

1 **The sequence of a male-specific genome region containing the sex**
2 **determination switch in *Aedes aegypti***

3

4 Joe Turner ^{1,2¶}, Ritesh Krishna ^{1#a¶}, Arjen E. van't Hof ^{1#b¶}, Elizabeth R. Sutton ^{2,3#c}, Kelly
5 Matzen ², Alistair C. Darby ^{1*}

6

7

8 1. Centre for Genomic Research, Institute of Integrative Biology, University of Liverpool,
9 Crown Street, Liverpool, L69 7ZB, UK.

10 2. Oxitec Ltd., 71 Innovation Drive, Milton Park, Abingdon, OX14 4RQ, UK.

11 3. Department of Zoology, University of Oxford, South Parks Road, Oxford, OX1 3PS, UK.

12

13 Current Addresses:

14 #a IBM Research UK, STFC Daresbury Laboratory, Warrington, WA4 4AD, UK.

15 #b Liverpool School of Tropical Medicine, Pembroke Place, Liverpool, L3 5QA, UK.

16 #c Sismic, West of Scotland Science Park, Glasgow, G20 0SP, UK.

17

18 * Corresponding author

19 acdarby@liverpool.ac.uk

20

21 ¶ Authors contributed equally to this work.

22

23 *Aedes aegypti* is the principal vector of several important arboviruses. Among the methods of
24 vector control to limit transmission of disease are genetic strategies that involve the release of
25 sterile or genetically modified non-biting males (Alphey 2014); which has generated interest
26 in manipulating mosquito sex ratios (Gilles *et al.* 2014; Adelman and Tu 2016). Sex
27 determination in *Ae. aegypti* is controlled by a non-recombining Y chromosome-like region
28 called the M locus (Craig *et al.* 1960), yet characterisation of this locus has been thwarted by
29 the repetitive nature of the genome (Hall *et al.* 2015). In 2015, an M locus gene named *Nix*
30 was identified that displays the qualities of a sex determination switch (Hall *et al.* 2015).
31 With the use of a whole-genome BAC library, we amplified and sequenced a ~200kb region
32 containing this male-determining gene. In this study, we show that *Nix* is comprised of two
33 exons separated by a 99kb intron, making it an unusually large gene. The intron sequence is
34 highly repetitive and exhibits features in common with old Y chromosomes, and we speculate
35 that the lack of recombination at the M locus has allowed the expansion of repeats in a
36 manner characteristic of a sex-limited chromosome, in accordance with proposed models of
37 sex chromosome evolution in insects.

38

39 At least 2.5 billion people live in areas where they are at risk of dengue transmission from
40 mosquitoes, principally *Ae. aegypti*, with an estimated 390 million infections per year
41 (Laughlin *et al.* 2012; Bhatt *et al.* 2013). Recently, the emergence of chikungunya and Zika
42 viruses further highlights the public health importance of *Ae. aegypti* (Musso *et al.* 2015;
43 Fauci and Morens 2016). Future mosquito control strategies may incorporate genetic
44 techniques such as the sustained release of sterile or transgenic “self-limiting” mosquitoes
45 (Alphey *et al.*, 2013; WHO: <https://goo.gl/FRqJ0d>). Given that only female mosquitoes bite
46 and spread disease, there has been substantial interest in manipulating mosquito sex
47 determination using these genetic techniques and others, including gene drive (Adelman and

48 Tu 2016; Hoang *et al.* 2016). Therefore, elucidating the genetic basis for sex determination
49 could, for instance, facilitate production of male-only cohorts for release, or allow
50 transformation of mosquitoes with sex-specific “self-limiting” gene cassettes.

51 Sex determination in insects is variable, and generally not well understood outside of model
52 species (Charlesworth and Mank 2010). Unlike the malaria mosquito *Anopheles gambiae* and
53 *Drosophila* species, *Ae. aegypti* does not have heteromorphic (XY) sex chromosomes (Craig
54 *et al.* 1960). Instead, the male phenotype is determined by a non-recombining M locus on one
55 copy of autosome 1 (Newton *et al.* 1978; Clements 1992; Toups and Hahn 2010). This locus
56 is poorly characterised because its highly repetitive nature has confounded attempts to study
57 it based on the existing genome assembly (Hall *et al.* 2015). The 1,376Mb *Ae. aegypti*
58 genome was assembled from Sanger sequencing reads in 2007 (Nene *et al.* 2007), which are
59 commonly not long enough to span the repetitive transposable elements that comprise a large
60 proportion of the genome (Koren and Phillippy 2015). Consequently, the current assembly is
61 still relatively low quality (Severson and Behura 2012). Furthermore, the fact that both male
62 and female genomic DNA was used for genome sequencing reduces the expected coverage of
63 the M locus to one quarter of the autosome 1 sequences, further obscuring candidate M locus
64 sequences (Hall *et al.* 2014).

65 Recently, a team of researchers was nevertheless able to identify *Nix*, a gene with male-
66 specific, early embryonic expression. Knockout of *Nix* using CRISPR/Cas9 results in
67 morphological feminisation of male mosquitoes along with feminisation of gene expression
68 and female splice forms of the conserved sex-regulating genes *doublesex* (*dsx*) and *fruitless*
69 (*fru*), strongly indicating that *Nix* is the upstream regulator of sexual differentiation (Hall *et*
70 *al.* 2015). The translated *Nix* protein contains two RNA recognition motifs and is
71 hypothesised to be a splicing factor, acting either directly on *dsx* and *fru* or on currently
72 unknown intermediates (Adelman and Tu 2016). A comparison of sexually dimorphic gene

73 expression in different mosquito tissue types also detected male-specific transcripts of *Nix*
74 (Matthews *et al.* 2016). An ortholog of *Nix* is present in *Ae. albopictus*, but it is not known if
75 the two are functionally homologous (Chen *et al.* 2015).

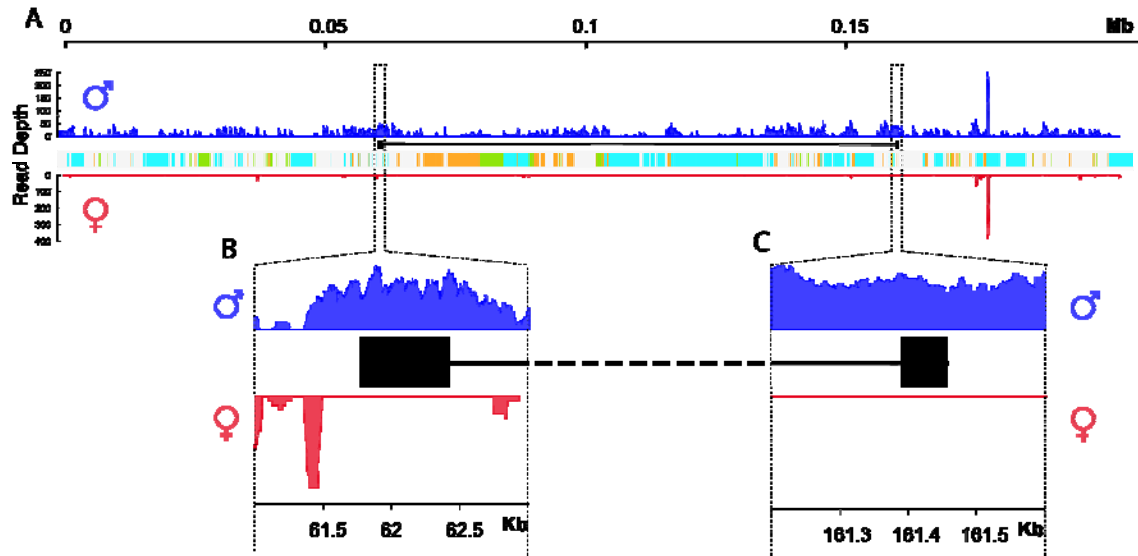
76 To date, *Nix* has only been characterised as an mRNA transcript. To fully understand this
77 gene's role in sex determination and to utilise this knowledge for vector control, it is essential
78 to decipher its genomic context. For this purpose, this study identifies and describes the
79 region of the M-locus in which *Nix* is located.

80

81 Four BAC clones positive for *Nix* assembled into a single region of 207 kb with no gaps and
82 a GC content of 40.2% (submitted to the NCBI as accession KY849907). The presence of the
83 *Nix* gene in the assembled BACS was confirmed by BLASTN. The whole gene was present
84 in tiled BACs, though not completely within individual BAC clones. Neither *Nix* nor the
85 complete region could be found in the AegL3 or Aag2 reference genome assemblies. While
86 *Nix* was originally identified in the genome-sequenced Liverpool strain (Hall *et al.* 2015),
87 PCR revealed that it is exclusively present in male genomic DNA from other geographically
88 varied *Ae. aegypti* populations (Figure S1), further strengthening the evidence that it is
89 wholly present in the M locus.

90 The *Nix* gene was found to be made up of two exons with a single intron of 99 kb (Figure 1).
91 Although large introns are not uncommon in *Ae. aegypti* (average intron length ~5000
92 bp)(Nene *et al.* 2007), this intron is at the extreme end of intron sizes observed (Figure S2),
93 especially considering the small size of its protein coding regions (<1000 bp). The gene
94 structure is confirmed by Illumina RNA-Seq data clearly showing reads spanning the intron
95 between the two exons (Figure 1). RepeatMasker identified approximately 55% of the
96 sequenced region as repetitive, and the intron region of *Nix* as 72% repetitive (Table S1).

97



98

99

100 **Figure 1: Structure and gene expression of the ~207 kb genomic region containing the**
101 ***Nix* gene.** *Nix* is shown as two black boxes representing the exons, joined by a black line
102 representing the intron. Colours on the central track of **A** represent the classes of repetitive
103 elements (orange: DNA transposons; cyan: Gypsy LTRs; green: Ty1/Copia LTRs). Blue
104 histograms represent the coverage of RNA-Seq reads from male samples on the y axis; red
105 histograms represent the coverage from female samples. **B** and **C** show enlargements of the
106 first and second exons of *Nix* in the dotted regions in **A**, respectively.

107

108 The genomic data from our assembled M locus region show that *Nix* is approximately 100 kb
109 in length – exceptionally long even for an insect, and one of the longest in the mosquito
110 genome. This is particularly unusual because *Nix* is expressed in early embryonic
111 development, before the onset of the syncytial blastoderm stage 3-4 hours after oviposition
112 (Hall *et al.* 2015), during which time most active genes have very short introns, or lack them
113 entirely. There is evidence of selection against intron presence in genes expressed in the early
114 *Ae. aegypti* zygote (Biedler *et al.* 2012). In *Drosophila*, the majority of early-expressed genes
115 have small introns and encode small proteins, suggesting that selection has favoured high
116 transcript turnover during early embryonic development due to the requirement for short cell
117 cycles and rapid division (Artieri and Fraser 2014). It might therefore be expected that
118 selection would limit the *Nix* intron's expansion to preserve efficient transcription in the
119 zygote.

120 One possible explanation is the expansion of repetitive DNA. The RepeatMasker results
121 reveal that the *Nix* region contains a high number of repetitive sequences, especially
122 retrotransposons (Figure 1; Table S1). The M locus has accumulated repeats in between
123 protein-coding DNA in a manner characteristic of a sex chromosome, which are prone to
124 degeneration by Muller's ratchet due to the lack of recombination (Muller 1964;
125 Charlesworth 1991; Kaiser and Bachtrog 2010). For instance, repetitive sequences comprise
126 almost the entire *Anopheles gambiae* Y chromosome, and these repetitive sequences show
127 rapid evolutionary divergence (Hall *et al.* 2016). Similarly, genes on the *Drosophila* Y
128 chromosome, such as those involved in spermatogenesis, have gigantic repetitive introns,
129 sometimes in the megabase range, that consequently make them many times larger than
130 typical autosomal genes (Carvalho *et al.* 2001; Bachtrog 2013).

131 It is therefore possible that the lack of recombination may pose constraints on the structure of
132 the M locus, and in the absence of strong selection the *Nix* gene has degenerated outside the
133 coding regions. Non-recombining sex loci such as the *Ae. aegypti* M locus may represent an
134 evolutionary precursor to differentiated sex chromosomes, which are thought to emerge when
135 sexually antagonistic alleles accumulate on either chromosome and favour reduced
136 recombination between the two homologs, eventually leading to degeneration and loss of
137 genes on the proto-Y (Charlesworth *et al.* 2005). Recent data appears to show that
138 recombination is reduced along autosome 1 even outside of the M locus (Fontaine *et al.*
139 2016), while the fully differentiated *Anopheles* X and Y chromosomes still display some
140 degree of recombination with each other (Hall *et al.* 2016). Thus, *Ae. aegypti* may be “further
141 along” this evolutionary trajectory than previously assumed.

142 The *Ae. aegypti* M locus provides an intriguing example of the complexity of evolutionary
143 forces acting on sex chromosomes, and further study of the locus will contribute to
144 understanding the evolution of sex determination in insects and address general questions
145 about the factors impacting gene and genome length. Importantly, these may also yield
146 insights that can be applied to increase the efficiency of genetic strategies for vector control.

147

148 **Methods**

149 **BAC library construction**

150 A BAC library of insert size 130 kb was constructed (Amplicon Express, USA) for an
151 estimated coverage of ~5x for autosomal regions (~2.5x for sex specific regions) from a
152 DNA pool of approximately 50 sibling males. The male siblings were from one family from
153 an Asian wild type laboratory strain after five generations of full-sib mating. Superpools and
154 matrixpools were supplied to allow PCR based screening of the BAC library.

155 **BAC library screening, isolation and sequencing**

156 The BAC library was PCR screened using primers (Nix1F 3'-
157 TTGAGTCTGAAAAGTCTATGCAA-5', Nix1R 3'-TCGCTCTTCCGTGGCATTGA-5',
158 Nix2F 3'-ACGTAGTCGGCAACTCGAAG-5', Nix2R 3'-
159 CTGGGACAAATCGAACGGAA-5') based on the complete coding sequence of *Nix*
160 (GenBank accession number KF732822). The first primer set was also used to screen for *Nix*
161 in the genomic DNA of six male and six female individuals each from two wildtype *Ae.*
162 *aegypti* strains.
163 Screening of the library resulted in four positive clones - two for each primer pair. These
164 BAC clones were propagated, extracted using a Maxiprep kit (Qiagen, UK), pooled before
165 SMRTbell library preparation (PacBio, USA), and sequenced on a single SMRTcell using
166 P6-C3 chemistry on the PacBio RS II platform (PacBio, USA).

167 **Data analysis**

168 The sequence data was trimmed to remove vector sequences and adaptors prior to assembly
169 with the CANU v1 assembler (Berlin *et al.* 2015), followed by sequence polishing with
170 QUIVER.
171 BLASTN was used to assess the uniqueness of the assembled *Nix* region compared to the
172 *Aedes aegypti* Liverpool reference genome AaegL3 and the newer Aag2 cell line assembly.
173 Illumina data generated from male and female genomic DNA (accession numbers
174 SRX290472 and SRX290470) and RNA (accession numbers SRX709698-SRX709703) were
175 mapped to a combined reference containing the assembled *Nix* region added to the AaegL3
176 genome. DNA samples were mapped with BOWTIE 2.2.1 (using default parameters with -I
177 200 and -X 500) and RNA-Seq data with TOPHAT 2.1.1 version (using default parameters).
178 RNA-Seq data was processed using the CUFFLINKS 2.2.1 pipeline to look for potential
179 genes and male/female specific expression from the region.

180 Genes were predicted using AUGUSTUS and the *Aedes aegypti* model (Nene *et al.* 2007),
181 repetitive regions described using REPEATMASKER 4.0.6 and the *Ae. aegypti* repeat
182 database.

183

184 **Supplementary Information** is available in the online version of the paper.

185

186 **Author contributions**

187 J.T., R.K. and A.E.v.H. contributed equally to this work. K.M. and A.C.D. designed the study
188 and obtained funding, with contribution from J.T.; K.M. provided mosquito samples; E.R.S.
189 and A.C.D. commissioned the BAC library construction; A. E. v. H. and J. T. screened the
190 BAC library and extracted DNA; A. E. v. H. performed BAC scaffolding; A.C.D. oversaw
191 sequencing and assembled the DNA sequence; R.K. performed the mapping and developed
192 computational strategies for data analysis; J.T. performed the repeat masking; J.T. and
193 A.C.D. wrote the paper, with contribution from A. E. v. H.; R.K. and A.C.D. produced the
194 figures.

195

196 **Acknowledgments**

197 This work was funded by BBSRC PhD training grant BB/M503460/1 (J.T. & A.C.D.) and a
198 BBSRC grant BB/M001512/1 (K.M. & A.C.D.).

199 The PacBio sequencing was conducted at the Centre for Genomics Research, University of
200 Liverpool with the assistance of Dr Margaret Hughes and Dr John Kenny.

201 We thank Dr Andrea Betancourt and Dr Ilik Saccheri for comments on the manuscript.

202

203 **References**

204 Adelman Z. N., Tu Z., 2016 Control of mosquito-borne infectious diseases: sex and gene

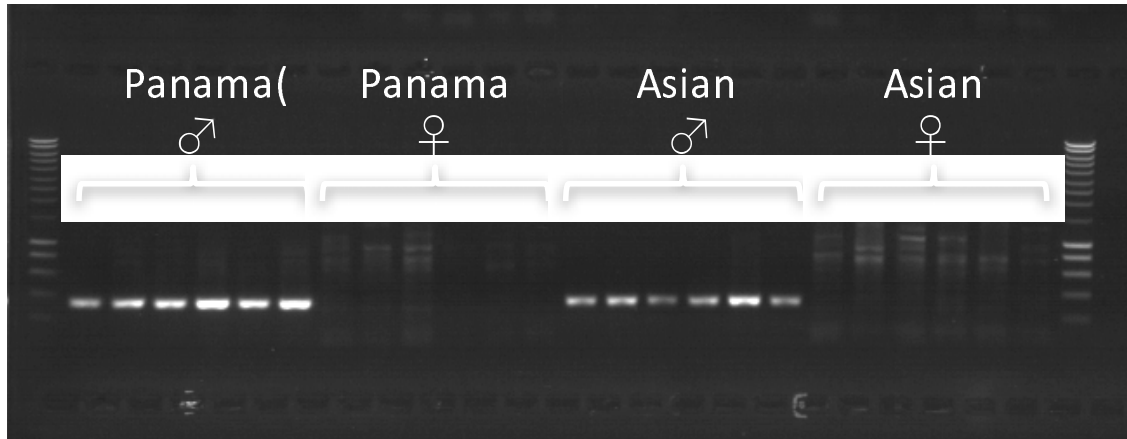
- 205 drive. Trends Parasitol. 32: 219–229.
- 206 Alphey L., McKemey A., Nimmo D., Neira Oviedo M., Lacroix R., *et al.*, 2013 Genetic
207 control of *Aedes* mosquitoes. Pathog. Glob. Health 107: 170–179.
- 208 Alphey L., 2014 Genetic control of mosquitoes. Annu. Rev. Entomol. 59: 205–224.
- 209 Artieri C. G., Fraser H. B., 2014 Transcript length mediates developmental timing of gene
210 expression across *Drosophila*. Mol. Biol. Evol. 31: 2879–2889.
- 211 Bachtrog D., 2013 Y-chromosome evolution: emerging insights into processes of Y-
212 chromosome degeneration. Nat. Rev. Genet. 14: 113–124.
- 213 Berlin K., Koren S., Chin C.-S., Drake J. P., Landolin J. M., *et al.*, 2015 Assembling large
214 genomes with single-molecule sequencing and locality-sensitive hashing. Nat.
215 Biotechnol. 33: 623–630.
- 216 Bhatt S., Gething P. W., Brady O. J., Messina J. P., Farlow A. W., *et al.*, 2013 The global
217 distribution and burden of dengue. Nature 496: 504–507.
- 218 Biedler J. K., Hu W., Tae H., Tu Z., 2012 Identification of early zygotic genes in the yellow
219 fever mosquito *Aedes aegypti* and discovery of a motif involved in early zygotic genome
220 activation. PLoS One 7: e33933.
- 221 Carvalho A. B., Dobo B. A., Vibranovski M. D., Clark A. G., 2001 Identification of five new
222 genes on the Y chromosome of *Drosophila melanogaster*. Proc. Natl. Acad. Sci. USA
223 98: 13225–13230.
- 224 Charlesworth B., 1991 Evolution of sex chromosomes. Science 251: 1030–1033.
- 225 Charlesworth D., Charlesworth B., Marais G., 2005 Steps in the evolution of heteromorphic
226 sex chromosomes. Heredity (Edinb). 95: 118–128.
- 227 Charlesworth D., Mank J. E., 2010 The birds and the bees and the flowers and the trees:
228 Lessons from genetic mapping of sex determination in plants and animals. Genetics 186:
229 9–31.

- 230 Chen X.-G., Jiang X., Gu J., Xu M., Wu Y., *et al.*, 2015 Genome sequence of the Asian Tiger
231 mosquito, *Aedes albopictus*, reveals insights into its biology, genetics, and evolution.
232 Proc. Natl. Acad. Sci.: 201516410.
- 233 Clements A. N., 1992 *The Biology of Mosquitoes*. Chapman & Hall, London.
- 234 Craig G. B., Hickey W. A., Vandehey R. C., 1960 An inherited male-producing factor in
235 *Aedes aegypti*. Science 132: 1887–1889.
- 236 Fauci A. S., Morens D. M., 2016 Zika Virus in the Americas — Yet Another Arbovirus
237 Threat. N. Engl. J. Med. 374: 601–604.
- 238 Fontaine A., Filipović I., Fansiri T., Hoffmann A. A., Rašić G., *et al.*, 2016 Cryptic genetic
239 differentiation of the sex-determining chromosome in the mosquito *Aedes aegypti*.
240 bioRxiv.
- 241 Gilles J. R. L., Schetelig M. F., Scolari F., Marec F., Capurro M. L., *et al.*, 2014 Towards
242 mosquito sterile insect technique programmes: exploring genetic, molecular, mechanical
243 and behavioural methods of sex separation in mosquitoes. Acta Trop. 132: S178-187.
- 244 Hall A. B., Timoshevskiy V. A., Sharakhova M. V., Jiang X., Basu S., *et al.*, 2014 Insights
245 into the preservation of the homomorphic sex-determining chromosome of *Aedes*
246 *aegypti* from the discovery of a male-biased gene tightly linked to the M-locus. Genome
247 Biol. Evol. 6: 179–191.
- 248 Hall A. B., Basu S., Jiang X., Qi Y., Timoshevskiy V. A., *et al.*, 2015 A male-determining
249 factor in the mosquito *Aedes aegypti*. Science 348: 1268–70.
- 250 Hall A. B., Papathanos P.-A., Sharma A., Cheng C., Akbari O. S., *et al.*, 2016 Radical
251 remodeling of the Y chromosome in a recent radiation of malaria mosquitoes. Proc.
252 Natl. Acad. Sci. 113: 201525164.
- 253 Hoang K. P., Teo T. M., Ho T. X., Le V. S., 2016 Mechanisms of sex determination and
254 transmission ratio distortion in *Aedes aegypti*. Parasit. Vectors 9: 49.

- 255 Kaiser V. B., Bachtrog D., 2010 Evolution of sex chromosomes in insects. *Annu. Rev. Genet.*
256 44: 91–112.
- 257 Koren S., Phillippy A. M., 2015 One chromosome, one contig: complete microbial genomes
258 from long-read sequencing and assembly. *Curr. Opin. Microbiol.* 23: 110–120.
- 259 Laughlin C. A., Morens D. M., Cassetti M. C., Costero-Saint Denis A., San Martin J. L., *et*
260 *al.*, 2012 Dengue research opportunities in the Americas. *J. Infect. Dis.* 206: 1121–1127.
- 261 Matthews B. J., McBride C. S., DeGennaro M., Despo O., Vosshall L. B., 2016 The
262 neurotranscriptome of the *Aedes aegypti* mosquito. *BMC Genomics* 17: 32.
- 263 Muller H. J., 1964 The relation of recombination to mutational advance. *Mutat. Res.* 1: 2–9.
- 264 Musso D., Cao-Lormeau V. M., Gubler D. J., 2015 Zika virus: following the path of dengue
265 and chikungunya? *Lancet* 386: 243–244.
- 266 Nene V., Wortman J. R., Lawson D., Haas B., Kodira C., *et al.*, 2007 Genome sequence of
267 *Aedes aegypti*, a major arbovirus vector. *Science* 316: 1718–1723.
- 268 Newton M. E., Wood R. J., Southern D. I., 1978 Cytological mapping of the M and D loci in
269 the mosquito, *Aedes aegypti* (L.). *Genetica* 48: 137–143.
- 270 Severson D. W., Behura S. K., 2012 Mosquito genomics: progress and challenges. *Annu.*
271 *Rev. Entomol.* 57: 143–166.
- 272 Toups M. A., Hahn M. W., 2010 Retrogenes reveal the direction of sex-chromosome
273 evolution in mosquitoes. *Genetics* 186: 763–766.
- 274
- 275
- 276

277 **Supplementary information**

278



279

280

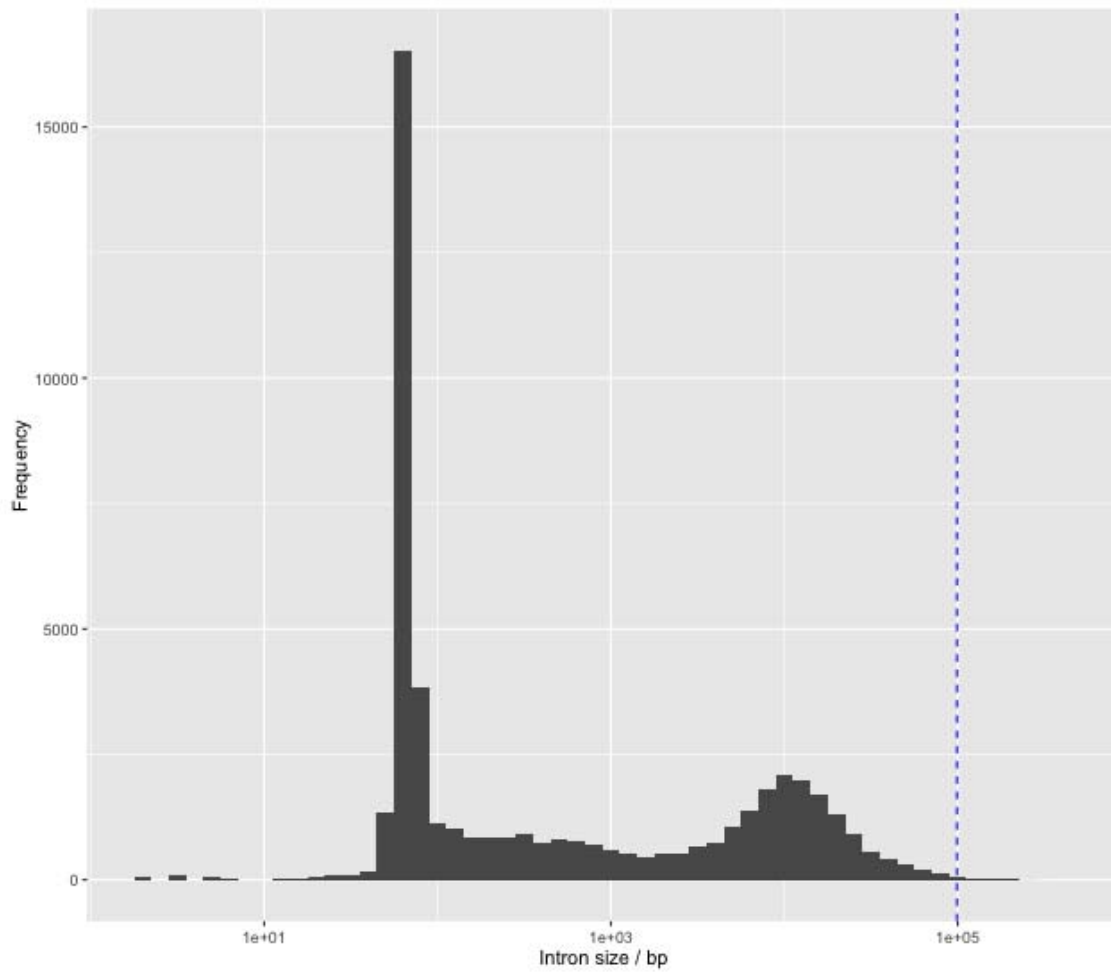
281 **Figure S1: PCR screening of the M locus gene *Nix* in male and female DNA of wild type**

282 *Aedes aegypti* strains. Primers used were Nix1F (3'-TTGAGTCTGAAAAGTCTATGCAA-

283 5') and Nix1R (3'-TCGCTCTTCCGTGGCATTGGA-5'), targeting *Nix* exon 1.

284

285
286



287
288
289
290
291
292

Figure S2: Intron size distribution in *Aedes aegypti* Liverpool reference genome

AaegL3. Blue dashed line indicates the size of the *Nix* intron relative other introns. X axis is transformed by \log_{10} .

293 **Table S1:** Types and abundance of repeats in the 207kb assembled M locus region and 99kb

294 *Nix* intron, identified by RepeatMasker using the *Aedes aegypti* repeat library.

295

Repeat Type	Entire region		<i>Nix</i> intron region	
	Number of elements	Percentage of sequence	Number of elements	Percentage of sequence
Retroelements	105	42.1%	49	51.0%
SINEs	8	0.81%	5	1.11%
Penelope	3	0.08%	2	0.20%
LINES	24	5.43%	6	6.85%
L2/CR1/Rex	4	0.13%	0	0%
R1/L0A/Jockey	13	3.87%	3	6.60%
RTE/Bov-B	3	1.33%	0	0%
L1/CIN4	1	0.02%	1	0.05%
LTR Elements	73	35.8%	38	43.0%
BEL/Pao	9	0.71%	3	0.87%
Ty1/Copia	16	11.3%	14	19.2%
Gypsy/DIRS1	48	23.8%	21	23.0%
DNA transposons	97	11.7%	69	20.1%
Tc1-IS630-Pogo	11	3.87%	11	9.04%
Other (Mirage, P-element, Transib)	1	0.06%	0	0%
Unclassified	6	0.48%	3	0.22%
Small RNA	8	0.81%	5	1.11%
Satellites	1	0.75%	0	0%
Simple repeats	19	0.34%	7	0.24%
Low complexity	3	0.07%	1	0.04%
Total repeats		55.4%		71.6%

296

