotransposons at Chromosome

820           Ends: Adaptive Response to Chromosome Maintenance. *Front Genet, 6*, 358.

821           doi: 10.3389/fgene.2015.00358

822    Staginnus, C., Gregor, W., Mette, M. F., Teo, C. H., Borroto-Fernandez, E. G., Machado,

823           M. L., . . . Schwarzacher, T. (2007). Endogenous pararetroviral sequences in

824           tomato (Solanum lycopersicum) and related species. *BMC Plant Biol, 7*, 24. doi:

825           1471-2229-7-24 [pii]

826    10.1186/1471-2229-7-24

827    Sullivan, M. J., Petty, N. K., & Beatson, S. A. (2011). Easyfig: a genome comparison

828           visualizer. *Bioinformatics, 27*(7), 1009-1010. doi: 10.1093/bioinformatics/btr039

829    btr039 [pii]

830    Thompson, P. J., Macfarlan, T. S., & Lorincz, M. C. (2016). Long Terminal Repeats:

831           From Parasitic Elements to Building Blocks of the Transcriptional Regulatory

832           Repertoire. *Mol Cell, 62*(5), 766-776. doi: 10.1016/j.molcel.2016.03.029

833    S1097-2765(16)30012-0 [pii]

834    Vicoso, B., & Bachtrog, D. (2015). Numerous transitions of sex chromosomes in Diptera.

835           *PLoS Biol, 13*(4), e1002078. doi: 10.1371/journal.pbio.1002078

836    PBIOLOGY-D-14-02861 [pii]

837    Watson, M., Schnettler, E., & Kohl, A. (2013). viRome: an R package for the
838            visualization and analysis of viral small RNA sequence datasets. *Bioinformatics,*
839            *29*(15), 1902-1903. doi: 10.1093/bioinformatics/btt297
840    btt297 [pii]
841    Whitehorn, J., Kien, D. T., Nguyen, N. M., Nguyen, H. L., Kyrylos, P. P., Carrington, L.
842            B., . . . Simmons, C. P. (2015). Comparative Susceptibility of Aedes albopictus
843            and Aedes aegypti to Dengue Virus Infection After Feeding on Blood of Viremic
844            Humans: Implications for Public Health. *J Infect Dis, 212*(8), 1182-1190. doi:
845            10.1093/infdis/jiv173
846    jiv173 [pii]
847    Yamanaka, S., Siomi, M. C., & Siomi, H. (2014). piRNA clusters and open chromatin
848            structure. *Mob DNA, 5*, 22. doi: 10.1186/1759-8753-5-22
849    1759-8753-5-22 [pii]
850
851

852

853    **LEGENDS FOR MAIN FIGURES**

854

855    **Figure 1. Contiguity of the *Aedes aegypti* genome is drastically improved in**
856    **the Aag2 assembly.** (A) Histogram of contig length vs. total amount of sequence
857    contained in each bin. The Aag2 assembly achieved the largest contig sizes (by
858    an order of magnitude) compared to previous *Aedes aegypti*-derived assemblies.
859    This large contig size also resulted in more overall sequence information/number
860    of bases. (B) Boxplots indicating number of contigs aligned between LVP
861    (Sanger) and Aag2 (PacBio). When aligned to each other, more contigs from
862    LVP are aligned to larger Aag2 contigs than vice versa. (C) Dot plot alignment of
863    multiple contigs from LVP (each denoted by a different color) to a single contig
864    from Aag2 (sequenced by Pac Bio). Expanded repeats in the Aag2 assembly
865    can be seen where the same LVP contig aligns to multiple Aag2 locations. In
866    some instances, portions of multiple LVP contigs align to a single locus within the
867    Aag2 contig (vertical 'spikes' of alignment). (D) More detailed dot plot alignments
868    between Aag2 and LVP assemblies. Regions in the LVP assembly are expanded
869    within the Aag2 assembly at numerous loci. Regions corresponding to putative
870    piClusters in the Aag2 assembly are highlighted with grey shaded boxes. Each
871    panel is labeled with the contig and position of the predicted piCluster shown.

872

873    **Figure 2. The repeat landscape of the *Aedes aegypti* genome is**
874    **predominantly made up of transposable elements.** (A) Circular plot of the
875    Aag2 assembly, showing the 10 largest contigs (black rectangles; ordered by
876    size) and transposable elements (circles; colored by TE class). Transposable
877    elements are prevalent throughout the congtigs (and the entire *aegypti* genome).
878    Rectangles representing contigs are staggered to indicate relative contig size. (B)
879    Density plot showing the distribution of TE content for contigs in the assembly.
880    Red line indicates the genome-wide TE content. (C) Density plot comparing TE
881    counts binned by TEs per contig between the Aag2, LVP(Sanger), and UCB
882    (Illumina) assemblies. The Aag2 assembly has significantly more TEs per single
883    contig than previous assemblies. (D) Pie chart representing TE class
884    representation in the Aag2 genome. Further detail can be found in Table 2. (E)
885    Stacked area plot showing the relative density of TEs along an example contig.
886    Local proportions are based on a window size of 25 kbps. (F) Stacked histogram
887    of Kimura divergence for classes of TEs found in the Aag2 assembly, expressed
888    as a function of percentage of the genome. A relatively recent expansion/active
889    phase of LTRs is evident (increase in LTRs at low Kimura divergence scores).
890    Kimura divergence scores are based on the accumulated mutations of a given
891    TE sequence compared to a consensus.

892

893    **Figure 3. Endogenous Viral Elements (EVEs) are found throughout the**
894    **Aag2 genome and are observed in piClusters.** (A) Circular plot showing
895    piRNAs mapping to the Aag2 genome, with those derived from EVEs highlighted
896    in red. Inner tracks show EVEs found in the Aag2 genome, colored by virus
897    family from which they derived. (B) i) Density plot of TEs along Aag2 contig
898    000015F. Local enrichment for particular TE classes can be seen. ii) Zoom of

899    indicated region in i). Transposons are indicated by colored arrows, categorized
900    by TE class.  '+' and '-' indicate strandedness of the transposons. EVEs are
901    shown as black triangles, organized by strandedness. piRNAs mapping to the
902    region are indicated at the bottom, in either sense or anti-sense direction.  Clear
903    clusters of EVEs and TEs can be seen, with piRNAs specifically mapping to
904    these clusters. The further zoom (iii) better shows the shared directionality of
905    EVEs (colored triangles) and their surrounding transposon sequences (black,
906    semi-transparent triangles). (C) Bar plot showing counts of EVEs derived from
907    different viral families.(D) Bar plot showing a larger number of piRNA reads can
908    be successfully mapped to the Aag2 (PacBio) assembly compared to the LVP
909    (Sanger) assembly.
910
911    **Figure 4. EVEs are primarily associated with LTR transposable elements.**
912    (A) Histograms showing counts of non-overlapping TEs closest to EVEs binned
913    by distance, both upstream (negative x-axis values) and downstream (positive x-
914    axis values). Positive y-axis counts refer to TE/EVE 'pairs' with the same
915    strandedness , while negative y-axis counts are EVEs where the closest TE has
916    the opposite strandedness. The "+/-" value indicates the ratio of TE/EVE pairs
917    with the same strandedness to those with the opposite strandedness. Total
918    counts represented in each histogram: All classes (n=766); LTR only (n=475); No
919    LTRs (n=291); Ty3/gypsy only (n=274); Pao Bel only (n=180); Ty1/copia only
920    (n=21). (B) Heatmap showing categories of TEs nearest EVEs, categorized by
921    the viral family from which the EVEs were derived. Only TEs with the same
922    strandedness as its nearest EVE are shown. A "*" indicates significant
923    enrichment by one-sided binomial test against the background prevalence of a
924    given TE category in the genome (eg among all LTRs nearest Chuviridae-derived
925    EVEs, Pao Bel elements are specifically enriched compared to the genome-wide
926    counts of Pao Bel among all LTRs). Color indicates proportion of a given TE
927    category nearest EVEs derived from the indicated viral family. Grey indicates the
928    element was not found to be the closest TE to any EVEs derived from the
929    indicated viral family. Only TE elements which made up at least 10% of the
930    dataset for a given viral family are shown. "Pao Bel elements" refers to
931    Chuviridae, while "Ty3/gypsy elements" corresponds to Flaviviridae and
932    Rhabdoviridae. Total sample size of all TEs analyzed for each dataset: AllTEs-
933    Rhabdoviridae (n=223),Flaviviridae (n=118), Chuviridae (n=112); RNA-
934    basedTEs- Rhabdoviridae (n=160),Flaviviridae (n=116), Chuviridae (n=95);
935    LTRs- Rhabdoviridae (n=123),Flaviviridae (n=106), Chuviridae (n=84); Pao Bel-
936    Chuviridae (n=67); Ty3/gypsy- Rhabdoviridae (n=99), Flaviviridae (n=77).
937
938    **Figure 5. EVEs found in the genomes of *Aedes aegypti* and *Aedes***
939    ***albopictus* are derived from overlapping families of viruses, but do not**
940    **share the same sequences.** (A) Phylogenetic relationship between 61 members
941    of Flaviviridae. EVEs present in (i) *Ae. aegytpi* or (ii) *Ae. albopictus* which align to
942    the indicated virus are marked with a colored circle. Size corresponds to
943    abundance of EVEs derived from given species. (B) Coverage plots of EVEs
944    derived from the viral families (i) Flaviviridae, (ii) Rhabdoviridae, and (iii)

945     Chuviridae. Each bar represents a single EVE, while its length and position
946     denotes the region of the indicated ORF from which its sequence is derived.
947     Length is expressed as a percentage of the total ORF, in order to normalize for
948     varying ORF lengths among different members of a given viral family. In (i), the
949     genome of CFAV is presented for reference. In (ii) and (iii), a generic genome is
950     presented to better illustrate where EVEs are derived from within the genome as
951     a whole (and within each specific ORF). (C) Violin plot of Kimura divergence
952     scores for EVEs' nearest-neighbor TEs. Only EVEs whose nearest neighbor TE
953     is non-overlapping and shares the same strandedness are shown. Both
954     upstream and downstream nearest neighbors are represented. Counts in each
955     category are as follows: Pao Bel-All (n=113,526); Pao Bel-Chuviridae (n=36);
956     Pao Bel-Flaviviridae (n=12); Pao Bel-Rhabdoviridae (n=14); Ty3/gypsy -All
957     (n=174,353); Ty3/gypsy -Chuviridae (n=10); Ty3/gypsy -Flaviviridae (n=43);
958     Ty3/gypsy-Rhabdoviridae (n=53). Differences between indicated distributions
959     were determined by Kolmogorov-Smirnov test. '*': p-value <0.05; '**': p-value
960     <5e-4; '***': p-value <5e-8.
961
962

**LEGENDS FOR SUPPLEMENTARY FIGURES**

**Figure 2-Figure Supplement 1. Transposons are distributed throughout the entire Aag2 genome.** (A) Similar plot to Figure 2A, but showing all contigs of the Aag2 assembly. Circular plot of the Aag2 assembly, with every contig (black rectangles; ordered by size) and transposable element (circles; colored by TE class). Transposable elements are prevalent throughout the entire *Ae. aegypti* genome. Rectangles representing contigs are staggered to indicate relative contig size.

**Figure 4-Figure Supplement 1. TEs which overlap EVEs are also overrepresented by LTR elements.** (A) Histograms of TEs which overlap EVEs, broken down by the indicated categories. The left bin represents TEs whose start is upstream, and end overlaps the EVE. The right bin indicates TEs whose end is downstream, and start overlaps an EVE. The middle bin indicates TEs whose coordinates surround an EVE. Positive count values indicate TE and EVEs with shared directionality, while negative values represent TE and EVEs with opposite directionality. Some EVEs showed multiple overlapping TEs, all of which are represented on the charts. (B) Heatmap, as in Figure 3, showing EVE-overlapping TE 'preference' for Rhabodoviridae, Flaviviridae, and Chuviridae-derived EVEs.

**Figure 3-Figure Supplement 1. EVEs are typically found within unidirectional piRNA clusters.** The left panels correspond to a region of Contig 000933F encoding 4 tandem, unidirectional piRNA clusters (as identified by proTRAC), each containing EVEs. Each cluster expresses piRNAs primarily anti-sense to the TEs/EVEs which define them. Similarly, a single large piRNA cluster on Contig 000044F is shown in the right panels. The shared directionality between TEs and EVEs (Figure 3B) is evident. Again, piRNA expression is almost exclusively in the antisense direction with respect to the TEs/EVEs.

**Figure 2-Figure Supplement 2. Kimura divergence scores of LTRs only show expansion of Pao Bel and Ty3/gypsy elements.** Bar plot of kimura scores assigned to LTRs only, categorized by TE family and expressed as percent of total genome (as in Figure 2E). At very low (0-1) Kimura divergence scores, Ty3/gypsy and Pao Bel exhibit a marked increase in proportion of the genome.

**Figure 5-Figure Supplement 1. Kimura divergence of EVE proximal TEs is distinct from TEs genome wide.** Density plots of kimura distributions of all Pao Bel or Ty3/gypsy TEs and EVE-proximal TEs. EVE-proximal TEs are further categorized by the viral family of the EVE it is nearest.
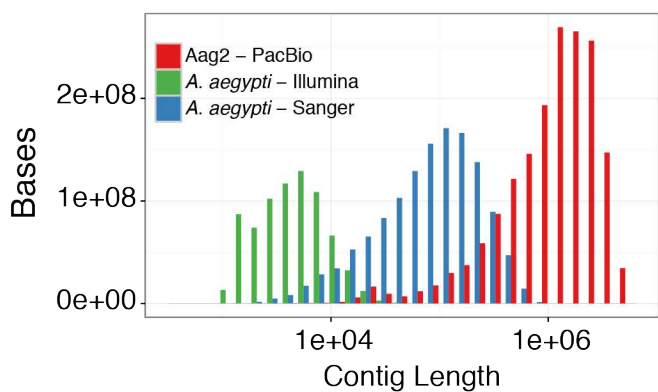
|                            | UCB           | LVP           | Aag2           |
| -------------------------- | ------------- | ------------- | -------------- |
| Sample                     | LVP strain    | LVP strain    | Aag2 cell line |
| Seq. Strategy              | Illumina      | Sanger        | PacBio         |
| Released                   | 5/2015        | 6/2006        | NA             |
| Coverage                   | 6.8x          | 7.6x          | ~50x           |
| Total sequence length      | 744,596,036   | 1,383,957,531 | 1,723,930,323  |
| Total assembly gap length  | 196,533,049   | 73,881,199    | 0              |
| Num. of contigs            | 961,292       | 36,204        | 3,752          |
| Contig N50                 | 989           | 82,618        | 1,420,116      |
| Contig L50                 | 151,087       | 4,346         | 368            |

| | Num. of elements | Length (Mbp) | Percent of genome |
|---|---|---|---|
| SINE | 28,301 | 4.4 | 0.25 |
| LINE | 558,382 | 259.9 | 15.07 |
| LTR | 495,204 | 163.9 | 9.51 |
| DNA | 1,184,522 | 309.0 | 17.93 |
| Other* | 725,958 | 233.7 | 13.55 |
| Total | 2,992,367 | 970.9 | 56.31 |

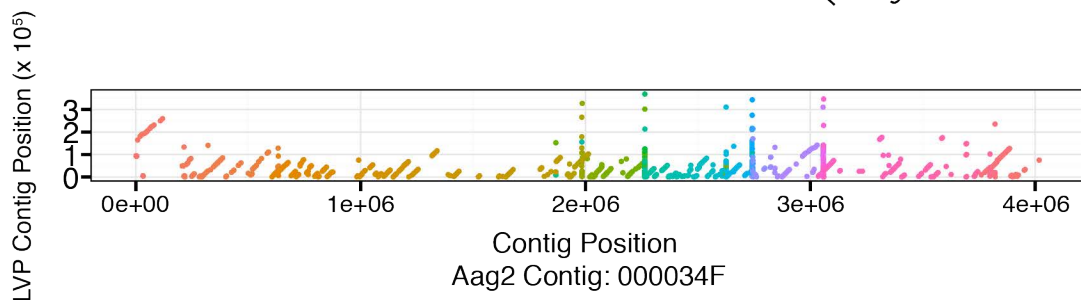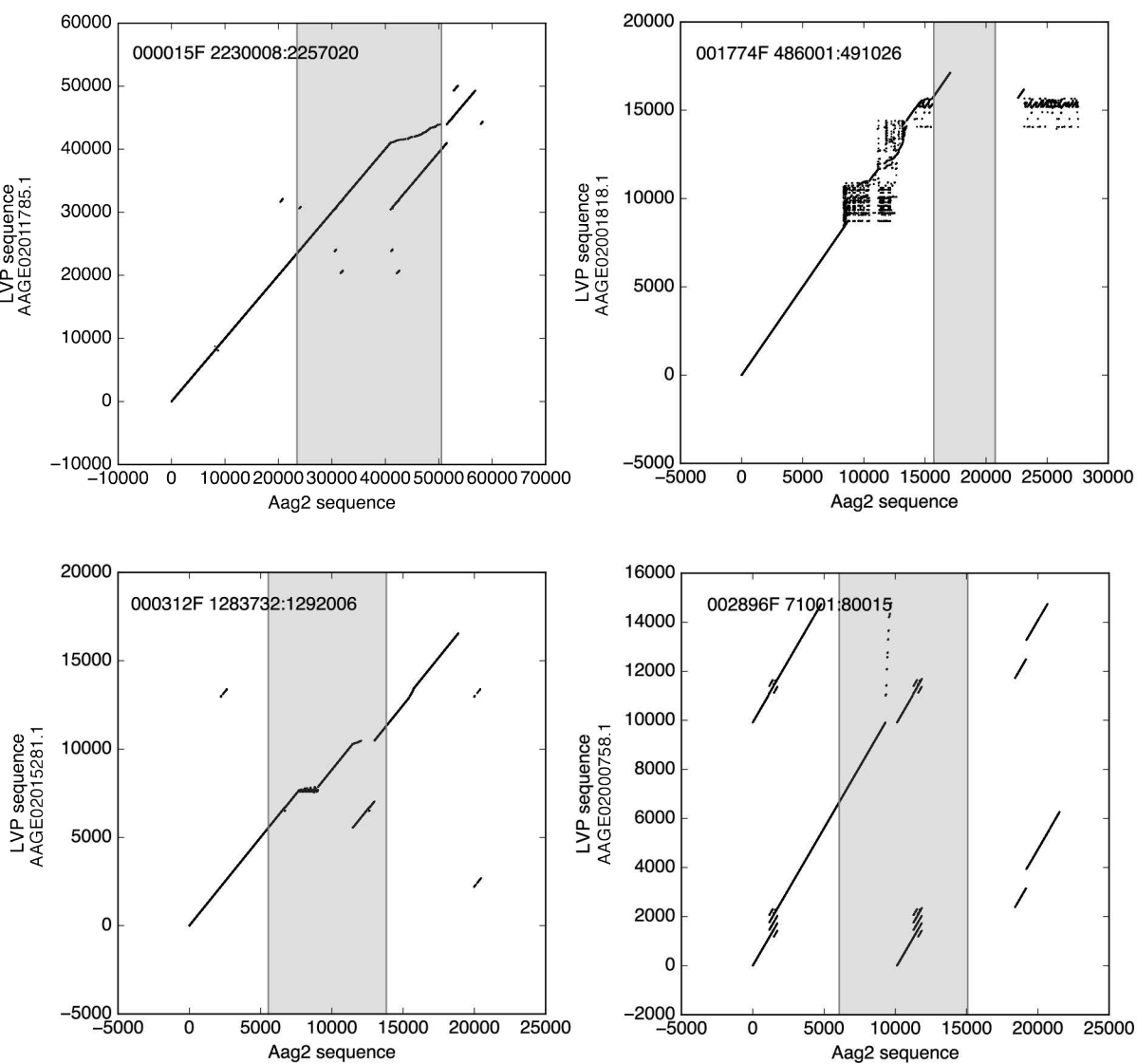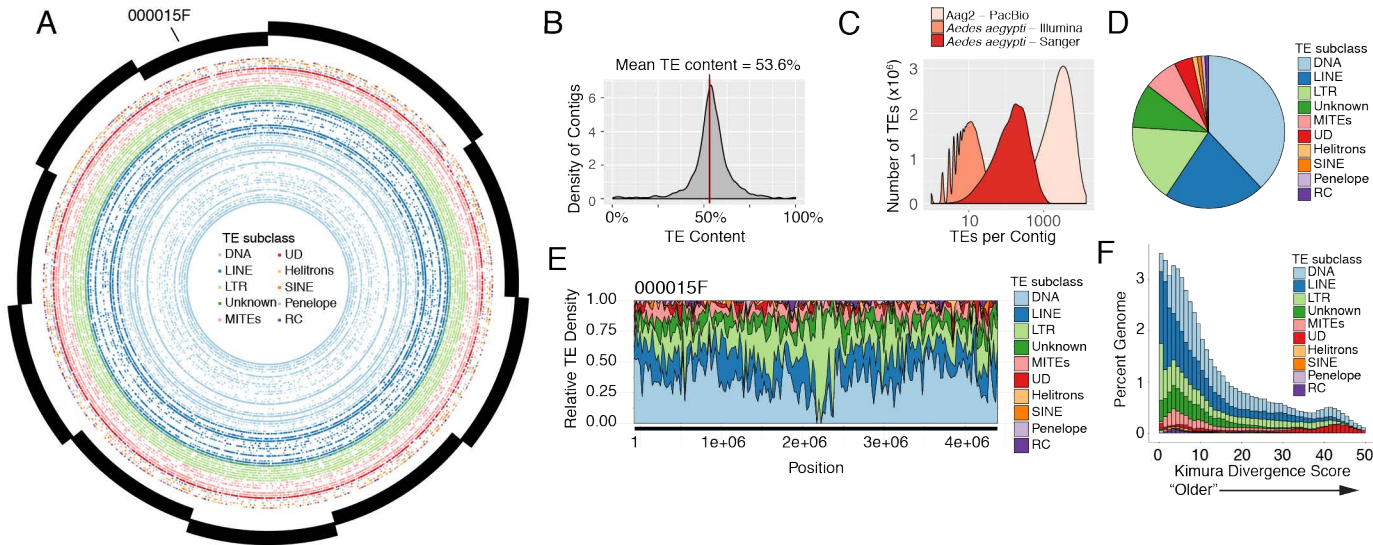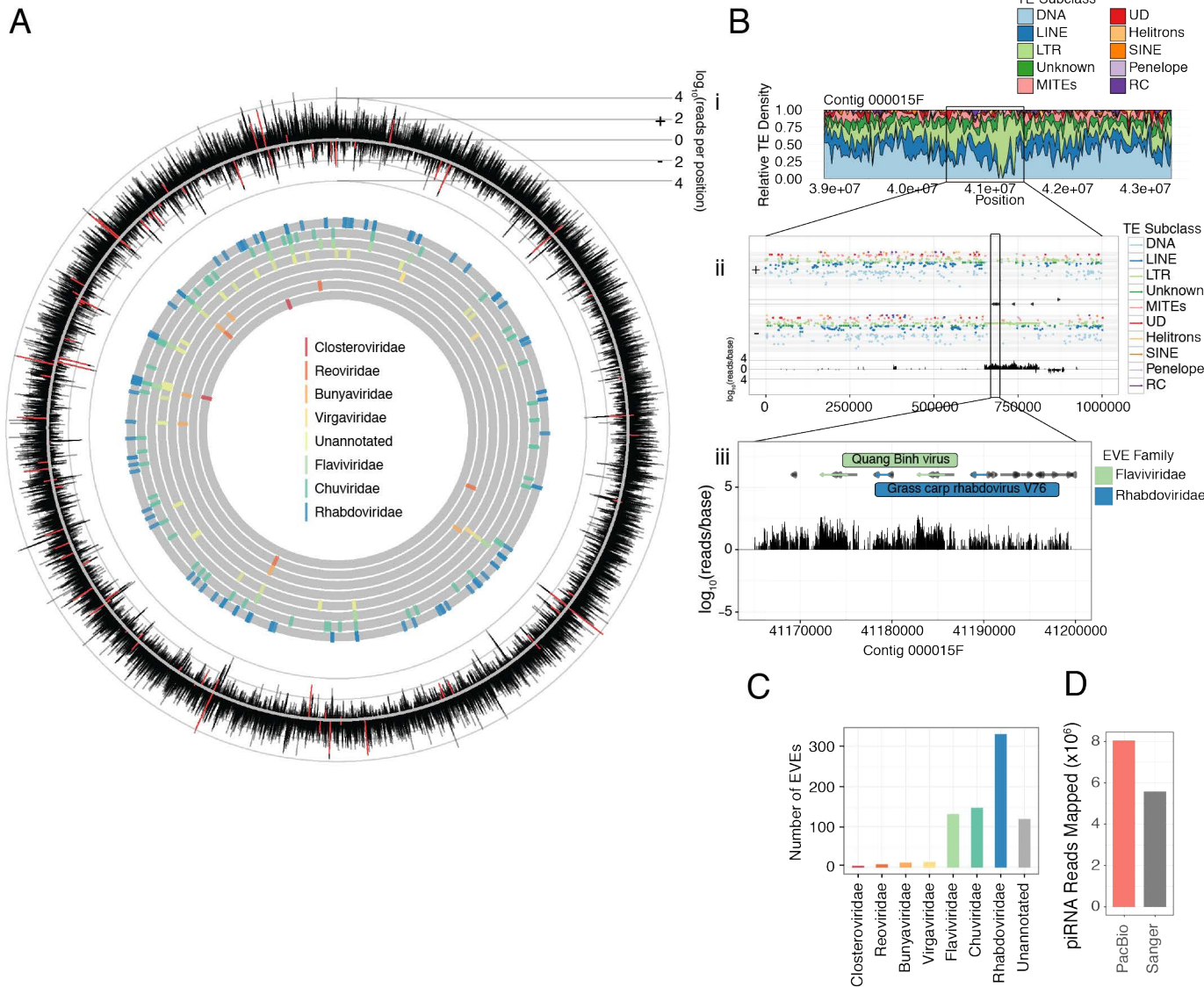*includes helitrons, MITEs, Penelope, RC, UD, and unknown elements

Fig. 1

Fig. 2



A  000015F

B  Mean TE content = 53.6%

C  Aag2 – PacBio / Aedes aegypti – Illumina / Aedes aegypti – Sanger

D  TE subclass: DNA, LINE, LTR, Unknown, MITEs, UD, Helitrons, SINE, Penelope, RC

E  000015F

F  TE subclass: DNA, LINE, LTR, Unknown, MITEs, UD, Helitrons, SINE, Penelope, RC

Fig. 3



A

B

i

ii

iii

C

D

TE Subclass
DNA
LINE
LTR
Unknown
MITEs
UD
Helitrons
SINE
Penelope
RC

Closteroviridae
Reoviridae
Bunyaviridae
Virgaviridae
Unannotated
Flaviviridae
Chuviridae
Rhabdoviridae

$\log_{10}$(reads per position)

Contig 000015F

Relative TE Density

Position

$\log_{10}$(reads/base)

Quang Binh virus

Grass carp rhabdovirus V76

EVE Family
Flaviviridae
Rhabdoviridae

$\log_{10}$(reads/base)

Contig 000015F

Number of EVEs

piRNA Reads Mapped ($\times 10^6$)

PacBio

Sanger

Fig. 4

Fig. 5

Fig. 2, Suppl 1

subclass
- DNA
- LINE
- LTR
- Unknown
- MITEs
- UD
- Helitrons
- SINE
- Penelope
- RC

Fig. 3, Suppl 1

# A

## Overlapping TE/EVE pairs



i — All Transposon Classes

ii — LTR Retrotransposons only

iii — Ty3/gypsy only

iv — Pao Bel only

v

Ty3/gypsy | Pao Bel — Fig. 5, Suppl 1

All TEs

EVE proximal TEs

virus family
- Chuviridae
- Flaviviridae
- Rhabdoviridae