

1 **High-Throughput Sequencing of Transposable**  
2 **Elements Insertions Reveals Genomic Evidence for**  
3 **Adaptive Evolution of the Invasive Asian Tiger**  
4 **Mosquito Towards Temperate Environment**

5 Clément Goubert<sup>1</sup>, Hélène Henri<sup>1</sup>, Guillaume Minard<sup>2,3</sup>, Claire Valiente Moro<sup>2</sup>, Patrick  
6 Mavingui<sup>2,4</sup>, Cristina Vieira<sup>1</sup>, and Matthieu Boulesteix<sup>1</sup>

7 <sup>1</sup>Université de Lyon, F-69622, Lyon, France; Université Claude Bernard Lyon 1, CNRS,  
8 Laboratoire de Biométrie et Biologie Evolutive, UMR5558, F-69100 Villeurbanne

9 <sup>2</sup>Université de Lyon, F-69622, Lyon, France; Université Lyon 1, Villeurbanne, France;  
10 CNRS, UMR 5557, Ecologie Microbienne, Villeurbanne, France; INRA, UMR 1418,  
11 Villeurbanne, France

12 <sup>3</sup>Metapopulation Research Center, Department of Biosciences, University of Helsinki,  
13 Helsinki, Finland

14 <sup>4</sup>Université de La Réunion, UMR PIMIT, INSERM 1187, CNRS 9192, IRD 249,  
15 Plateforme Technologique CYROI, Sainte-Clotilde, La Réunion.

16 **Corresponding author:** Matthieu Boulesteix, Laboratoire de Biométrie et Biologie  
17 Evolutive, UMR CNRS 5558, INRIA, VetAgro Sup, Université Claude Bernard Lyon 1,  
18 Villeurbanne, France, +334 72 43 29 16, [matthieu.boulesteix@univ-lyon1.fr](mailto:matthieu.boulesteix@univ-lyon1.fr)

## 19 **Abstract**

20 Invasive species represent unique opportunities to evaluate the role of local  
21 adaptation during colonization of new environments. Among these, the Asian tiger  
22 mosquito, *Aedes albopictus*, is a threatening vector of several human viral dis-  
23 eases, including dengue, chikungunya and the emerging Zika feverssometimes  
24 been considered as the reflect of a great "ecological plasticity". However, no  
25 study has been conducted to assess the role of adaptive evolution in the ecologi-  
26 cal success of *Ae. albopictus* at the molecular level. In the present study we per-  
27 formed a genomic scan to search for potential signatures of selection leading to  
28 local adaptation in a hundred of field collected mosquitoes from native popula-  
29 tions of Vietnam and temperate invasive populations of Europe. High throughput  
30 genotyping of transposable element insertions generated more than 120 000  
31 polymorphic loci, which in their great majority revealed a virtual absence of struc-  
32 ture between bio-geographic areas. Nevertheless, 92 outlier loci show a high  
33 level of differentiation between temperate and tropical populations. The majority  
34 of these loci segregates at high insertion frequencies among European popula-  
35 tions, indicating that this pattern could have been caused by recent events of  
36 adaptive evolution in temperate areas. Six outliers were located near putative dia-  
37 pause effector genes, suggesting fine tuning of this critical pathway during local  
38 adaptation.

39

40 *Keywords:* Invasive species, *Aedes albopictus*, local adaptation, genome scan,  
41 Transposable Elements, diapause

42

43

## 44 **Author Summary**

45 to evaluate the importance of local adaptation in the ecological success of inva-  
46 sive species, especially in the case of major disease vectors. In this context, we  
47 investigated whether adaptation has facilitated the invasion of the Asian tiger  
48 mosquito *Aedes albopictus* in temperate environments. This species, that already  
49 transmits Yellow fever and Chikungunya viruses, is also competent for the Zika  
50 virus, posing dramatic sanitary consequences given the current species distribu-  
51 tion. We genotyped the insertion polymorphism of mobile genetic elements, used  
52 as dense genetic markers, in more than a hundred field collected individuals from  
53 temperate and tropical populations. We identified several very divergent markers  
54 between these populations that points to recent targets of natural selection. In  
55 depth analyses of our results suggests that the diapause pathway could have  
56 been this way tuned during the invasion of temperate environments.

## 57 **Introduction**

58 Biological invasions represent unique opportunities to study fast evolutionary changes  
59 such as adaptive evolution. Indeed, settlement in a novel area represents a biological  
60 challenge that invasive species have successfully overcome. The underlying processes  
61 could be studied at the molecular level, particularly to gather empirical knowledge about  
62 the genetics of invasions, a field of study that has produced extensive theoretical  
63 predictions, but for which there is still little evidence in nature (1). Some of the main  
64 concerns are to disentangle the effects of neutral processes during colonization, such as  
65 founder events or allele surfing at the migration front, from adaptive evolution (i.e. local  
66 adaptation, (1–3)).

67 Adaptation can arise either through the appearance and spread of a new beneficial  
68 mutation, the spread of a favorable allele from standing genetic variation, or from  
69 hybridization in the introduction area (1, 4, 5). Detection of the footprint of natural  
70 selection is however dependent on the availability of informative genetic markers, which  
71 should provide a substantial coverage of the genome to allow selection scans and be  
72 easily and confidently scored across many individuals. Unfortunately, invasive  
73 organisms are rarely model species, making the development of a reliable and efficient  
74 marker challenging.

75 The Asian tiger mosquito, *Aedes (Stegomyia) albopictus* (Diptera:Culicidae) is currently  
76 one of the most threatening invasive species (Invasive Species Specialist Group);  
77 originating from South-Eastern Asia, this species is one of the primary vectors of

78 Dengue and Chikungunya viruses, and is also involved in the transmission of other  
79 threatening arboviruses (6), in particular the newly emerging Zika virus (7–9). *Ae.*  
80 *albopictus* has now settled in every continent except Antarctica, and is found both under  
81 tropical and temperate climates (10). While this species is supposed to originate(11), the  
82 native area of *Ae. albopictus* encompasses contrasted environments including  
83 temperate regions of Japan and China, offering a large potential of fit towards newly  
84 colonized environments. For example, the induction of photoperiodic diapause in  
85 temperate areas, that has a genetic basis in *Ae. albopictus* (12, 13), is decisive to  
86 ensure invasive success in Europe or Northern America. Indeed it allows the sensible  
87 populations to survive through winter at the larval stage into the eggs. Such a trait  
88 appears governed by a “genetic toolkit” involving numerous genes and metabolic  
89 networks, for which however the genetic polymorphism between diapausing and non-  
90 diapausing strains remains to be elucidated(14). In addition, the colonization of new  
91 areas that look similar at first glance can still involve *de novo* adaptation: indeed, even  
92 environment sharing climatic variables are not necessarily similar regarding edaphic and  
93 biotic interactions(1). Hence, this suggests that whatever are the native and settled might  
94 be possible to find evidence of adaptive evolution in invasive populations of *Ae.*  
95 *albopictus*.

96 To better understand the invasive success of this species, we genotyped 140 field  
97 individuals, collected from three Vietnamese (native tropical area) and five European  
98 (invasive temperate area) populations, aiming to identify genomic regions involved in  
99 local adaptation. To do so, we developed new genetic markers, based on high

100 throughput genotyping of the insertion of Transposable Elements (TEs), which are highly  
101 prevalent and polymorphic in the *Ae. albopictus* genome.

102 To distinguish between neutral demographic effects and adaptive evolution, we first  
103 performed population genetic analyses to reveal the global genetic structure of the  
104 studied populations. We then performed a genomic scan for selection and identified 92  
105 candidate loci under directional selection, among which several can be located in the  
106 neighborhood of diapause related genes.

107

## 108 **Results**

### 109 **Rationale for using TE markers in identifying regions under natural selection**

110 Amplification of TE insertions is particularly efficient to obtain a large number of genetic  
111 markers throughout one genome (15), especially if few genomic resources are available  
112 (16), which was case until recently for the Asian tiger mosquito. We hypothesized that  
113 some TE insertion sites could be located at the neighborhood of targets of natural  
114 selection and thus could reach high level of differentiation between native and invasive  
115 populations if selective sweeps occurred during local adaptation. In addition, some TEs  
116 could also insert near or inside coding regions and thus could be directly involved in  
117 environmental adaptation (17), eventually contributing to the success of invasive species  
118 (18).

119 Such genetic elements represent at least one third of the genome of *Ae. albopictus* and

120 include recently active families that could reach thousands of copies in one genome  
121 (19). Using such markers, represent a seducing alternative to other methods of diversity  
122 reduction, such as RAD-sequencing (20), that could be less efficient in species with high  
123 TE load (21). In mosquitoes, TEs have been shown to be powerful markers both for  
124 population structure analysis (22–25) and genome scans (15).

125

126 **High throughput TE insertion genotyping.** A total of 140 individuals were collected in  
127 Europe (invasive temperate populations) and Vietnam (native tropical), and screened for  
128 their insertion polymorphism of five highly repeated families of TEs using paired-end  
129 Illumina sequencing (see Material and Methods). Briefly, individual insertions of five TE  
130 families (IL1, L2B, RTE4, RTE5 and Lian1) were genotyped using Transposon Display ()  
131 (26), a TE insertion specific PCR method, combined with Illumina sequencing of all TD  
132 amplification products. TEs used for this study are non-LTR (LINE) retrotransposons that  
133 usually do not show a specific insertion site preference (31), making them likely to be  
134 well dispersed in the genome. Sequencing produced a total of 102,319,300 paired-end  
135 reads (2x101bp). After quality and specificity filtering 24,332,715 reads were suitable for  
136 analyses. Because of read coverage variation between individuals, we applied a read  
137 sampling procedure before the recovery of individual insertion loci by clustering; to  
138 ensure the consistency of this procedure, sampling and subsequent analysis were  
139 performed independently three times. On average, a total number of 128,491  
140 polymorphic insertion loci were available for each of the three sampling replicates. The  
141 mean number of loci per individual and per TE family ranged from  $1025 \pm 290$  s.d. (IL1

142 family, mean and s.d. averaged over the three replicates) to  $3266 \pm 766$  s.d. (RTE5  
143 family). Details are given in S1 Table. While our read sampling procedure could have  
144 artificially lowered the mean insertion frequency of the loci, this effect should be small  
145 because in our final datasets, the TE insertion frequencies (i. e. the number of  
146 individuals that share an insertion) are not correlated with the mean number of read per  
147 individual at the considered locus (S1 Figure).

148

149 **Population structure.** Principal Coordinate Analyses (PcoAs) were performed  
150 independently for each of the five TEs (Figure 1). Among the three main Principal  
151 Coordinates (PCs), individuals tend to be grouped according to their respective  
152 populations with little overlap between groups. However, the three main PCs represent  
153 only a small fraction of total genetic variation (< 10%), suggesting a weak genetic  
154 structuring between the populations. Overall, individuals from Vietnamese populations  
155 (HCM, TA, VT) tend to be grouped together in a single cluster, at the exception of 13 to  
156 14 individuals from HCM when using L2B and RTE5 TE families (S2 Figure), along with  
157 six individuals of VT with the RTE4 TE family (Figure 1) that can not be clearly  
158 distinguished from European samples. BCN individuals (Spain) represent the most  
159 homogeneous group, well differentiated from Vietnamese and French individuals (SP,  
160 CGN, NCE and PLV).

161 In agreement with PCoAs, Analyses of Molecular Variance (AMOVAs (27)) attributed  
162 very few genetic variances among groups (Vietnam-Europe) and between populations  
163 within groups (Table 1). In the studied populations, most of the genetic variance was



164 distributed among individuals within groups.

165 Measures of genetic differentiation among pairs of populations were consistent with

166 PCoAs and AMOVAs (S1 File): the BCN population shows the highest  $F_{ST}$  values with

167 the other populations for each of the five TEs ( $0.051 < F_{ST} < 0.148$ ), while Vietnamese

168 populations were the most closely related ( $0.011 < F_{ST} < 0.032$ ). While VT is located

169 100 km away from TA and HCM (both sampled in the same city, Hô Chi Minh, Vietnam)

170 the  $F_{ST}$  values are very similar between the three Vietnamese populations, suggesting

171 no influence of geography at this scale. CGN and NCE, sampled in the same urban area

172 (Nice agglomeration), are also little or not significantly differentiated depending on the

173 TE family. The previously identified intermediate pattern of HCM with some European

174 populations at L2B and RTE5 loci (PCoAs analyses) is also found at the  $F_{ST}$  level,

175 especially regarding the low differentiation with the PLV population for these markers

176 ( $0.011 < F_{ST} < 0.020$ ).

177

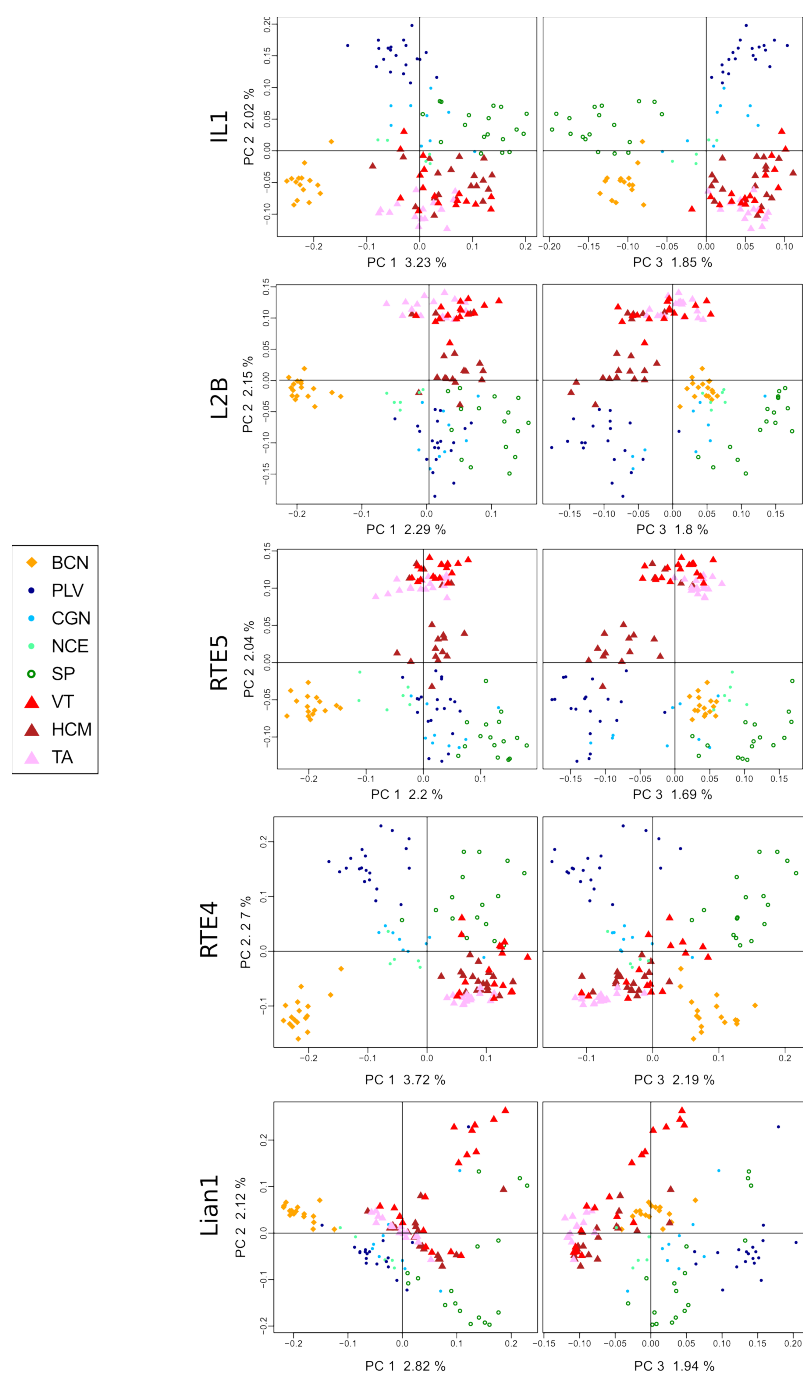
178

179

180

181

182



183 **Fig 1. Principal Coordinate Analysis (PCoAs).** Projection of individuals over the three first principal coordinates  
 184 (PC) of PCoAs for each of the 5 TE families and for the first replicate (M1, see Material and Methods). Proportion of  
 185 inertia represented by each axes is noted in %. circles: European populations; triangles: Vietnamese populations.  
 186 Results for other sampling replicate can be found in S2 Figure.

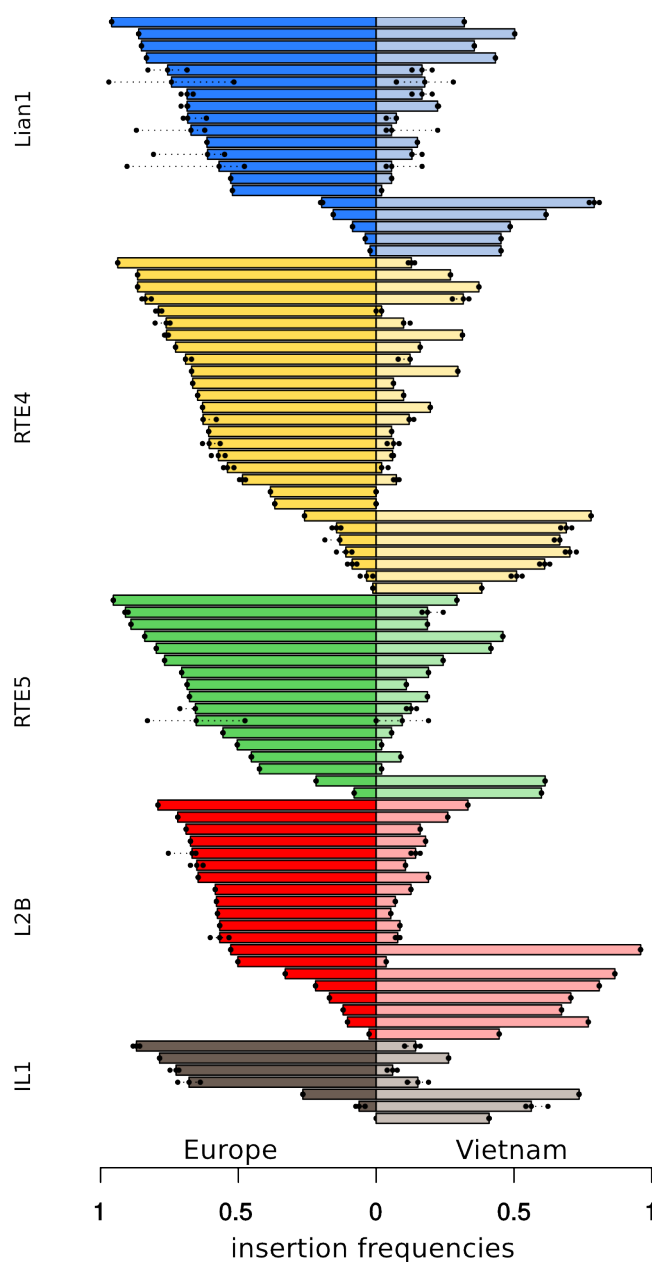
187 **Table 1. Analyses of Molecular Variance (AMOVAs).** Results for the three replicates (M1,M2,M3) of read sampling  
188 for the five TE families (IL1, L2B, RTE5, RTE4, Lian1). Values are given in percentage of the total genetic variance

189

	IL1	L2B	RTE5	RTE4	Lian1
Among groups	[0.59-0.70]	[1.22-1.29]	[1.08-1.10]	[1.97-2.04]	[0.67-0.74]
Among populations within groups	[5.15-5.37]	[3.58-3.63]	[3.36-3.40]	[6.67-6.78]	[4.47-4.56]
Within populations	[94.04-94.16]	[95.08-95.18]	[95.51-95.55]	[91.18-91.30]	[94.77-94.81]

intervals reports min and max values among the 3 sampling replicates

191 **Genomic scan.** Research of outlier loci for both selection signature using Bayescan  
192 (non-hierarchical island model), and for significant  $F_{CT}$  (between Europe-Vietnam group  
193 differentiation) identified 92 candidate insertion loci (Figure 2). Most of these insertions  
194 are found in both areas (no private allele), except for RTE4\_6 and RTE4\_7 that were not  
195 found in Vietnam. In addition, a majority of outliers corresponds to high frequency  
196 insertions in Europe, while the same trend is not observed at 92 randomly chosen loci  
197 among those having the same minimum insertion frequency ( $\geq 20$  individuals/locus)  
198 between Europe and Vietnam (Figure 3). PCR amplification of the outlier loci were  
199 carried out on a representative panel of 47 individuals to validate the insertion pattern  
200 detected by TD (see Material and Method). For loci where the amplification was  
201 successful, the insertion pattern observed by PCR always confirmed the results from TD  
202 (S3 Figure).

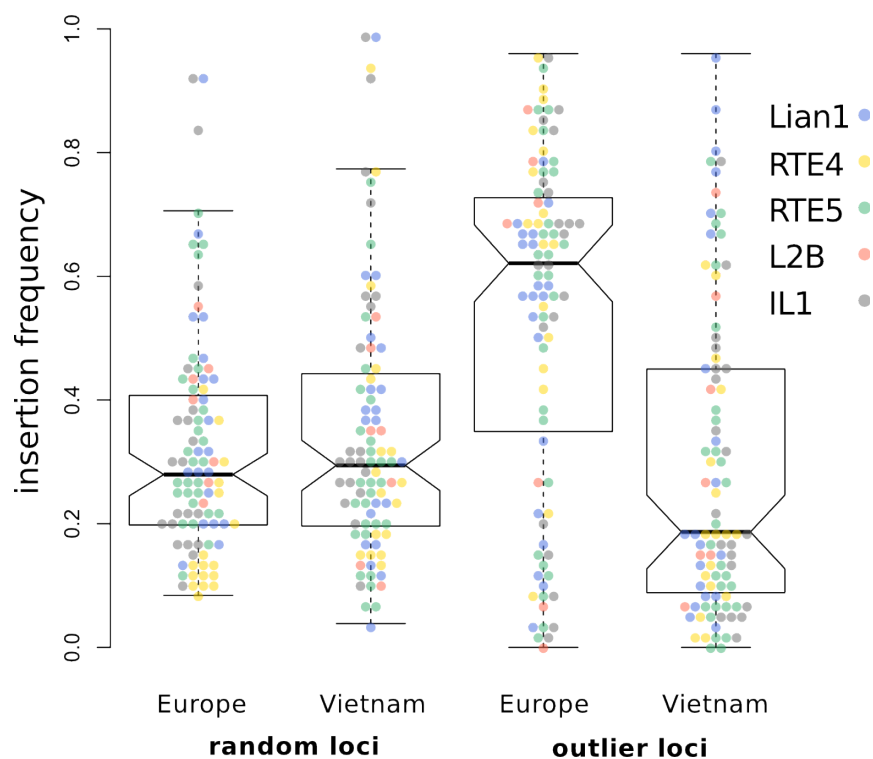


203 **Fig 2. Insertion frequencies in Europe and Vietnam for the 92 outlier loci.** Bars represent the median value from  
204 the three reads sampling replicates and dots the values from the other replicates (if outlier found in replicate). Colors  
205 correspond to each of the 5 TE families.

206 From 92 outlier loci, 21 could be attributed to a unique position on the *Ae. albopictus*  
207 genome (28). Annotation and distance to surrounding genes are reported on S2 File. We

208 found that six outliers (S2 File, sheet 3, highlighted) loci are located on contigs that  
209 harbor genes previously identified in *Ae. albopictus* as being differentially expressed  
210 between diapause-induced and non diapause-induced samples , and for which  
211 orthologs in *Drosophila melanogaster* are known to be part of well-identified functional  
212 networks (14). All these six loci are found to be outliers because of their high insertion  
213 frequencies in Europe compared with Vietnam.

214



215 **Fig 3. Comparison of outliers frequencies with randomly selected loci.** Insertion frequencies of 92 randomly  
216 chosen loci among those having the same minimum insertion frequency ( $\geq 20$  individuals) as outliers compared to the  
217 92 outlier loci. Random loci were taken from the first replicate (M1) and values for outliers are median values obtained  
218 among the three replicates. Non-overlapping notches indicate a significant difference between the true medians (thick  
219 dark horizontal bars).

## 220 Discussion

221 The goal of our study was to identify genomic regions involved in adaptive evolution of  
222 *Ae. albopictus* thanks to the development of new genetic markers. Through high-  
223 throughput genotyping of five TE families insertion polymorphisms, we identified up to  
224 128,617 polymorphic loci among a hundred of individuals from eight sampling sites. The  
225 estimated genome size of *Ae. albopictus* exceeds one billion base-pairs (15, 29, 30).  
226 Accordingly, the amount of markers scored in this study offers a comfortable genomic  
227 density of one marker every 10 kb.

228 We provide here a new and cost efficient method to quickly generate a large amount of  
229 polymorphic markers without extensive knowledge about one species genome.  
230 Specifically, this strategy could be extremely valuable for species with a large genome  
231 size, where TE density could severely compromise the development of more classical  
232 approaches, such as the very popular RAD-sequencing (16).

233 The genetic structure of the studied populations showed strong consistency between  
234 sampling replicates of individuals' reads, demonstrating the robustness of the method in  
235 spite of an initial substantial coverage variation among individuals. Population genetics  
236 analyses revealed a very low level of genetic structuring between European and  
237 Vietnamese populations. Among the studied populations, AMOVAs showed that most of  
238 the genetic variation is distributed between individuals within populations (> 90%), and  
239 as suggested by pairwise  $F_{ST}$  and PCoAs, only a small part (< 10%) of the genetic  
240 variance is due to differentiation between populations. The genetic differentiation we

241 measured is indeed as high among European populations as it is between populations  
242 from Europe and Vietnam.

243 This singular population structure is in agreement with previous results gathered in *Ae.*  
244 *albopictus* using different collections of allozymes, mtDNA or microsatellites markers  
245 (31–35). Moreover, a recent analysis performed with a set of 11 microsatellites on  
246 individuals from the same populations (at the exception of BCN) showed a similar  
247 distribution of genetic variation among hierarchical levels (36). These results  
248 demonstrate the reliability of our markers and confirm that a non-hierarchical island  
249 model can likely fit the global genetic structure. The observed genetic diversity is  
250 compatible with a scenario of multiple and independent introductions, as already  
251 suggested for *Ae. albopictus* (37–40). However, as previously suggested, this pattern  
252 could also be the result of founder events that may occur during colonization and/or a  
253 restriction of gene flow between populations consecutive to their introduction. Answering  
254 such a question would require an extended sampling all over the native area.

255 Outlier analysis revealed 92 loci with high posterior probabilities of being under positive  
256 selection between European and Vietnamese populations. When possible, the PCR  
257 amplification of the outlier loci using a set of representative individuals always confirmed  
258 a shift of insertion frequencies toward either the European or the Vietnamese sampling  
259 sites. This suggests that in spite of a reduced coverage, introduced by sampling in the  
260 dataset, the scored insertion polymorphisms are reliable. In addition, our method of  
261 analysis is likely to be conservative: the Bayesian outliers were selected for their  
262 consistency with a significant  $F_{CT}$  between European temperate and Vietnamese tropical

263 populations, which avoid retaining outliers that we were not looking for, for example  
264 those due to a population specific event.

265 We were able to assign a unique position for 21 of the outlier loci on the *Ae. albopictus*  
266 genome. As expected by the *Ae. albopictus* genomic composition (15, 28, 41), an  
267 important part of the other outlier loci were located in repeated regions (44,6% of total),  
268 despite our efforts to remove *a priori* loci occurring in known transposable elements.

269 Since the *Ae albopictus* genome publication is very recent, no gene set or other genome  
270 annotation are currently available. We thus took advantage of the *Ae. albopictus*  
271 transcriptome data to annotate regions surrounding the detected outliers.(14). We found  
272 six outliers located on contigs which also harbors genes that are differentially expressed  
273 between individuals induced for diapause and controls (14), two of them, being located  
274 either in an intron (RTE4\_7442) or within 3kb (RTE4\_17015) of these candidate genes.

275 It is worth mentioning that these two genes belong to the the same functional group  
276 (GO:0005576 extracellular region). Diapause is a critical developmental stage found  
277 only in temperate populations of the Asian tiger mosquito. Interestingly, this functional  
278 pathway has been shown to benefit from fast adjustments thanks to local adaptation.

279 For instance, Urbanski *et al.* (42) showed that invasive American populations originating  
280 from Japan have rapidly evolved a new adaptive clinal response to diapause induction,  
281 independent from that observed in the native area. Thus, adaption in the temperate  
282 regions could have led to several selective sweeps on gene or regulatory sequences  
283 involved in this critical pathway, allowing the settlement of the mosquito in new  
284 temperate areas.



285 Interestingly, and as it is the case for these six outliers, we found significantly more  
286 outlier loci with a high frequency in Europe and low frequency in Vietnam than the  
287 opposite pattern. This was unexpected regarding our initial assumptions: a favored allele  
288 selected in one or another environment has *a priori* no reason to be more often  
289 associated with the presence or the absence of a TE insertion at linked sites. However,  
290 we found that the majority of the sequenced TE insertions segregates at low  
291 frequencies (around 10% of all individuals). When considering the linked region of one  
292 polymorphic TE insertion, if a favorable mutation appears in an individual where the  
293 insertion is absent, the increase of frequency of this “absence” haplotype will thus, most  
294 of the time, have a modest effect on the genetic differentiation at this marker, since it is  
295 already segregating at high frequency. By contrast, if a favorable mutation appears in a  
296 TE “presence” haplotype, the increase in frequency of the linked TE insertion would lead  
297 to high  $F_{ST}$  ( $F_{CT}$ ) values. In absence of an alternative explanation, our outlier loci could  
298 thus indicate in which subset of populations the adaptive mutation occurred, and in the  
299 present case, this would have happened more frequently in the temperate populations.

300 Two scenarios, not mutually exclusive, could be invoked in the light of our data. A simple  
301 case would be a direct adaptive evolution in European invasive population that  
302 originated from tropical regions of the native area. A second hypothesis, could be that  
303 invasive temperate populations came from Northernmost territories of the native area  
304 such as northern China or Japan where *Ae. albopictus* populations are already cold-  
305 adapted. It would be thus interesting to know whether the observed signature of  
306 selection results from more “ancient” adaptations in the native area, or if it originates

307 from more recent fine tuning of cold-related traits in the invasive areas. A recent study  
308 (33) suggested, using new variable *COI* mtDNA sequences and historical species range  
309 modeling, that Northern areas of the native range of *Ae. albopictus* would be the latest  
310 to have been colonized after a range expansion from Southern refugia following the last  
311 glacial around 21,000 years ago (34). The authors suggested that *Ae. albopictus* may  
312 have followed the human populations during their expansion from South to North in this  
313 area, that began approximately 15,000 years ago. Thus wherever the origin of the  
314 invasive individuals sampled in Europe, it is likely that they are representatives of  
315 populations that had recently undergone a shift of selective pressure from tropical to  
316 temperate climatic conditions. This could explain why so many outliers are associated  
317 with high insertion frequency in Europe, and that candidate genes in the diapause  
318 pathway are found in the neighborhood of some of these outliers. An easy way to  
319 distinguish between these possibilities would be to search if the same outlier insertions  
320 are present in several temperate populations from the native area.

321 It is important to note that the results presented here only are restricted to a subset of  
322 the Asian tiger mosquito populations located in temperate and tropical environments. It  
323 is thus probable that some of the outliers detected could be specific to this particular  
324 comparison and do not reflect the global pattern of differentiation between tropical and  
325 temperate populations. Research of the same outliers between other tropical and  
326 temperate populations from the native and non-native areas would be extremely  
327 valuable to extrapolate our results at a larger scale. Should the same outlier insertions  
328 be found at high frequencies in temperate locations – such as in USA, Japan or China

329 —, extended investigations about the origin of invasive populations would help clarify if  
330 those similarities are due to an ancestral sweep or parallel sweeps that occurred  
331 independently in several populations. This study already provides for some candidate  
332 loci a set of functional primers that could be directly used to answer this question in any  
333 DNA sample of *Ae. albopictus*.

334 We report here the first leads supporting adaptive evolution at the molecular level in the  
335 Asian tiger mosquito. Progress in the annotation of published genomes, and the looming  
336 availability of supplementary genomic resources will allow to gain the most from these  
337 results. We hope that this work will contribute to unravel the implication of adaptive  
338 processes during the invasion of disease vectors.

## 339 **Material and Methods**

### 340 **Biological samples**

341 A total number of 140 flying adult females *Ae. albopictus* were collected in the field at  
342 eight sampling sites in Europe and Vietnam during the summers of 2012 and 2013 (S2  
343 Table). Individuals were either sampled using a single trap or using aspirators through  
344 the sampling site within a 50 meters radius. When traps were used, live mosquitoes  
345 were collected after a maximum of two days.

### 346 **High throughput Transposon Display (TD) genotyping.**

347 Insertion polymorphism of five transposable elements families: I Loner Ele1 (IL1), Loa  
348 Ele2B (L2B), RTE4, RTE5 and Lian 1 identified by Goubert *et al.* (15) in *Ae. albopictus*

349 were characterized. These TE families were chosen according to their high estimate of  
350 copy number (from 513 to 4203 copies), high identity between copies, and a “copy and  
351 paste” mode of transposition (all these TEs are non-LTR Class I retrotransposons). The  
352 protocol was developed combining methods from previous studies (26, 43–45) with high  
353 throughput Illumina sequencing of TD products (S4 Figure).

354 **DNA Extraction and TD adapted ligation.** Total DNA was extracted from whole adult  
355 bodies following the Phenol-Chloroform protocol described by Minard *et al.* (36).  
356 Individual extracted DNA ( $\approx 75\text{ng}$ ) was then used for enzymatic digestion in a total  
357 volume of  $20\ \mu\text{L}$ , with HindIII enzyme ( $10\text{U}/\mu\text{L}$ ) and buffer R (Thermo Scientific) for 3  
358 hours at  $37\text{C}$ . The enzyme was inactivated at  $80\text{C}$  for 20 minutes. TD adapters were set  
359 up hybridizing Hindlink with MSEB oligonucleotides ( $100\ \mu\text{M}$ , see S3 Table) in  $20\text{X}$  SSC  
360 and  $1\text{M}$  Tris in a total volume of  $333\ \mu\text{L}$  after 5mn of initial denaturation at  $92\text{C}$  and 1h at  
361 room temperature for hybridization of the two parts. Once ready, TD adapters were then  
362 ligated to  $20\ \mu\text{L}$  of the digested DNA mixing  $2\ \mu\text{L}$  of TD adapter with  $10\text{U}$  T4 ligase and  
363  $5\text{X}$  buffer (Fermentas) in a final volume of  $50\ \mu\text{L}$  for 3 hours at  $23\text{C}$ .

365 **Library construction.** For each individual, and for each of the five TE families, TE  
366 insertions were amplified by PCR (PCR 1) in a Biorad Thermal Cycler (either C1000 or  
367 S1000), in a final volume of  $25\ \mu\text{L}$ . Mixture contained  $2\ \mu\text{L}$  of digested-ligated DNA with  
368  $1\ \mu\text{L}$  dNTPs ( $10\text{mM}$ ),  $0.5\ \mu\text{L}$  TD-adapter specific primer (LNP,  $10\ \mu\text{M}$ , see S3 Table) and  
369  $0.5\ \mu\text{L}$  of TE specific primer ( $10\ \mu\text{M}$ ),  $1\text{U}$  AccuTaq polymerase ( $5\text{U}/\mu\text{L}$ ) with  $10\text{X}$  buffer  
370 and Dimethyl-Sulfoxyde (Sigma). Amplification was performed as follows: denaturation

371 at 98C for 30 seconds then 30 cycles including 94C for 15 seconds, hybridization at 60C  
372 for 20 seconds and elongation at 68C for 1 minute; final elongation was performed for 5  
373 minutes at 68C. For L2B and RTE5 TEs, a nested PCR was performed in order to  
374 increase specificity in the same PCR conditions using internal forward TE primers and  
375 LNP (S3 Table). PCR 1 primers include a shared tag sequence that was used for  
376 hybridization of the individual indexes by PCR 2.

377 For each TE, three independent PCR 1 were performed from the same digestion  
378 product. PCR 1 products (3 PCR \* 5 TE per individual) were then purified using volume-  
379 to-volume Agencourt AMPure XP beads (20 $\mu$ L PCR 1 + 20 $\mu$ L beads) and eluted in 30 $\mu$ L  
380 Resuspension buffer. After nanodrop quantification, equimolar pools containing the 3\*5  
381 PCR products per individual were made using Tecan EVO200 robot. Individual pools  
382 were then size selected for fragment ranging from 300 to 600 bp using Agencourt  
383 AMPure XP beads as follow: first magnetic beads were diluted in H<sub>2</sub>O with a ratio of 1:  
384 0.68 then add to 0.625 X PCR products in order to exclude long fragments. A second  
385 purification was performed using a non-diluted bead: DNA ratio of 1:8.3 to exclude small  
386 fragments.

387 Samples multiplexing was performed using home made 6 bp index (included in SRA  
388 individual name), which were added to the R primer (S3 Table) during a second PCR  
389 (PCR 2) with 12 cycles in ABI 2720 Thermal Cycler. Mixture contained 15ng PCR  
390 products, 1 $\mu$ l of dNTPs (10mM), 0.5 $\mu$ l MTP Taq DNA Polymerase (5U/ $\mu$ l, Sigma), 5 $\mu$ l  
391 10X MTP Taq Buffer and 1.25 $\mu$ l of each tagged-primer (20 $\mu$ m) in a final volume of 50 $\mu$ l.  
392 Amplification was performed as follows: denaturation at 94C for 60 seconds then 12

393 cycles including denaturation at 94C for 60 seconds, hybridization at 65C for 60 seconds  
394 and elongation at 72C for 60 seconds; final elongation was performed for 10 minutes at  
395 72C. PCR 2 products were purified using Agencout AMPure XP beads: DNA ratio of  
396 1:1.25 to obtain libraries. Finally, TD products were paired-end sequenced on an  
397 Illumina Hiseq 2000 (1 lane) at the GeT-PlaGe core facility (Genome and Transcriptome,  
398 Toulouse) using TruSeq PE Cluster Kit v3 (2x100 bp) and TruSeq SBS Kit v3.

### 399 **Bioinformatic treatment of TD sequencing.**

400 The different steps of the informatics treatment from the raw sequencing dataset to  
401 population binary matrices for presence/absence of TE insertions per individual are  
402 described in S5 Figure. A total number of 102,319,300 paired-end 101bp Illumina reads  
403 were produced by sequencing PCR products. First, the paired-end reads of each  
404 individual were quality checked and trimmed using UrQt v. 1.0.17 (46) with standard  
405 parameters and a *t* quality threshold of 10. Reads pairs were then checked and trimmed  
406 for Illumina adapter contamination using cutadapt (47). Specific amplification of TE  
407 insertions was controlled by checking for the expected 3' TE sequence on the R1 read  
408 using Blat (48) with an identity threshold of 0.90. Only reads with an alignment-  
409 length/read-length ratio  $\geq 0.90$  were then retained. R2 reads for which the R1 mate  
410 passed this filter were then selected for the insertion loci construction, after the removal  
411 of the TD adapter on the 5' start using cutadapt and the removal of reads under 30 bp.  
412 Selected reads were separated in each individual according to the TE families for loci  
413 construction.

414 In order to correct the inter-individual coverage variations, we performed a sampling of

415 the cleaned reads. First, for each TE family, distribution of the number of read per  
416 individual was drawn, and individuals with less reads than the first decile of this  
417 distribution were removed; then cleaned reads of the remaining individuals were  
418 randomly sampled at the value of the first decile of coverage (this value varies among  
419 TE families). For each TE, the sampled reads of each retained individual were clustered  
420 together using the CD-HIT-EST program (49) to recover insertion loci. During this all-to-  
421 all reads comparison, the alignments must had a minimum of 90 percent identity, and  
422 the shortest sequence should be 95% length of the longest, global identity was used and  
423 each read was assigned at its best cluster (instead of the first that meet the threshold).  
424 In a second step, the reference reads of each locus within individual, given by CD-HIT-  
425 EST, were clustered with all the reference reads of all individuals, using the same  
426 threshold, in order to build the locus catalog including list of loci of all individuals and the  
427 coverage for each locus in each individual. After this step, insertion loci that matched  
428 known repeats of the Asian tiger mosquito (15) were discarded; alignments were  
429 performed with Blastn (52) using default parameters.

430 Since the quality control removed a substantial number of reads for the construction of  
431 TE insertions catalog, the raw R2 reads (TD adapter removed), that could have been  
432 discarded in a first attempt were then mapped over the catalog in order to increase the  
433 scoring sensibility. Before mapping, the raw R2 reads were also sampled at the first  
434 decile of individual coverage (as described previously). At this step, individuals that have  
435 been removed from at least two TE families for loci construction were definitively  
436 removed from the whole analysis. Mapping was performed over all insertion loci of all TE

437 families in a single run in order to prevent multiple hits. Blat was used with an identity  
438 threshold of 90 percent. Visual inspection of alignment quality over 30 sampled loci per  
439 TE family was performed in order to ensure the quality of scoring.

440 In order to check if the sampling procedure would affect our results, the read sampling  
441 procedures and subsequent analysis were performed independently 3 times (replicates  
442 M1, M2 and M3).

443

#### 444 **Genetic analyses and Genomic scan.**

445 Population structure analyses were performed independently for each TE family.

446 Principal Coordinate Analysis (PCoAs) were performed to identify genetic clusters using  
447 the ade4 package (50) of R vers. 3.2.1 (R development core team 2015). S7 coefficient  
448 of Gower and Legendre was used as a genetic distance since it gives more weight to  
449 shared insertions. Shared absences were not used because they do not give information  
450 about the genetic distance between individuals. Pairwise populations  $F_{ST}$  were computed  
451 using Arlequin 3.5 (27); significance of the index was assessed over 1000 permutations  
452 using a significance threshold of 0.05.

453 The genomic scan was performed in two steps for each of the sampling replicates of  
454 each TE. First, Bayescan 2.1 (51) was used to test for each locus deviation from  
455 neutrality. Bayescan consider a fission/island model where all subpopulations derive  
456 from a unique ancestral population. In this model, variance in allele frequencies between  
457 subpopulations is expected to be due either to the genetic drift that occurred



458 independently in each subpopulation or to selection that is a locus-specific parameter.  
459 The differentiation at each locus in each subpopulation from the ancestral population is  
460 thus decomposed into a  $\beta$  component (shared by all loci in a subpopulation) and is  
461 related to genetic drift, and a  $\alpha$  component (shared for a locus by all subpopulations)  
462 due to selection. Using a Bayesian framework, Bayescan tests for each locus the  
463 significance of the  $\alpha$  component. Rejection of the neutral model at one locus is done  
464 using posterior Bayesian probabilities and controlled for multiple testing using false  
465 discovery rate. In addition, Bayescan manage uncertainty about allele frequency from  
466 dominant data such as the TD polymorphism, leaving the  $F_{IS}$  to freely vary during the  
467 estimation of parameters. Bayescan was used with default values except for the prior  
468 odds that were set to 100 (more compatible with datasets with a large number of loci,  
469 see Bayescan manual), and a significance  $q$ -value threshold of 0.05 was used to retain  
470 outliers loci. In a second step, only outliers loci suggesting divergent directional selection  
471 between, Europe and Vietnam were considered. To identify them, locus by locus  
472 Analyses of Molecular Variance (AMOVAs) were performed using Arlequin 3.5 for each  
473 TE family. Significance of the  $F_{CT}$  (inter group differentiation) between Vietnamese and  
474 European populations was assessed performing 10,000 permutations with a significance  
475 threshold of 0.05. For each dataset, Bayescan outliers were crossed with significant  $F_{CT}$   
476 loci to retain candidate loci.  
477 To identify the genomic environment of the candidate loci, the outlier sequences  
478 (reference R2 read) were mapped onto the assembled genome of *Ae. albopictus*  
479 (28) using Blastn. Blastn alignments were performed with default parameters and sorted  
480 according to alignment score and after visual inspection of each alignment. Outlier loci

481 with multiple identical hits were discarded. To identify genes surrounding the mapped  
482 outliers, the complete transcriptome of *Ae. albopictus* (including eggs, larvae and adult  
483 females, downloaded at [www.albopictusexpression.org](http://www.albopictusexpression.org) [Armbruster *et al.*]) was mapped  
484 over the reference genome using blat with default parameters; after alignment, one best  
485 hit was retained per transcript according to the best alignment score. When a transcript  
486 had multiple best hits, all positions for the transcript were considered.

487

#### 488 **PCR validation and Outlier analyses.**

489 Pairs of primers were designed for each outlier locus in order to be used in standardized  
490 conditions. Forward primer was located in the TE end of the concerned family and  
491 reverse primer was set from the outlier locus (pairs of primers for successfully amplified  
492 insertions are provided in S3 Table). Primer pairs were first tested on a set of 10  
493 individuals in order to assess their specificity using 1/50 dilution of starting DNA from the  
494 TD experiment. Validated primers were then used to check the insertions polymorphism  
495 in 47 representatives individuals from the 8 populations studied in the TD experiment  
496 using 1/50 dilutions of the starting DNA (not all individuals could be used because of  
497 DNA limitations). All PCRs were conducted in a final volume of 25 $\mu$ L using 0.5 $\mu$ L of  
498 diluted DNA, 0.5 $\mu$ L of each primer (10 $\mu$ M), 1 $\mu$ L of dNTPs (10mM) and 1U of DreamTaq  
499 Polymerase with 1X green buffer (ThermoFisher Scientific). Amplification was performed  
500 as follows: denaturation at 94C for 2 minutes then 34 cycles including denaturation at  
501 94C for 30 seconds, hybridization at 60C for 45 seconds and elongation at 72C for 45  
502 seconds; final elongation was performed for 10 minutes at 72C. After 45 minutes

503 migration of the PCR product on 1X electrophoresis agarose gel, CG and MB assessed  
504 insertion polymorphism independently.

505

## 506 **Acknowledgements**

507 We thank Van Tran-Van, Christophe Bellet, Grégory Lambert, Huynh Kim Ly Khanh and  
508 Trang Huynh who made possible and contributed to the samplings in France and  
509 Vietnam. Library construction and sequencing was made in collaboration with Clémence  
510 Genthon and Olivier Bouchez. We are grateful to Valèria Romero Soriano and her family  
511 for their help during sampling in Sant Cugat dèl Vallès. We thank Manon Vigneron for  
512 PCR validation experiments We also thank Rita Rebollo who provided insightful  
513 comments and English revision of the manuscript. This work was performed using the  
514 computing facilities of the CC LBBE/PRABI.

## 515 **References**

- 516 1. Colautti RI, Lau JA (2015) Contemporary evolution during invasion: evidence for differentiation,  
517 natural selection, and local adaptation. *Mol Ecol* 24(9):1999–2017.
- 518 2. Lande R (2015) Evolution of phenotypic plasticity in colonizing species. *Mol Ecol*.  
519 doi:10.1111/mec.13037.
- 520 3. Peischl S, Excoffier L (2015) Expansion load: recessive mutations and the role of standing genetic  
521 variation. *Mol Ecol*. doi:10.1111/mec.13154.

- 522 4. Handley L-J, et al. (2011) Ecological genetics of invasive alien species. *BioControl* 56(4):409–428.
- 523 5. Bock DG, et al. (2015) What we still don't know about invasion genetics. *Mol Ecol*:n/a–n/a.
- 524 6. Paupy C, Delatte H, Bagny L, Corbel V, Fontenille D (2009) *Aedes albopictus*, an arbovirus vector:  
525 from the darkness to the light. *Microbes Infect* 11(14-15):1177–85.
- 526 7. Grard G, et al. (2014) Zika Virus in Gabon (Central Africa) – 2007: A New Threat from *Aedes*  
527 *albopictus*? *PLoS Negl Trop Dis* 8(2):e2681.
- 528 8. Marcondes CB, Ximenes M de FF de M (2015) Zika virus in Brazil and the danger of infestation by  
529 *Aedes* (*Stegomyia*) mosquitoes. *Rev Soc Bras Med Trop* (AHEAD). doi:10.1590/0037-8682-0220-  
530 2015.
- 531 9. Chouin-Carneiro T, et al. (2016) Differential Susceptibilities of *Aedes aegypti* and *Aedes albopictus*  
532 from the Americas to Zika Virus. *PLoS Negl Trop Dis* 10(3):e0004543.
- 533 10. Bonizzoni M, Gasperi G, Chen X, James AA (2013) The invasive mosquito species *Aedes*  
534 *albopictus*: current knowledge and future perspectives. *Trends Parasitol* 29(9):460–468.
- 535 11. Hawley WA (1988) The biology of *Aedes albopictus*. *J Am Mosq Control Assoc Suppl* 1:1–39.
- 536 12. Hawley W, Reiter P, Copeland R, Pumpuni C, Craig G (1987) *Aedes albopictus* in North America:  
537 probable introduction in used tires from northern Asia. *Science* (80- ) 236(4805):1114–1116.
- 538 13. Hanson SM, Craig GB (1994) Cold Acclimation, Diapause, and Geographic Origin Affect Cold  
539 Hardiness in Eggs of *Aedes albopictus* (Diptera: Culicidae). *J Med Entomol* 31(2):192–201.
- 540 14. Poelchau MF, Reynolds J a, Elsik CG, Denlinger DL, Armbruster P a (2013) Deep sequencing

- 541 reveals complex mechanisms of diapause preparation in the invasive mosquito, *Aedes albopictus*.  
542 *Proc Biol Sci* 280(1759):20130143.
- 543 15. Goubert C, et al. (2015) De novo assembly and annotation of the Asian tiger mosquito (*Aedes*  
544 *albopictus*) repeatome with dnaPipeTE from raw genomic reads and comparative analysis with the  
545 yellow fever mosquito (*Aedes aegypti*). *Genome Biol Evol*:evv050–.
- 546 16. Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA (2007) Rapid and cost-effective  
547 polymorphism identification and genotyping using restriction site associated DNA (RAD) markers.  
548 *Genome Res* 17(2):240–8.
- 549 17. Davey JW, et al. (2012) Special features of RAD Sequencing data: implications for genotyping. *Mol*  
550 *Ecol*. doi:10.1111/mec.12084.
- 551 18. Biedler J, Tu Z (2003) Non-LTR retrotransposons in the African malaria mosquito, *Anopheles*  
552 *gambiae*: unprecedented diversity and evidence of recent activity. *Mol Biol Evol* 20(11):1811–1825.
- 553 19. Boulesteix M, et al. (2007) Insertion polymorphism of transposable elements and population  
554 structure of *Anopheles gambiae* M and S molecular forms in Cameroon. *Mol Ecol* 16(2):441–452.
- 555 20. Santolamazza F, et al. (2008) Insertion polymorphisms of SINE200 retrotransposons within  
556 speciation islands of *Anopheles gambiae* molecular forms. *Malar J* 7(1):163.
- 557 21. Esnault C, et al. (2008) High genetic differentiation between the M and S molecular forms of  
558 *Anopheles gambiae* in Africa. *PLoS One* 3(4):e1968.
- 559 22. Bonin A, et al. (2008) A MITE-based genotyping method to reveal hundreds of DNA polymorphisms  
560 in an animal genome after a few generations of artificial selection. *BMC Genomics* 9:459.

- 561 23. Monden Y, Yamaguchi K, Tahara M (2014) Application of iPBS in high-throughput sequencing for  
562 the development of retrotransposon-based molecular markers. *Curr Plant Biol*.  
563 doi:10.1016/j.cpb.2014.09.001.
- 564 24. Casacuberta E, González J (2013) The impact of transposable elements in environmental  
565 adaptation. *Mol Ecol*:1503–1517.
- 566 25. Stapley J, Santure AW, Dennis SR (2015) Transposable elements as agents of rapid adaptation  
567 may explain the genetic paradox of invasive species. *Mol Ecol* 24(9):2241–52.
- 568 26. Roy AM, et al. (1999) Recently integrated human Alu repeats: finding needles in the haystack.  
569 *Genetica* 107(1-3):149–161.
- 570 27. Excoffier L, Lischer HEL (2010) Arlequin suite ver 3.5: a new series of programs to perform  
571 population genetics analyses under Linux and Windows. *Mol Ecol Resour* 10(3):564–7.
- 572 28. Chen X-G, et al. (2015) Genome sequence of the Asian Tiger mosquito, *Aedes albopictus*, reveals  
573 insights into its biology, genetics, and evolution. *Proc Natl Acad Sci U S A* 112(44):E5907–5915.
- 574 29. Rao PN, Rai KS (1987) Inter and intraspecific variation in nuclear DNA content in *Aedes*  
575 mosquitoes. *Heredity (Edinb)* 59(2):253–258.
- 576 30. Kumar A, Rai KS (1990) Intraspecific variation in nuclear DNA content among world populations of  
577 a mosquito, *Aedes albopictus* (Skuse). *Theor Appl Genet* 79(6):748–52.
- 578 31. Black WICI V, Hawley WA, Rai KS, Craig GB (1988) Breeding structure of a colonizing species:  
579 *Aedes albopictus* (Skuse) in peninsular Malaysia and Borneo. *Heredity (Edinb)* 61(March):439–  
580 446.

- 581 32. Kambhampati S, Black WC, Rai KS (1991) Geographic origin of the US and Brazilian *Aedes*  
582 *albopictus* inferred from allozyme analysis. *Heredity (Edinb)* 67 ( Pt 1)(September 1990):85–93.
- 583 33. Zhong D, et al. (2013) Genetic analysis of invasive *Aedes albopictus* populations in Los Angeles  
584 County, California and its potential public health impact. *PLoS One* 8(7):e68586.
- 585 34. Gupta S, Preet S (2014) Genetic differentiation of invasive *Aedes albopictus* by RAPD-PCR:  
586 Implications for effective vector control. *Parasitol Res* 113(6):2137–2142.
- 587 35. Manni M, et al. (2015) Molecular markers for analyses of intraspecific genetic diversity in the Asian  
588 Tiger mosquito, *Aedes albopictus*. *Parasit Vectors* 8(1):188.
- 589 36. Minard G, et al. (2015) French invasive Asian tiger mosquito populations harbor reduced bacterial  
590 microbiota and genetic diversity compared to Vietnamese autochthonous relatives. *Front Microbiol*  
591 6. doi:10.3389/fmicb.2015.00970.
- 592 37. Urbanelli S, Bellini R, Carrieri M, Sallicandro P, Celli G (2000) Population structure of *Aedes*  
593 *albopictus* (Skuse): the mosquito which is colonizing Mediterranean countries. *Heredity (Edinb)* 84  
594 ( Pt 3)(November 1999):331–337.
- 595 38. Birungi J, Munstermann LE (2002) Genetic Structure of *Aedes albopictus* (Diptera: Culicidae)  
596 Populations Based on Mitochondrial ND5 Sequences: Evidence for an Independent Invasion into  
597 Brazil and United States. *Ann Entomol Soc Am* 95(1):125–132.
- 598 39. Takumi K, et al. (2009) Introduction, scenarios for establishment and seasonal activity of *Aedes*  
599 *albopictus* in The Netherlands. *Vector Borne Zoonotic Dis* 9(2):191–6.
- 600 40. Becker N, et al. (2013) Repeated introduction of *Aedes albopictus* into Germany, July to October

- 601           2012. *Parasitol Res* 112(4):1787–90.
- 602   41.    Dritsou V, et al. (2015) A draft genome sequence of an invasive mosquito: an Italian *Aedes*  
603           *albopictus*. *Pathog Glob Health*:2047773215Y0000000031.
- 604   42.    Urbanski J, et al. (2012) Rapid adaptive evolution of photoperiodic response during invasion and  
605           range expansion across a climatic gradient. *Am Nat* 179(4):490–500.
- 606   43.    Munroe DJ, et al. (1994) IRE-bubble PCR: a rapid method for efficient and representative  
607           amplification of human genomic DNA sequences from complex sources. *Genomics* 19(3):506–14.
- 608   44.    Akkouche A, et al. (2012) tirant, a newly discovered active endogenous retrovirus in *Drosophila*  
609           *simulans*. *J Virol* 86(7):3675–81.
- 610   45.    Carnelossi EAG, et al. (2014) Specific activation of an I-like element in *Drosophila* interspecific  
611           hybrids. *Genome Biol Evol* 6(7):1806–17.
- 612   46.    Modolo L, Lerat E (2015) UrQt: an efficient software for the Unsupervised Quality trimming of NGS  
613           data. *BMC Bioinformatics* 16(1):137.
- 614   47.    Martin M (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.  
615           *EMBnet.journal* 17(1):10.
- 616   48.    Kent WJ (2002) BLAT---The BLAST-Like Alignment Tool. *Genome Res* 12(4):656–64.
- 617   49.    Li W, Godzik A (2006) Cd-hit: a fast program for clustering and comparing large sets of protein or  
618           nucleotide sequences. *Bioinformatics* 22(13):1658–9.



619 50. Dray S, Dufour A-B (2007) The ade4 Package: Implementing the Duality Diagram for Ecologists. *J*  
620 *Stat Softw* 22(4):1–20.

621 51. Foll M, Gaggiotti O (2008) A genome-scan method to identify selected loci appropriate for both  
622 dominant and codominant markers: a Bayesian perspective. *Genetics* 180(2):977–93.

623

624

## 625 **Data Accessibility**

626 Paired-end raw sequences are available through SRA at NCBI under SRP070185  
627 (Bioproject PRJNA312147)

628 Final presence/absence matrices (including replicates) are available at Dryad  
629 (doi:10.5061/dryad.9p925) at <http://datadryad.org/review?doi=doi:10.5061/dryad.9p925>

630

## 631 **Author contributions**

632 CG, CV and MB conceived the experiments and conducted the analyses. CG and HH  
633 developed and performed the molecular experiments. GM, CVM and PM conducted the  
634 sampling in France and Vietnam. All authors contributed to the final version of the  
635 manuscript.

636

## 637 **SI Figures**

638 **S1 Figure. Insertion frequencies and locus coverage.** Relationship between mean  
639 locus coverage (mean number of read per individual at one locus when present) and  
640 insertion frequency (number of individuals that share the locus), M1 replicate (others  
641 have similar results). Red dots are outlier loci.

642 **S2 Figure. PcoAs.** Projection of individuals over the three first principal coordinates  
643 (PC) of principal Coordinates Analyses (PCoAs) computed over all loci for each of the  
644 TE family (replicates M2 and M3). Proportion of inertia represented by each axis is  
645 noted in %. ◦ • European populations; triangles: Vietnamese populations

646 **S3 Figure. Insertion frequencies in Europe and Vietnam.** Insertion Frequencies of 12  
647 outlier loci for which specific PCRs were performed. Darker bars are the results of  
648 bioinformatic analysis and lighter are the frequencies obtained over 47 individuals used  
649 for PCR validations.

650 **S4 Figure. Library preparation for high throughput sequencing of TD products.**  
651 Genomic DNA was sheared using *HindIII*. Blue parts are copies of TE interspersed in  
652 the genome. Y shaped TD adapters were then ligated to cohesive ends. For each TE  
653 family and for each individual, PCR 1 was performed using a TE specific primer  
654 annealing to the end of each TE copy; subsequent elongation completed the  
655 complementary 3' end of the MSEB part of the adapter, allowing annealing of the LNP  
656 primer. For each TE family, three independent PCR were performed, and all PCR 1  
657 products from one individual were pooled in one unique sample. Size selection using

658 magnetic bead was then done for each individual before normalization. Finally, PCR2  
659 was performed for each individual pool in order to ligate indexes and Illumina adapters.

660 **S5 Figure. Bioinformatic workflow.** From sequencing to outlier analysis. Details are  
661 given in Material and Methods. R1 reads include the terminal end of one TE copy (in  
662 blue) and R2 reads include flanking region (black) and the TD adapter (white and  
663 green). After cleaning, R2 flanking regions are sampled to account for the coverage  
664 heterogeneity among individuals (sampling 1). Individual clustering of these reads,  
665 reference sequences (colored) of each insertion loci are then clustered between  
666 individuals to build the insertion catalog. Raw R2 reads are then sampled and mapped  
667 over the catalogs (1 catalog per TE family) and individual coverage per locus is  
668 calculated. Loci are then filtered for insertion frequency and coverage before population  
669 genetics and outlier analysis on 1/0 matrices. Steps surrounded in purple dashed line  
670 are replicated three times. At the end, all outliers from all replicates are recovered  
671 (candidate loci total).

672

## 673 **SI Tables**

674 **S1 Table.** Total number of loci recovered per TE family and insertion frequencies.

675 **S2 Table.** Sampling information about the *Ae. albopictus* population studied.

676 **S3 Table.** Oligonucleotide sequences used during TD library construction and outliers  
677 validation experiments.

678

## 679 **SI files**

680 **S1 File. Estimate of pairwise  $F_{ST}$  for the three replicates of read sampling.** For each

681 TE family the first table is the  $F_{ST}$  estimate and the second the pairwise  $P$ -value of  $F_{ST}$

682 estimate. Population names are: 1=BCN; 2=CGN; 3=NCE; 4=PLB; 5=SP; 6=HCM;

683 7=TA; 8=VT

684 **S2 File. Outlier genomic position and annotation.** Table for the contigs harboring out-

685 lier loci (three spreadsheets).