

1 **Candidate-gene based GWAS identifies reproducible DNA markers for**  
2 **metabolic pyrethroid resistance from standing genetic variation in East**  
3 **African *Anopheles gambiae***

4 David Weetman<sup>1†\*</sup>, Craig S. Wilding<sup>2†</sup>, Daniel E. Neafsey<sup>3</sup>, Pie Müller<sup>1,4,5</sup>, Eric Ochomo<sup>6,7</sup>,  
5 Alison T. Isaacs<sup>1</sup>, Keith Steen<sup>1</sup>, Emily J. Rippon<sup>1</sup>, John C. Morgan<sup>1</sup>, Henry D. Mawejje<sup>8</sup>, Daniel  
6 J. Rigden<sup>9</sup>, Loyce M. Okedi<sup>10</sup> and Martin J. Donnelly<sup>1,11</sup>

7 <sup>1</sup>Department of Vector Biology, Liverpool School of Tropical Medicine, Pembroke Place, Liverpool,  
8 UK

9 <sup>2</sup>School of Natural Sciences and Psychology, Liverpool John Moores University, Liverpool, UK

10 <sup>3</sup>Genome Sequencing and Analysis Program, Broad Institute of MIT and Harvard, Cambridge, USA

11 <sup>4</sup>Epidemiology and Public Health Department, Swiss Tropical and Public Health Institute, Basel,  
12 Switzerland

13 <sup>5</sup>University of Basel, Basel, Switzerland

14 <sup>6</sup>School of Public Health and Community Development, Maseno University, Maseno, Kenya

15 <sup>7</sup>KEMRI/CDC Research and Public Health Collaboration, Kisumu, Kenya

16 <sup>8</sup>Infectious Diseases Research Collaboration, Kampala, Uganda

17 <sup>9</sup>Institute of Integrative Biology, University of Liverpool, Liverpool, UK

18 <sup>10</sup>National Livestock Resources Research Institute, Tororo, Uganda

19 <sup>11</sup>Malaria Programme, Wellcome Trust Sanger Institute, Hinxton, Cambridge, UK

20 \*Correspondence: David Weetman. email: david.weetman@lstmed.ac.uk. fax: +44 (0) 151 705 3369

21 †these authors contributed equally to this work

22 **Keywords:** mosquito, insecticide resistance, soft sweep, malaria, permethrin, deltamethrin, lambda-  
23 cyhalothrin, humidity, bioassay, Uganda, Kenya

24

25 **Running title:** Metabolic resistance DNA markers in *An. gambiae*

26

27 **Abstract**

28 Metabolic resistance to pyrethroid insecticides is widespread in *Anopheles* mosquitoes and is a  
29 major threat to malaria control. DNA markers would aid predictive monitoring of resistance, but few  
30 mutations have been discovered outside of insecticide-targeted genes. Isofemale family pools from a  
31 wild Ugandan *Anopheles gambiae* population, from an area where operational pyrethroid failure is  
32 suspected, were genotyped using a candidate-gene enriched SNP array. Resistance-associated SNPs  
33 were detected in three genes from detoxification superfamilies, in addition to the insecticide target  
34 site (the Voltage Gated Sodium Channel gene, *Vgsc*). The putative associations were confirmed for  
35 two of the marker SNPs, in the P450 *Cyp4j5* and the esterase *Coae1d* by reproducible association  
36 with pyrethroid resistance in multiple field collections from Uganda and Kenya, and together with  
37 the *Vgsc*-1014S (*ldr*) mutation these SNPs explained around 20% of variation in resistance.  
38 Moreover, the >20 Mb 2La inversion also showed evidence of association with resistance as did  
39 environmental humidity. Sequencing of *Cyp4j5* and *Coae1d* detected no resistance-linked loss of  
40 diversity, suggesting selection from standing variation. Our study provides novel, regionally-  
41 validated DNA assays for resistance to the most important insecticide class, and establishes both 2La  
42 karyotype variation and humidity as common factors impacting the resistance phenotype.

## 43 Introduction

44 Insecticide resistance is a strongly selected, highly heritable trait that represents an excellent model  
45 system for the study of contemporary evolution<sup>1</sup>. Disease-transmitting mosquitoes are often subject  
46 to intense insecticidal selection pressure and in many instances there are critical financial and health  
47 implications of insecticide resistance evolution. In sub-Saharan Africa, where over 90% of all malaria  
48 fatalities occur, cases have fallen dramatically, due largely to the massive scale-up of insecticidal  
49 interventions such as treated bednets and indoor residual spraying<sup>2</sup>. Few insecticides are licenced  
50 for vector control and extensive application in public health, agriculture and household pesticide  
51 formulations has selected for widespread resistance to pyrethroids, by far the most important class  
52 of insecticides which may be used on bednets<sup>3</sup>. Whilst epidemiological evidence for an impact of  
53 resistance on malaria remains sparse<sup>4</sup> there exists the worrying spectre that the recent gains in  
54 malaria control may be lost<sup>5</sup>. Preservation of pyrethroid efficacy is thus a primary objective for  
55 malaria control programmes, which need to apply effective insecticide resistance management.  
56 Insecticide resistance management depends on timely information on changes in resistance and the  
57 potential for cross resistance between insecticide options. High throughput molecular diagnostics  
58 can provide more sensitive, rapid and geographically-widespread assessments of changes in  
59 resistance than phenotypic assessments<sup>6</sup>. Knowledge of the genetic mechanisms underlying  
60 metabolic cross-resistance, which is far more difficult to predict a priori than cross-resistance at  
61 shared target sites<sup>7</sup>, is vital to allow implementation of insecticide combinations which will slow  
62 resistance development or allow susceptibility to return. Thus, improved knowledge of pyrethroid  
63 resistance mechanisms and provision of DNA markers for rapid predictive diagnosis<sup>8</sup> and their  
64 evolution is an important goal to aid future insecticide deployment strategies..

65 Insecticide resistance occurs through four general mechanisms: insecticide target site alterations;  
66 elevated insecticide metabolism/sequestration; which are the primary focus here, and, potentially  
67 important but currently less-well described, reduced cuticular penetration and behavioural  
68 avoidance<sup>9</sup>. Target site resistance typically occurs through amino acid-altering substitutions, with  
69 the same mutations often found across diverse insect taxa<sup>10,11</sup>, or through copy number variation of  
70 alleles harbouring resistance mutations<sup>12,13</sup>. In *Anopheles* mosquitoes, and insects generally,  
71 pyrethroids and the organochlorine DDT target the voltage-gated sodium channel (*Vgsc*), within  
72 which knockdown resistance (*kdr*) mutations occur, which can interfere with insecticide binding,  
73 preventing the normal insecticidal effect of repetitive nerve firing, paralysis and death<sup>14,15</sup>. In the  
74 primary malaria mosquito *Anopheles gambiae* *kdr* mutations exhibit limited variation in flanking  
75 haplotypes<sup>1,16,17</sup>. This conforms to a model of insecticide resistance evolution via hard selective

76 sweeps from rare origins<sup>18</sup>, further supported by observed rapid increases of *kdr* variants toward  
77 fixation over a few years<sup>16,19,20</sup>. Similarly pronounced genomic footprints at the other major  
78 insecticide target site genes, acetylcholinesterase, AChE<sup>13</sup> and GABA<sup>21</sup> suggest that hard sweeps  
79 may be the norm for major resistance mutations in these essential neurotransmission genes. Such  
80 genomic signatures can greatly aid detection of major target site variants in genome scans, though  
81 selection toward fixation progressively reduces minor allele frequency and, as a direct consequence,  
82 statistical power to detect phenotypic association<sup>6,22</sup>.

83

84 With evidence of strong selection of target site mutations, it is perhaps surprising that metabolic  
85 resistance is usually considered the greater threat to vector control<sup>23</sup>. Indeed, the only widely  
86 accepted case of pyrethroid failure in malaria control to date was attributed to a local resurgence of  
87 *An. funestus*<sup>24,25</sup> resistant to insecticide as a result of elevated expression of cytochrome P450  
88 enzymes that metabolize pyrethroids<sup>26,27</sup>. Transcriptomic experiments conducted in *Anopheles spp.*,  
89 have repeatedly linked overexpression of *CYP6* subfamily P450 genes to insecticide resistance and  
90 shown that some genes have the ability to metabolize pyrethroids or even multiple insecticides<sup>28-30</sup>.  
91 In contrast, there are few studies in *Anopheles* which have identified DNA variation associated with  
92 metabolic resistance. The exceptions have targeted a single candidate gene for DDT resistance<sup>31,32</sup>  
93 or a pair of duplicated P450 genes<sup>33</sup>. A paucity of resistance-associated DNA variants might seem to  
94 indicate a relatively greater role for gene overexpression in metabolic insecticide resistance  
95 evolution. Yet no regulatory variants for insecticide resistance-linked differential expression have  
96 been identified in *Anopheles* populations. This implies that allelic variants involved in metabolic  
97 resistance, whether altering proteins or regulatory regions, remain to be discovered.

98

99 Significant methodological challenges are involved in the identification of metabolic SNPs, such as  
100 phenotyping and genotyping sufficient mosquitoes to attain statistical power; population  
101 substructure within samples<sup>22</sup> and typically extremely short linkage disequilibrium in the *An.*  
102 *gambiae* genome<sup>34</sup>, which reduces correlations between functional and marker polymorphisms<sup>35</sup>.  
103 An additional problem for identification of markers for metabolic resistance may result from their  
104 flexibility of function and exceptionally high polymorphism, at least in the largest detoxification gene  
105 superfamilies<sup>36</sup>. Both of these properties contrast markedly with highly-conserved target site genes  
106 subject to purifying selection<sup>13</sup>, and suggest that selection at detoxification genes is more likely to  
107 act upon standing genetic variation and result in soft selective sweeps, making their detection far  
108 more problematic<sup>18</sup>.

109

110 The aim of our study was to apply a candidate gene-enriched single nucleotide polymorphism (SNP)  
111 genotyping array to identify putatively pyrethroid resistance associated polymorphisms segregating  
112 within a wild population of *An. gambiae* from the Uganda-Kenya border area, an important focus of  
113 East African pyrethroid resistance<sup>37-39</sup>, with association established by reproducibility in additional  
114 field populations.

115 The study proceeded via the following steps:

- 116 1. Collection of wild blood-fed females and rearing of progeny as isofemale lines.
- 117 2. Genotyping of adult females individually from the collections using the SNP array to investigate  
118 sources of within population structure
- 119 3. Bioassays and genotyping (using the same SNP array) of the each isofemale lines as a pools and  
120 association analysis based on relationship between family resistance phenotype and SNP genotypes
- 121 4. Identification of putative candidates with additional analysis first within a laboratory colony  
122 established from the same location and subsequently in several distinct wild populations
- 123 5. Sequencing and functional modelling analyses to attempt to identify causal SNPs and interpret the  
124 nature of selection

125 Our results suggest an important role for metabolic gene variants in pyrethroid resistance, which  
126 may be partially diagnosed using the SNPs identified as resistance-associated, and also indicate  
127 association of environmental humidity with pyrethroid resistance.

128

## 129 **Results**

### 130 ***Resistance phenotype***

131 Isofemale families of *An. gambiae* from Tororo, Uganda were bioassayed using WHO permethrin  
132 papers for 60 min for males (N=76 families) and 75 min for females (N=98 families), which provided  
133 a full spectrum of family mortality levels ranging from zero to 100%, with a median of 76% (Figure  
134 S1). Humidity and temperature varied during the testing period in the field insectary (Figure S2), and  
135 across all families, mortality exhibited a significant inverse correlation with humidity ( $r=-0.43$ ,  
136  $P=4\times 10^{-6}$ ; Figure S3), which was subsequently treated as a covariate in association analysis. Within  
137 each family, male and female mortalities were strongly correlated ( $r=0.65$ ;  $P=2\times 10^{-9}$ ), and remained  
138 so following correction for humidity (partial  $r=0.61$ ;  $P=3\times 10^{-6}$ ), suggesting a strong familial basis to

139 resistance. Females were used for subsequent DNA analysis taking the proportionate mortality of  
140 the family from the bioassays (Figure S1) as a quantitative resistance phenotype.

141

#### 142 ***Within population structure***

143 We first examined evidence for within-population structure using individual genotypes from adult  
144 female *An. gambiae*, which comprised a collection of the wild caught mothers of the isofemale lines  
145 and others which did not produce sufficient offspring for use. These females were not classified by  
146 phenotype and were used solely for assessment of possible sources of within population structure  
147 and relatedness, which might impact subsequent association analysis using the pooled isofemale  
148 lines. Clustering analysis suggested three partitions, differentiation between which showed a strong  
149 localisation to the large 2La inversion region, with very limited inter-cluster divergence elsewhere in  
150 the genome (Figure 1). Levels of differentiation within the 2La region were consistent with  
151 comparison among two homokaryotypic clusters and a heterokaryotypic cluster, and linkage  
152 disequilibrium was also extreme (Figure S4). This suggests that 2La inversion polymorphism was the  
153 only major source of structure within the sampled population, and that the sample did not comprise  
154 of groups of closely-related individuals. Though clear that 2La polymorphism is a major source of  
155 genomically-localised subdivision it was unclear how it might influence association analyses because  
156 these females were not phenotyped. Therefore this was evaluated using principal components  
157 analysis as part of the association analysis of isofemale lines.

158

#### 159 *Association analysis using isofemale lines*

160 For pooled genotyping data from the isofemale lines, estimated allele frequencies were highly  
161 correlated with those calculated from individual genotypes after correction for SNP-specific dye bias  
162 ( $R^2=0.98$ ; Figure S5). Principal components analysis was used to identify covariates indicative of  
163 stratification in the pooled data based on 195 control (non-candidate gene) SNPs. The first principal  
164 component (PC1) was significantly correlated with resistance phenotype (partial correlation  
165 controlling for humidity:  $r=-0.33$ ,  $P=0.001$ ; Table S1). PC1 primarily reflected variation within the 2La  
166 inversion region. All control (non-candidate) SNPs correlating strongly with PC1 ( $r \geq 0.5$ ,  $N=13$ ) were  
167 located therein, and the mean correlation for 2La SNPs ( $r=0.41$ ,  $N=32$ , was significantly higher  
168 (Mann-Whitney U-test,  $Z=5.30$ ,  $P<0.001$ ) than for control SNPs distributed throughout the rest of the  
169 genome ( $r=0.11$ ,  $N=163$ ). A second principal component was also correlated with mortality ( $r=0.22$ ,  
170  $P=0.03$ ; Table S1), but showed no localisation of SNPs within the genome, and no difference in  
171 correlation between SNPs in 2La and elsewhere ( $Z=1.45$ ,  $P=0.15$ ). Owing to the concordance with  
172 individual genotyping results (above), in which the 2La region appeared to be the major source of

173 stratification within the pooled dataset, we opted to consider PC1 further in order to apply  
174 correction for 2La variation. Elevated probabilities within the 2La region were evident when plotting  
175 all SNPs, but were removed by correction for PC1 as a covariate (Figure S6A), whereas if SNPs from  
176 the 2La region were removed *a priori*, correction for PC1 was not required to align the vast majority  
177 of SNPs with their expectation (Figure S6B). However, because correction for 2La variation (via PC1)  
178 might obscure SNPs which contributed to, or acted in addition to, the resistance association of 2La,  
179 and also because correction is likely to be imperfect, we considered results both with and without  
180 correction for PC1.

181

182 For the association analysis the relatively limited number of families analysed limited study power,  
183 with calculations suggested only moderate-large effects might be detected (Materials and Methods).  
184 Consequently, we did not apply formal thresholds for significance in the analysis, but rather  
185 identified candidate SNPs for further investigation from those with high ranking  $-\log P$  values, with  
186 demonstration of association based on subsequent replication of results. Most of the SNPs with  
187 lowest P-values were on chromosome 2L, of which the two showing strongest evidence for  
188 association were a synonymous variant in *Cyp4j5* and a non-synonymous variant in *Cyp4j10* located  
189 2kb apart (Figure 2A; Table S2). Systematic elevation of P-values was clear within the 2La inversion,  
190 especially toward the ends, and following correction for PC1, P-values of most SNPs located within  
191 2La reduced substantially. SNPs with the lowest P-values post-correction were from chromosome 2L  
192 but outside the 2La region; within the *Vgsc* target site gene and a carboxylesterase gene, *Coae1d*  
193 (Figure 2B).

194

195 The candidate P450 and carboxylesterase gene SNPs were at high minor allele frequency (MAF, 0.35-  
196 0.49). However, those in the *Vgsc* exhibited much lower MAF ( $\leq 0.12$ ), possibly as a result of a  
197 selective sweep of the *kdr 1014S* mutation, which is close to fixation (frequency=0.94). To address  
198 the issue of potentially poor detectability of association of low MAF SNPs that might be influenced  
199 by strong selection in the wild, we performed an additional (individual-mosquito) association study  
200 using the class II pyrethroid deltamethrin. This experiment used bioassay survivors and dead from a  
201 recently-founded colony from the same location, which had not been previously been exposed to  
202 insecticide (Figure S7A). In the colony, MAFs of *Vgsc*-L1014S (MAF=0.37) and other *Vgsc* SNPs were  
203 much higher than in the wild population (Table S3) and L1014S exhibited the strongest resistance  
204 association of all the SNPs typed (Figure S7B). The *Vgsc* SNP exhibiting the lowest probability in the  
205 previous analysis (Figure 2B) and the SNP in *Coae1d* were also strongly associated with resistance  
206 phenotype in the colony.

207

### 208 **Marker association repeatability**

209 We chose four SNPs for further testing, the linked *Cyp4j5* and *Cyp4j10* variants, the *Coeae1d* SNP  
210 and the *Vgsc*-L1014S variant. When tested sequentially using stepwise regression, insectary humidity  
211 explained the highest proportion of the variance in permethrin resistance phenotype among the  
212 families, but each of the four SNPs explained significant additional variance, together totalling 22%  
213 among isofemale lines, with *Cyp4j5* the strongest predictor (Table 1). Novel TaqMan assays were  
214 designed for the *Cyp4j10*, *Coeae1d* and *Cyp4j5* SNPs. For the latter we used a non-synonymous SNP  
215 (not present on the array) in highest LD with the original marker, a leucine-phenylalanine change at  
216 codon 43 of the gene (L43F) (see Materials and Methods for details of source). These assays were  
217 combined with that for *Vgsc*-L1014F for screening in up to seven independent field samples  
218 phenotyped for susceptibility to three pyrethroid insecticides (Figure 3). The known causal marker  
219 *Vgsc*-L1014S was significantly resistance-associated in two of the samples and the *Coeae1d* marker  
220 in three. The *Cyp4j10* SNP failed to approach significance in three test samples and screening was  
221 discontinued, but the *Cyp4j5*-L43F marker was highly reproducible, with significant resistance  
222 association in five out of seven independent sample sets from both Uganda and Kenya (Figure 3).  
223 Average effect sizes and their variability are shown in Table 2 and Table S4. For *Vgsc*-L1014S, odds  
224 ratios were highly variable, ranging from approximately 1 to over 10; were more consistent for  
225 *Cyp4j5* (range 1.5-4.7) and lower but again quite consistent (range 0.7-2.7) for *Coeae1d* (Table S4).  
226 For *Cyp4j5* we also investigated the predictive utility for a more extreme phenotype by comparing  
227 deltamethrin survivors of a 2 hour exposure with those killed by a one hour exposure; interestingly  
228 this yielded the highest odds ratio observed for this marker (OR=6.9, P=0.0001, Table S4).

229 Temporal analysis of variation in marker frequencies in Tororo in 2013-2014 revealed relative  
230 stability of frequencies of each of the three metabolic gene markers (Figure S8). Interestingly,  
231 genotype frequencies of the *Cyp4j5* and *Cyp4j10* markers were strongly correlated (Table S5) as  
232 expected from their very close proximity in the genome, suggesting that the initial association of the  
233 *Cyp4j10* marker for pyrethroid resistance may have in part arisen from co-variation with *Cyp4j5*.  
234 Neither *Cyp4j5* nor *Coeae1d* marker genotype variation was correlated with 2La variation, despite  
235 apparent allele frequency covariation (Figure S5) highlighting their independent predictive utility  
236 (Table S5).

237

### 238 **Investigation of possible causal SNPs**

239 The resistance-associated variant at base 20288132 (VectorBase nucleotide numbering) in the  
240 *Coeae1d* gene is a synonymous variant. Therefore, to investigate whether any non-synonymous

241 SNPs showed evidence of association which might be underpinning this, we sequenced a single  
242 individual from each of 24 families from the extremes of the mortality distribution (see [Figure S1A](#);  
243 [Table S6A](#)). The *Coeae1d* gene is extremely polymorphic with a large number of missense variants  
244 (N=116), many of which are at high frequency (31% with MAF>0.1; 19% with MAF>0.2). In this  
245 methodologically-independent replication of the initial association analysis, the 20288132 marker  
246 SNP exhibited a significant level of association (OR=2.9,  $\chi^2= 5.20$ , P=0.022) but this was not exceeded  
247 by any function-altering polymorphism within the gene, and only marginally so by two other  
248 synonymous SNPs ([Figure S9](#); [Table S6B](#)). We also genotyped the 2La inversion in the same samples,  
249 which showed significant association of the derived (2L<sup>a</sup>) karyotype with permethrin resistance  
250 (OR=4.0,  $\chi^2= 4.85$ , P=0.028).

251

252 Sequencing of the *Cyp4j5* gene in Ugandan individuals revealed a split between resistant and  
253 susceptible haplotypes around the L43F variant, including similarity of resistant alleles from distinct  
254 locations (Jinja and Oyam) in southern and northern Uganda ([Figure 4](#)). However, no marked  
255 difference in diversity, as would be expected with a strong selective sweep, was evident. We used  
256 structural modelling to evaluate the possible functional consequences of the L43F, and other non-  
257 synonymous substitutions in *CYP4J5*. Of the variants ([Table S7](#)) that could be modelled (noting that  
258 available templates bore only approximately 20% sequence identity to CYP4J5, limiting model  
259 quality), the substitution at L43F was not modelled reliably but seems unlikely to have an impact on  
260 protein stability. Of the other SNPs only one, a rare Alanine to Threonine substitution at codon 69,  
261 was predicted to potentially fall near enough to the active site for impact. However, this SNP was  
262 not detected as being in LD with the original marker, nor with the L43F variant; in fact the  
263 correlation with each was negative suggesting, if anything, weak linkage of the resistance associated  
264 allele at each marker with the wild type allele at A69T. Strong repeatability of the L43F mutation as a  
265 resistance predictor (above) suggests that either this is in strong LD with an, as yet unidentified,  
266 causal variant, or contributes directly to the resistance phenotype in a way that was not revealed by  
267 the current modelling approach.

## 268 **Discussion**

269 Using a candidate-GWAS approach we have identified two novel SNPs in genes from major  
270 metabolic superfamilies that are reliably and reproducibly associated with resistance to class I and II  
271 pyrethroids. We also detected a strong relationship between bioassay survival and humidity in our  
272 initial experiment and evidence for association of 2La inversion polymorphism with resistance.  
273 Moreover, although the genes containing the associated SNP are located within and close to,

274 respectively, independent temporal dynamics of the SNP and 2La polymorphisms in Ugandan field  
275 samples suggests that additional variants resistance-associated variants remain to be discovered  
276 within 2La. The SNP markers identified are unlikely to be causal substitutions and we speculate that  
277 linkage disequilibrium between each marker SNP and one or more haplotypes, each likely to contain  
278 multiple non-synonymous changes in a resistance-conferring haplotype, may underpin the  
279 associations observed. Certainly the pattern of polymorphism deviates from that arising from  
280 selection on a single, relatively recent mutation<sup>13,16,17</sup> and, particularly if haplotypes are driving the  
281 resistance association, selection from standing variation appears more likely<sup>40</sup>.

282 In the Eastern Ugandan Tororo field population strong association of the *Vgsc*-1014S (*kdr*) variant  
283 with both permethrin and DDT resistance has been demonstrated previously<sup>37</sup> and we therefore  
284 suspected that the relatively weak association with permethrin-survival across families in our study  
285 might be attributed to limited statistical power arising from low polymorphism (1014S frequency  $\approx$   
286 95%). Results from the Tororo colony within which the frequency of the wild type 1014L allele was  
287 much higher, supported this hypothesis: L1014S was the SNP most strongly associated with survival  
288 following deltamethrin exposure. This is significant because: (i) the association of *Vgsc*-1014S has  
289 not previously been tested against a background of many other candidates; (ii) a significant link with  
290 class II resistance has proved difficult to establish<sup>37,41</sup>; (iii) it suggests that 1014S, and probably other  
291 *kdr* mutations continue to play a pivotal role in resistance in the presence of other resistance-linked  
292 mutations. Thus, it is important to recognise that even following rise toward fixation where *kdr*  
293 mutations explain relatively little of the variation in resistance *within* a population, they provide a  
294 strong baseline of resistance, which can be elevated further – rather than being superseded - by the  
295 presence of other resistance-linked variants<sup>6</sup>. An additional *Vgsc* mutation (N1575Y) detected in *A.*  
296 *gambiae* and *A. coluzzii* across West Africa<sup>17</sup> illustrates this point. The 1575Y allele occurs only on a  
297 *Vgsc*-1014F background and though the SNP itself exhibits a moderate odds ratio ( $\approx 2$ ), this  
298 effectively doubles that of *kdr1014F*, and so translates into a large aggregate odds ratio ( $\approx 12$ )<sup>7,17</sup>.

299 Polymorphism of the 2La inversion has long been linked with important phenotypic variation in *An.*  
300 *gambiae*, primarily with adaptation to environmental aridity<sup>42,43</sup>, but also with refractoriness to  
301 *Plasmodium falciparum*<sup>44</sup>. Our data suggest that insecticide resistance can also be associated with  
302 the 2La inversion, with the major permethrin resistance-associated principal component (PC1)  
303 identified in our family pool study, correlating strongly with variation at SNPs within 2La. Owing to  
304 the known relationship between the 2La inversion and aridity, it is important to note that in the  
305 Tororo discovery dataset the association of PC1 remained significant following correction for  
306 humidity, i.e. both humidity and 2La polymorphism were determinants of resistance phenotype. In

307 addition, when 2La polymorphism was typed directly, during our investigation of *Coeae1d* variation,  
308 the 2L+<sup>a</sup> karyotype was significantly associated with permethrin survival. Genomic stratification in  
309 the Tororo population was attributable to 2La polymorphism, therefore we opted to investigate  
310 further the most associated SNPs that emerged prior to, and post-correction using PC1 as a proxy for  
311 2La variation. The resultant analyses showed that identification of top candidate SNPs solely post-  
312 correction for 2La would have been a mistake, because the strongest and most reproducible  
313 candidate SNP (*Cyp4j5*-L43F) would probably have been missed at the discovery phase. This is  
314 noteworthy and reflects a difference in genomewide stratification as expected in an admixed  
315 population, which would necessitate *a priori* correction to avoid false positives, and genomically-  
316 localised stratification, as seen here, which is a driver of the phenotype and unrelated to mating  
317 barriers.

318 Interestingly, when typed in a series of field samples collected from near Tororo between 2013  
319 through 2014, neither *Cyp4j5* (which was correlated with PC1) nor *Coeae1d* (uncorrelated with PC1)  
320 genotype frequencies were correlated within 2La karyotype variation, despite visual resemblance of  
321 summary data (see [Figure S4](#)). This is an important finding because it means that these two SNPs and  
322 2La can serve as independent diagnostic markers, and perhaps most crucially because evolutionary  
323 dynamics of *Coeae1d* and especially *Cyp4j5* (located within 2La) need not be constrained by  
324 dependence on 2La inversion polymorphism dynamics, which might be subject to multiple and  
325 potentially contrasting selection pressures. Whether such constraints are currently operative within  
326 the East African area we covered in our study is unclear, but recent indirect evidence suggests that  
327 insecticidal pressure might be the dominant selective factor for 2La variation. In *A. gambiae* s.s.  
328 samples from a 17-year time series collected in Western Kenya, Matoke-Muhia et al.<sup>45</sup> detected a  
329 sharp decline in frequency of the previously numerically-dominant 2La karyotype in favour of the  
330 2L+<sup>a</sup> karyotype, which correlated strongly ( $r=-0.96$ ) with the changing coverage of pyrethroid-treated  
331 bednets in the area. To our knowledge our study is the first to provide direct association data  
332 directly linking pyrethroid resistance with 2La inversion polymorphism, adding corroborative  
333 evidence to the field survey results from Western Kenya.

334 An important objective is now to identify the genes/ variants in 2La underpinning pyrethroid  
335 resistance association (noting that *Cyp4j5* apparently is not the major cause). Contrasting with most  
336 of the *An. gambiae* genome<sup>22,34,46</sup>, LD within the 2La inversion region can be extreme, especially for  
337 a few megabases within the breakpoints around 22 and 42Mb<sup>22</sup>. Such strong LD creates the  
338 problem, commonly encountered in GWAS of human populations outside of Africa<sup>35</sup>, of how to  
339 isolate markers from large statistically-elevated haplotype blocks. Increasing sample size may be one

340 answer, but will prove difficult in the near-term, because low LD throughout the majority of the  
341 genome argues strongly for the use of high density genotyping-by-sequencing approaches for which  
342 highly powered studies remain expensive. Instead, we suggest that investment of greater resources  
343 into *post hoc* investigation of repeatability of phenotype associations for smaller, filtered, sets of  
344 markers in wild populations could permit the signals from multiple linked markers to be teased  
345 apart. Polymorphic inversions are common in insects and low genomic LD in wild populations likely  
346 to be also, therefore such study design issues are of relevance to many taxa.

347 Our study demonstrates that pyrethroid resistance in eastern Uganda/western Kenya populations  
348 has a multivariate basis. In addition to the genetic polymorphism associations in the original  
349 discovery phase of the project, our data show a strong impact of humidity during bioassays on  
350 resistance, as reported previously in houseflies *Musca domestica*<sup>47</sup>. Although humidity varied  
351 naturally, rather than experimentally, in our study it explained over 17% of the variation among  
352 families in resistance and is very likely to be an important ecological determinant of insecticide  
353 resistance phenotypes under natural conditions. Of the genetic component of variation, both target  
354 site and metabolic variants play a role, with the two metabolic markers identified explaining a  
355 significant proportion (19%) of within-population variance. In a microarray study of the same Tororo  
356 population sampled the following year, relatively few significantly overexpressed genes were  
357 detected and only two detoxification genes<sup>48</sup>, though the design did not include a fully susceptible  
358 sample which will tend to accentuate discovery of overexpressed genes<sup>7</sup>. Nevertheless, available  
359 evidence suggests that the variants we identified are likely to link with mutations causing qualitative  
360 rather than quantitative protein variation, and we suggest that in Tororo this may be a more  
361 important determinant of pyrethroid resistant phenotypes. Owing to the paucity of agnostic DNA-  
362 based studies, which contrasts with the plethora of insecticide resistance-targeted gene expression  
363 experiments, whether this is population-specific, or the case more generally, remains to be  
364 determined.

### 365 **Conclusion**

366 With increased selective pressure arising from massive scale-up of long lasting insecticidal net (LLIN)  
367 distributions over the last decade, levels of pyrethroid resistance are likely to increase further and  
368 threat to control appears imminent. Identification of the genetic mechanisms leading to heightened  
369 resistance, and provision of simple assays for rapid diagnosis will aid attempts to manage resistance  
370 and retain pyrethroids as a viable option. Our study has highlighted the importance of allelic  
371 variants, outside the insecticide target site, as determinants of a significant fraction of variation in  
372 the resistance phenotype in a wild *An. gambiae* populations. Whole genome resequencing should

373 permit identification of a spectrum of DNA polymorphisms underpinning resistance, and facilitate  
374 development and widespread application of DNA diagnostics as an accurate proxy for *a priori*  
375 assessment of phenotypic resistance.

## 376 **Materials and Methods**

### 377 **Sample collections**

378 Indoor-resting adult females were collected by aspiration from houses in Tororo, Uganda  
379 (00°40'41.6"N, 34°11'11.6"E) in November 2008 and transported to insectaries at the nearby  
380 National Livestock Resources Research Institute (NaLiRRI) facility. Morphological examination of  
381 almost 900 adult females revealed that the vast majority of samples ( $\approx$ 98%) were *Anopheles*  
382 *gambiae s.l.*, with the remainder *An. funestus*, which were discarded. Adult female *An. gambiae s.l.*  
383 were held individually in plastic cups covered with fine netting, and contained moist cotton wool at  
384 the base covered by a disc of filter paper onto which eggs were laid. A cotton wool pad on top of  
385 each cup provided 10% sucrose solution *ad libitum*. Once eggs were laid, each filter paper disc was  
386 removed to a plastic bowl containing water, the adult female was sacrificed and preserved over silica  
387 gel. Following hatching, larvae were fed twice-daily with finely-ground Tetramin fish food. Pupae  
388 were transferred to netted cups containing a small volume of water, with one cup per isofemale  
389 family. On eclosion females and males were separated daily and isofemale line-specific bioassays  
390 were performed separately on 10-20 (median N=15) 3-5 day old adult females (or males) in WHO  
391 tubes using standard 0.75% permethrin papers. Informed by bioassay data from preliminary tests on  
392 mixed F<sub>1</sub> offspring from multiple families, we used an insecticide exposure time of 75 min (and 60  
393 min for males), which aimed to produce an intermediate level of mortality (ca. 50%) for the  
394 population as a whole. Humidity and temperature were recorded during each assay. Approximately  
395 24 h after exposure, all individuals were counted as alive or dead and preserved individually over  
396 silica gel.

397

398 Establishment and phenotype-classification of a colony for secondary association analysis  
399 Eggs were initially transported to the Liverpool School of Tropical Medicine from the Tororo  
400 collection site in December 2008 to found a colony; with a second supplementary collection added  
401 in May 2009. The colony was maintained without any insecticide exposure until September 2009 at  
402 which point 3-5 day old adult females were exposed to varying concentrations of the pyrethroid  
403 insecticide deltamethrin (Greyhound chemicals, UK) using WHO bioassay tubes and custom papers  
404 produced using the appropriate quantity of deltamethrin dissolved in acetone, added to silicone oil  
405 and spread over Whatman papers. Following the usual 24 h holding period, 20 adult females that

406 died at low deltamethrin concentrations (<0.01%) and 36 surviving high deltamethrin concentrations  
407 ( $\geq 0.05\%$ ; note 0.05% is the standard WHO resistance diagnostic concentration) were defined as  
408 susceptible and resistant classes, respectively and preserved individually over silica gel.

409

#### 410 **SNP array genotyping**

411 DNA from *An. gambiae s.l.* individuals was extracted, quantified and typed to species using  
412 previously-described methods<sup>49</sup>. More than 96% of samples were *An. gambiae* (formerly *An.*  
413 *gambiae s.s. S-form*); the remainder were *An. arabiensis*, which were not considered further. The  
414 Illumina Goldengate 1536-SNP array used is described in detail elsewhere<sup>22,49</sup> but briefly was  
415 enriched for SNPs at 266 candidate detoxification or target site genes, with approximately 20% of  
416 SNPs, located in intergenic regions or non-candidate genes. Although the SNP-array was populated  
417 by candidate genes, it should be noted that none of the SNPs present were known candidates in this  
418 geographical area. The genotype of the known resistance mutation, *Vgsc (kdr) L1014S* was also  
419 determined by direct sequencing of PCR products amplified using primers ex20+ F/R<sup>16</sup>. Mothers of  
420 isofemale lines and additional females collected but not producing eggs were genotyped  
421 individually, as were females from the laboratory colony. Female offspring from families were  
422 genotyped as single-family pools, using an equimolar amount of DNA from each individual to provide  
423 a final concentration per family pool of 50 ng/ $\mu$ l.

424

#### 425 **SNP array data analysis**

426 *Individuals*: Genotyping arrays for 216 adult females (N=98 mothers of the isofemale lines and  
427 N=118 other females from house collections) and N=56 Tororo colony specimens were scored using  
428 Beadstudio v3.2 (Illumina Inc). Haploview 4.1<sup>50</sup> was used to compute minor allele frequencies,  
429 linkage disequilibrium among SNPs as  $r^2$ , observed and expected heterozygosities, and to test for  
430 deviation from Hardy-Weinberg expectations. Individual cluster analysis of genotypes at the control  
431 (intergenic and non-candidate) SNP loci was performed by BAPS 5.2<sup>51,52</sup>, with multiple runs  
432 performed to obtain the optimum clustering solution.

433

434 *Pools*: Alleles labelled with different dyes seldom give a perfectly equal signal in heterozygotes and  
435 such dye bias may be assay-specific. Dye bias is typically not a major concern for individual  
436 genotyping but can greatly impact accurate estimation of allele frequencies from pooled DNA<sup>53</sup>.  
437 Therefore prior to analysis of pooled genotypes we computed the dye bias ( $k$ ) for each SNP for which  
438 a minimum of two individual heterozygotes were available, as the mean cy3: cy5 dye ratio. For  
439 example, if heterozygotes exhibit a mean ratio of 0.5 (i.e. a twofold cy5 bias) then  $k=0.5$ . We

440 excluded SNPs exhibiting extreme or poorly estimated values of  $k$ , which can give rise to poor  
441 estimates of allele frequencies in pooled analysis<sup>53</sup>, based on two exclusion criteria: (1)  $k >$  upper  
442 95<sup>th</sup> percentile of the range of  $k$  values; (2) standard deviation of  $k$  across individuals  $>$  upper 95<sup>th</sup>  
443 percentile of the range of  $k$  value standard deviations. Of the 1536 SNPs on the array 895 could be  
444 scored reliably, were polymorphic and exhibited  $k$  values meeting criteria 1 and 2: only these SNPs  
445 were used in the analyses. Raw X and Y (cy3 and cy5) signal data were extracted from Beadstudio  
446 v3.2 and used as the basis for all analyses of pooled DNA. Allele frequencies,  $p(X)$  and  $p(Y)$ , were  
447 computed as:

448

449 1.  $p(X) = \text{raw X (cy3) signal} / [\text{raw X (cy3) signal} + k \cdot \text{raw Y (cy5) signal}]$

450 2.  $p(Y) = 1 - p(X)$

451

452 Principal components analysis (PCA) using allele frequency data from 195 control SNPs was used to  
453 investigate possible within-sample stratification using the pooled genotype data. Association  
454 analysis was performed on rank-transformed data by computing correlations between allele  
455 frequency and proportionate bioassay mortality for each family, or partial correlations to correct for  
456 possible covariates (temperature, humidity and principal components; see Results).

457 Power calculation suggested that our analysis could detect a slope for mortality vs allele frequency  
458 in the analysis of 0.46 (as significantly different from zero), with 80% power and a Bonferroni-  
459 corrected alpha level ( $0.05/N$  SNPs), with the 894 SNPs successfully genotyped and our dataset of  
460 pooled genotypes from 98 families. This suggests that moderate-large effect sizes associated with  
461 any individual SNP would be detectable. Whilst this was compatible with our aim to discover  
462 markers with useful predictive power for phenotypes (i.e. those of larger effect) rather than a  
463 comprehensive set of variants, owing to the limited power we do not apply formal significance  
464 thresholds in the analysis but rather base assessment of association on reproducibility of candidates  
465 selected for further testing.

466 IBM SPSS 22 was used to perform subsequent stepwise regression analyses for the pooled genotype  
467 data and multiple logistic regression for the individual data from the reproducibility testing  
468 populations.

#### 469 **Associated gene sequencing**

470 To attempt to identify the possible causal SNP and evaluate evidence of selection the *Cyp4j5* gene  
471 was sequenced in several individuals from Jinja and Oyam, distinct locations from southern and  
472 northern Uganda, respectively (locations given in subsection below). Specimens were genotyped

473 using the *CYP4J5* marker (2L: 25635973) and six homozygotes for each allele *Cyp4j5*-43L and *Cyp4j5*-  
474 43F, the latter resistance associated, were selected for cloning and sequencing. The entire gene was  
475 amplified using primers 4J5F1 5'-AGCACACGGTAAGGATGTC-3' and 4J5R1 5'-  
476 GCGGAGAAACGTAACCCATA-3' with a Phusion PCR kit (Thermo Scientific) prior to cloning in *E.coli*  
477 using a pJET1.2/blunt cloning vector. Due to the size of the product (~3.5 kb), two internal primers  
478 (4J5SEQ2 5'-ATTGCCGACTGTAGCTCGAT-3' and CYP4J5-2 5'-GGCTTCTTTGGGACACACAT-3') were used  
479 in addition to the pJET plasmid specific primers (pJETF and pJETR). Clones were sequenced by  
480 Macrogen, Korea and edited using CodonCode Aligner v4 (CodonCode Corporation). Gene trees  
481 were produced in MEGA 7<sup>54</sup> using the neighbor-joining method with 1000 bootstraps. Evolutionary  
482 distances were computed using the Maximum Composite Likelihood method using all codon  
483 positions. All positions containing gaps and missing data were eliminated yielding a total of 1596  
484 positions in the final dataset.

485

486 The *Coeae1d* gene was sequenced in a single individual from 24 families from each end of the  
487 spectrum of mortalities, and each individual thus designated as resistant or susceptible. Initial PCR  
488 amplification of a 2047bp fragment used primers Agap5756his (5'ATCGTCAACGTGCTCAGTCA3') and  
489 Agap5756exp (5'AGCACAGCGACTAACTCTTGC3') in the following mixture 5µl KapaTaq 10X buffer, 1µl  
490 (10mM dNTPs), 10 picomoles of each primer, 1.5U KapaTaq polymerase, 40.7 µl water & 1µl DNA  
491 template. Cycling conditions were 95°C for 5 min, followed by 40 cycles of 95°C for 30sec, 57°C for  
492 30 sec, 72°C for 90 sec, with a final step of 72°C for 10 min. Products were cleaned with QIAquick  
493 PCR Purification Kit (Qiagen) and sent for sequencing by Macrogen, Korea using the following  
494 overlapping sequencing primers covering 1951 bp within the amplicon, which included the whole  
495 gene (5756SEQ7i 5'AACTCCTCGAAACCTCACCTA3'; 5756SEQ1 5'ATGCAACCCTTCTAGCG3';  
496 5756SEQ2i 5'ACCCATTGCTCAAACA3'; 5756SEQ3i 5'GGTTTCGAGGAGTTTTGTTGT3'; 5756SEQ5  
497 5'TCACTCCGGCAGCTTCTTA3'; 5756SEQ4i 5'CAACTCAACTCTTTACAGACATGC3'). Sequences were  
498 aligned to a reference sequence from VectorBase in CodonCode Aligner v4, edited and  
499 polymorphisms identified. Chi-square tests were conducted for each SNP identified to investigate  
500 association with permethrin resistance phenotype, by comparing allele frequencies between  
501 resistant and susceptible individuals. We also genotyped polymorphism of the 2La inversion in these  
502 samples using a PCR-based method<sup>55</sup> for association analysis based on chi-square tests as before.

### 503 **Functional modelling**

504 Templates for modelling of *Anopheles gambiae* CYP4J5 were ranked by a search of the Protein Data  
505 Bank<sup>56</sup> using HHsearch<sup>57</sup>. The best scoring 10 structures were used for comparative modelling using

506 the recent high-resolution methods implemented in Rosetta<sup>58</sup> to generate 30 models. The top  
507 model was used for interpreting the likely consequences of observed SNPs using structure-based  
508 methods FoldX<sup>59,60</sup> and PoPMuSiC<sup>61</sup>. This analysis was confined to regions that comparison  
509 between models indicated were predicted reliably. SNPs positioned in regions not considered  
510 reliably modelled were interpreted using PolyPhen2<sup>62</sup>. The structure of deltamethrin, obtained  
511 from the Toxin and Toxin Target Database<sup>63</sup>, was manually positioned using PyMOL  
512 ([www.pymol.org](http://www.pymol.org)) in the CYP4J5 model in a similar position to ritonavir bound to CYP3A4<sup>64</sup>, one of  
513 the template structures. Its position relative to amino-acids altered by SNPs was assessed. Since the  
514 available templates bore only low (~20%) sequence identity with *Anopheles gambiae* CYP4J5 and  
515 apo-structures were among the templates, the binding cavity of the CYP4J5 model was not predicted  
516 accurately enough to allow for automated docking. Preliminary functional modelling of COEAE1D  
517 was performed but owing to a very large number of non-synonymous variants (see Results) was  
518 discontinued.

#### 519 **Replication study samples and assays**

520

521 We designed novel TaqMan assays for two candidate CYP450 markers identified in the family pooled  
522 genotyping and colony individual genotyping experiments and the *Coeae1d* marker. For the *Cyp4j5*  
523 marker, non-target polymorphism in the assay region prevented design for the associated SNP,  
524 therefore we sought to design an assay for a non-synonymous SNP within the gene in highest linkage  
525 disequilibrium (LD) with the marker. Data were obtained from whole genome sequences of 40  
526 isofamily pools of females, collected in the same way as before from the same locations one year  
527 later<sup>48</sup>. LD was approximated by computing family-wise Pearson correlations for each of the variants  
528 identified with the original marker (Table S7). The SNP with the highest correlation coefficient (and  
529 similar minor allele frequency to the original marker) was a leucine to phenylalanine substitution at  
530 codon 43, for which an assay was successfully designed.

531

532 For independent replication of marker: phenotype associations, samples bioassayed for pyrethroid  
533 susceptibility were obtained from a number of spatially and/or temporally distinct collection  
534 locations in the East African region. Oyam (02°14' N 32°23' E), Apac (01°59' N 32°32'E) Jinja<sup>65</sup>,  
535 (00°46'N 34°01'E) Nagongera<sup>66</sup> all in Uganda and the Busia/Bungoma area of Western Kenya<sup>38</sup>. All  
536 were collected as larvae with rearing as above and bioassays-screened using standard WHO  
537 methods. Samples were genotyped for SNPs using a TaqMan SNP genotyping assay for *Vgsc* L1014S  
538<sup>67</sup>, and novel TaqMan assays for the: (1) *CYP4J10* marker (2L: 25,636,722) primers were (forward 5'-  
539 ATCACGGTGCAGATCGT-3') and (reverse 5'- GCTTCAAGAATCTGAATCGC-3') and SNP specific probes

540 (G: 5'-6FAM- CATGTCACACGTCCACA-3' and A: 5'-VIC- CATGTCACACATCCACA-3'); (2) *COEAE1D*  
541 marker (2L: 20,288,132) primers were (forward 5'- GAGAGTGCAGGAGCTAAGGC-3') and (reverse 5'-  
542 CTCCACTTTGACAGATCACTCGAT-3') and SNP specific probes (G: 5'-6FAM- CCTATCTGCATTACCTTT-3'  
543 and A: 5'-VIC- CCTATCTGCACTACCTTT-3'); (3) *CYP4J5* marker (2L: 25635973) primers were (forward  
544 5'- AGCCTGCGCGTGTGATA-3') and (reverse 5'- CTTCTTCTCTGTGGTTCGTTG-3') and SNP specific  
545 probes (G: 5'-6FAM- TTGCCGGAAGGCAGT-3' and A: 5'-VIC- TTGCCGGAAGGCAGT-3'). Probes carried  
546 non-fluorescent quenchers at the 3' end. Assays were performed in a 10µl volume containing 1x  
547 Sensimix (Bioline), 1x primer/probe mix and 1µl template DNA with a temperature profile of 95°C for  
548 10 min followed by 40 cycles of 92°C for 15s and 60°C for 1 min on an Agilent MX3005 real-time PCR  
549 machine. VIC and FAM fluorescence was captured at the end of each cycle and genotypes called  
550 from endpoint fluorescence using Agilent MXPro software. Genotype-phenotype associations were  
551 examined using  $\chi^2$  tests and odds ratios.

552

553

## 554 References

- 555 1 Clarkson, C. S. *et al.* Adaptive introgression between Anopheles sibling species eliminates a  
556 major genomic island but not reproductive isolation. *Nature Communications* **5**, e4248,  
557 doi:10.1038/ncomms5248 (2014).
- 558 2 Bhatt, S. *et al.* The effect of malaria control on Plasmodium falciparum in Africa between  
559 2000 and 2015. *Nature* **526**, 207-+, doi:10.1038/nature15535 (2015).
- 560 3 Ranson, H. & Lissenden, N. Insecticide resistance in African Anopheles mosquitoes: a  
561 worsening situation that needs urgent action to maintain malaria control. . *Trends Parasitol*  
562 (2016).
- 563 4 Kleinschmidt, I. *et al.* Design of a study to determine the impact of insecticide resistance on  
564 malaria vector control: a multi-country investigation. *Malar. J.* **14**, 13, doi:10.1186/s12936-  
565 015-0782-4 (2015).
- 566 5 WHO-GMP. *Global plan for insecticide resistance management in malaria vectors (GPIRM)*.  
567 (World Health Organization, 2012).
- 568 6 Weetman, D. & Donnelly, M. J. Evolution of insecticide resistance diagnostics in malaria  
569 vectors. *Transactions of The Royal Society of Tropical Medicine and Hygiene* **109**, 291-293,  
570 doi:10.1093/trstmh/trv017 (2015).
- 571 7 Donnelly, M. J., Isaacs, A. & Weetman, D. Identification, validation, and application of  
572 molecular diagnostics for insecticide resistance in malaria vectors. *Trends in Parasitology* **32**,  
573 197-206, doi:doi:10.1016/j.pt.2015.12.001 (2016).
- 574 8 Donnelly, M. J., Isaacs, A. & Weetman, D. Identification, validation, and application of  
575 molecular diagnostics for insecticide resistance in malaria vectors. *Trends Parasitol*,  
576 doi:doi:10.1016/j.pt.2015.12.001 (2016).
- 577 9 Hemingway, J., Hawkes, N. J., McCarroll, L. & Ranson, H. I. The molecular basis of insecticide  
578 resistance in mosquitoes. *Insect Biochemistry and Molecular Biology* **34**, 653-665 (2004).
- 579 10 Weill, M. *et al.* Insecticide resistance: a silent base prediction. *Current Biology* **14**, R552-R553  
580 (2004).
- 581 11 Davies, T. G. E., Field, L. M., Usherwood, P. N. R. & Williamson, M. S. A comparative study of  
582 voltage-gated sodium channels in the Insecta: implications for pyrethroid resistance in  
583 Anopheline and other Neopteran species. *Insect Molecular Biology* **16**, 361-375 (2007).

- 584 12 Remnant, E. J. *et al.* Gene duplication in the major insecticide target site, Rdl, in *Drosophila*  
585 *melanogaster*. *Proceedings of the National Academy of Sciences of the United States of*  
586 *America* **110**, 14705-14710, doi:10.1073/pnas.1311341110 (2013).
- 587 13 Weetman, D. *et al.* Contemporary evolution of resistance at the major insecticide target site  
588 gene Ace-1 by mutation and copy number variation in the malaria mosquito *Anopheles*  
589 *gambiae*. *Molecular ecology* **24**, 2656-2672, doi:10.1111/mec.13197 (2015).
- 590 14 Ranson, H. *et al.* Identification of a point mutation in the voltage-gated sodium channel gene  
591 of Kenyan *Anopheles gambiae* associated with resistance to DDT and pyrethroids. *Insect*  
592 *Molecular Biology* **9**, 491-497 (2000).
- 593 15 Martinez-Torres, D. *et al.* Molecular characterization of pyrethroid knockdown resistance  
594 (*kdr*) in the major malaria vector *Anopheles gambiae s.s.* *Insect Molecular Biology* **7**, 179-184  
595 (1998).
- 596 16 Lynd, A., Weetman, D., Barbosa, S., Yawson, A.E., Mitchell, S., Pinto, J., Hastings, I. and  
597 Donnelly, M.J. Field, genetic and modelling approaches show strong positive selection acting  
598 upon an insecticide resistance mutation in *Anopheles gambiae s.s.* *Molecular Biology and*  
599 *Evolution* **27**, 1117-1125 (2010).
- 600 17 Jones, C. *et al.* Footprints of positive selection associated with a novel mutation (N1575Y) in  
601 the voltage gated sodium channel of *Anopheles gambiae*. *Proceedings of the National*  
602 *Academy of Sciences of the United States of America* **109**, 6614-6619 (2012).
- 603 18 Messer, P. W. & Petrov, D. A. Population genomics of rapid adaptation by soft selective  
604 sweeps. *Trends Ecol. Evol.* **28**, 659-669, doi:10.1016/j.tree.2013.08.003 (2013).
- 605 19 Mathias, D. *et al.* Spatial and temporal variation in the *kdr* allele L1014S in *Anopheles*  
606 *gambiae s.s.* and phenotypic variability in susceptibility to insecticides in Western Kenya.  
607 *Malaria Journal* **10**, e10 (2011).
- 608 20 Norris, L. C. *et al.* Adaptive introgression in an African malaria mosquito coincident with the  
609 increased usage of insecticide-treated bed nets. *Proceedings of the National Academy of*  
610 *Sciences of the United States of America* **112**, 815-820, doi:10.1073/pnas.1418892112  
611 (2015).
- 612 21 Lawniczak, M. K. N. *et al.* Widespread Divergence Between Incipient *Anopheles gambiae*  
613 Species Revealed by Whole Genome Sequences. *Science* **330**, 512-514 (2010).
- 614 22 Weetman, D. *et al.* Association mapping of insecticide resistance in wild *Anopheles gambiae*  
615 populations: major variants identified in a low-linkage disequilibrium genome. *PLoS ONE* **5**,  
616 e13140 (2010).
- 617 23 Ranson, H. *et al.* Pyrethroid resistance in African anopheline mosquitoes: what are the  
618 implications for malaria control? *Trends Parasitol* **27**, 91-97 (2011).
- 619 24 Hargreaves, K. *et al.* *Anopheles funestus* resistant to pyrethroid insecticides in South Africa.  
620 *Medical and Veterinary Entomology* **14**, 181-189, doi:10.1046/j.1365-2915.2000.00234.x  
621 (2000).
- 622 25 Maharaj, R., Mthembu, D. J. & Sharp, B. L. Impact of DDT re-introduction on malaria  
623 transmission in KwaZulu-Natal. *Samj South African Medical Journal* **95**, 871-874 (2005).
- 624 26 Wondji, C. S. *et al.* Mapping a Quantitative Trait Locus (QTL) conferring pyrethroid resistance  
625 in the African malaria vector *Anopheles funestus*. *BMC Genomics* **34** (2007).
- 626 27 Wondji, C. S. *et al.* Two duplicated P450 genes are associated with pyrethroid resistance in  
627 *Anopheles funestus*, a major malaria vector. *Genome Res.* **19**, 452-459,  
628 doi:10.1101/gr.087916.108 (2009).
- 629 28 Mitchell, S. *et al.* Identification and validation of a gene causing cross-resistance between  
630 insecticide classes in *Anopheles gambiae* from Ghana. *Proceedings of the National Academy*  
631 *of Sciences of the United States of America* **109**, 6147-6152 (2012).
- 632 29 Müller, P. *et al.* Field-caught permethrin-resistant *Anopheles gambiae* overexpress CYP6P3, a  
633 P450 that metabolises pyrethroids. *PLoS genetics* **4**, e1000286 (2008).

634 30 Edi, C. V. *et al.* CYP6 P450 enzymes and ACE-1 duplication produce extreme and multiple  
635 insecticide resistance in the malaria mosquito *Anopheles gambiae*. *PLoS genetics* **10**,  
636 e1004236-e1004236, doi:10.1371/journal.pgen.1004236 (2014).

637 31 Mitchell, S. N. *et al.* Metabolic and Target-Site Mechanisms Combine to Confer Strong DDT  
638 Resistance in *Anopheles gambiae*. *PLoS ONE* **9**, doi:10.1371/journal.pone.0092662 (2014).

639 32 Riveron, J. M. *et al.* A single mutation in the GSTe2 gene allows tracking of metabolically  
640 based insecticide resistance in a major malaria vector. *Genome Biology* **15**, 20,  
641 doi:10.1186/gb-2014-15-2-r27 (2014).

642 33 Ibrahim, S. S. *et al.* Allelic Variation of Cytochrome P450s Drives Resistance to Bednet  
643 Insecticides in a Major Malaria Vector. *PLoS genetics* **11**, e1005618,  
644 doi:10.1371/journal.pgen.1005618 (2015).

645 34 Neafsey, D. E. *et al.* SNP Genotyping Defines Complex Gene-Flow Boundaries Among African  
646 Malaria Vector Mosquitoes. *Science* **330**, 514-517, doi:10.1126/science.1193036 (2010).

647 35 Teo, Y. Y., Small, K. S. & Kwiatkowski, D. P. Methodological challenges of genome-wide  
648 association analysis in Africa. *Nature Reviews Genetics* **11**, 149-160, doi:10.1038/nrg2731  
649 (2010).

650 36 Wilding, C. S., Weetman, D., Steen, K. & Donnelly, M. J. High, clustered, nucleotide diversity  
651 in the genome of *Anopheles gambiae* revealed by SNP discovery through pooled-template  
652 sequencing: implications for high-throughput genotyping protocols. *BMC Genomics* **10**, e320  
653 (2009).

654 37 Ramphul, U. *et al.* Insecticide resistance and its association with target-site mutations in  
655 natural populations of *Anopheles gambiae* from eastern Uganda. *Transactions of the Royal  
656 Society of Tropical Medicine and Hygiene* **103**, 1121-1126, doi:10.1016/j.trstmh.2009.02.014  
657 (2009).

658 38 Ochomo, E., Bayoh, N.M., Kamau, L., Atieli, F., Vulule, J., Ouma, C., Ombok, M., Njagi, K., Soti,  
659 D., Mathenge, E., Muthami, L., Kinyari, T., Subramaniam, K., Kleinschmidt, I., Donnelly, M.J.  
660 and Mbogo, C. . Pyrethroid susceptibility of malaria vectors in four districts of western  
661 Kenya. *Parasites and Vectors* **7**, e310 (2014).

662 39 Verhaeghen, K. *et al.* Spatio-temporal patterns in *kdr* frequency in permethrin and DDT  
663 resistant *Anopheles gambiae* s.s. from Uganda. *American Journal of Tropical Medicine and  
664 Hygiene* **82**, 566-573, doi:10.4269/ajtmh.2010.08-0668 (2010).

665 40 Karasov, T., Messer, P. W. & Petrov, D. A. Evidence that adaptation in *Drosophila* is not  
666 limited by mutation at single sites. *PLoS Genet.* **6**, 10, doi:10.1371/journal.pgen.1000924  
667 (2010).

668 41 Reimer, L. *et al.* Relationship between *kdr* mutation and resistance to pyrethroid and DDT  
669 insecticides in natural populations of *Anopheles gambiae*. *Journal of Medical Entomology* **45**,  
670 260-266 (2008).

671 42 Coluzzi, M., Sabatini, A., Petrarca, V. & Di Deco, M. A. Chromosomal differentiation and  
672 adaptation to human environments in the *Anopheles gambiae* complex. *Transactions of the  
673 Royal Society of Tropical Medicine and Hygiene* **73**, 483-497 (1979).

674 43 Cheng, C. D. *et al.* Ecological genomics of *Anopheles gambiae* along a latitudinal cline: a  
675 population-resequencing approach. *Genetics* **190**, 1417-+, doi:10.1534/genetics.111.137794  
676 (2012).

677 44 Riehle, M. M. *et al.* Natural malaria infection in *Anopheles gambiae* is regulated by a single  
678 genomic control region. *Science* **312**, 577-579 (2006).

679 45 Matoke-Muhia, D. *et al.* Decline in frequency of the 2La chromosomal inversion in  
680 *Anopheles gambiae* (s.s.) in Western Kenya: correlation with increase in ownership of  
681 insecticide-treated bed nets. *Parasites & vectors* **9**, 334, doi:10.1186/s13071-016-1621-3  
682 (2016).

683 46 Harris, C., Rousset, F., Morlais, I., Fontenille, D. & Cohuet, A. Low linkage disequilibrium in  
684 wild *Anopheles gambiae* s.l. populations. *BMC Genetics* **11**, doi:10.1186/1471-2156-11-81  
685 (2010).

686 47 Bong, L. J. & Zairi, J. Temporal fluctuations of insecticides resistance in *Musca domestica* Linn  
687 (Diptera: Muscidae) in Malaysia. *Tropical Biomedicine* **27**, 317-325 (2010).

688 48 Wilding, C. S. *et al.* Parallel evolution or purifying selection, not introgression, explains  
689 similarity in the pyrethroid detoxification linked GSTE4 of *Anopheles gambiae* and *An.*  
690 *arabiensis*. *Molecular Genetics and Genomics* **290**, 201-215, doi:10.1007/s00438-014-0910-9  
691 (2015).

692 49 Weetman, D., Wilding, C. S., Steen, K., Pinto, J. & Donnelly, M. J. Gene flow-dependent  
693 genomic divergence between *Anopheles gambiae* M and S forms. *Molecular Biology and*  
694 *Evolution* **29**, 279-291, doi:10.1093/molbev/msr199 (2012).

695 50 Barrett, J. C., Fry, B., Maller, J. & Daly, M. J. W. Haploview: analysis and visualization of LD  
696 and haplotype maps. *Bioinformatics* **21**, 263-265 (2005).

697 51 Corander, J., Marttinen, P., Siren, J. & Tang, J. Enhanced Bayesian modelling in BAPS  
698 software for learning genetic structures of populations. *BMC Bioinformatics* **9**, 14,  
699 doi:10.1186/1471-2105-9-539 (2008).

700 52 Corander, J. & Marttinen, P. Bayesian identification of admixture events using multilocus  
701 molecular markers. *Mol. Ecol.* **15**, 2833-2843, doi:10.1111/j.1365-294X.2006.02994.x (2006).

702 53 Visscher, P. M. & Le Hellard, S. Simple method to analyze SNP-based association studies  
703 using DNA pools. *Genet. Epidemiol.* **24**, 291-296, doi:10.1002/gepi.10240 (2003).

704 54 Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics Analysis  
705 Version 7.0 for Bigger Datasets. *Mol. Biol. Evol.* **33**, 1870-1874, doi:10.1093/molbev/msw054  
706 (2016).

707 55 White, B. J. *et al.* Molecular karyotyping of the 2LA inversion in *Anopheles gambiae*.  
708 *American Journal Of Tropical Medicine And Hygiene* **76**, 334-339 (2007).

709 56 Rose, P. W. *et al.* The RCSB Protein Data Bank: new resources for research and education.  
710 *Nucleic Acids Res.* **41**, D475-D482, doi:10.1093/nar/gks1200 (2013).

711 57 Soding, J., Biegert, A. & Lupas, A. N. The HHpred interactive server for protein homology  
712 detection and structure prediction. *Nucleic Acids Research* **33**, W244-248,  
713 doi:10.1093/nar/gki408 (2005).

714 58 Song, Y. F. *et al.* High-Resolution Comparative Modeling with RosettaCM. *Structure* **21**, 1735-  
715 1742, doi:10.1016/j.str.2013.08.005 (2013).

716 59 Schymkowitz, J. *et al.* The FoldX web server: an online force field. *Nucleic Acids Res.* **33**,  
717 W382-W388, doi:10.1093/nar/gki387 (2005).

718 60 Guerois, R., Nielsen, J. E. & Serrano, L. Predicting changes in the stability of proteins and  
719 protein complexes: A study of more than 1000 mutations. *J. Mol. Biol.* **320**, 369-387,  
720 doi:10.1016/s0022-2836(02)00442-4 (2002).

721 61 Dehouck, Y., Kwasigroch, J., Gilis, D. & Rooman, M. PoPMuSiC 2.1: a web server for the  
722 estimation of protein stability changes upon mutation and sequence optimality. *BMC*  
723 *Bioinformatics* **12**, e151 (2011).

724 62 Adzhubei, I. A. *et al.* A method and server for predicting damaging missense mutations. *Nat.*  
725 *Methods* **7**, 248-249, doi:10.1038/nmeth0410-248 (2010).

726 63 Lim, E. *et al.* T3DB: a comprehensively annotated database of common toxins and their  
727 targets. *Nucleic Acids Res.* **38**, D781-D786, doi:10.1093/nar/gkp934 (2010).

728 64 Sevrioukova, I. F. & Poulos, T. L. Structure and mechanism of the complex between  
729 cytochrome P4503A4 and ritonavir. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 18422-18427,  
730 doi:10.1073/pnas.1010693107 (2010).

731 65 Maweje, H. D. *et al.* Insecticide resistance monitoring of field-collected *Anopheles gambiae*  
732 s.l. populations from Jinja, eastern Uganda, identifies high levels of pyrethroid resistance.

733            *Medical and Veterinary Entomology* **27**, 276-283, doi:10.1111/j.1365-2915.2012.01055.x  
734            (2013).  
735    66        Kilama, M. *et al.* Estimating the annual entomological inoculation rate for *Plasmodium*  
736            *falciparum* transmitted by *Anopheles gambiae s.l.* using three sampling methods in three  
737            sites in Uganda. *Malaria Journal* **13**, 111, doi:10.1186/1475-2875-13-111 (2014).  
738    67        Bass, C. *et al.* Detection of knockdown resistance (*kdr*) mutations in *Anopheles gambiae*: a  
739            comparison of two new high-throughput assays with existing methods. *Malaria Journal* **6**,  
740            e111 (2007).  
  
741  
  
742

743 **Acknowledgements**

744 The project described was supported by Award Numbers U19AI089674 and R01AI082734 from the  
745 National Institute of Allergy and Infectious Diseases (NIAID) and by the Innovative Vector Control  
746 Consortium. The content is solely the responsibility of the authors and does not necessarily  
747 represent the official views of the NIAID or NIH.

748 **Author Contributions**

749 DW, CSW, PM, LMO and MJD conceived/designed the research. DW, CSW, PM, EO, ATI, KS, EJR, JCM,  
750 HDM conducted the experiments. DEN, DJR and LMO contributed new reagents and/or analytical  
751 tools. DW, CSW, AIT, DEN, DJR, MJD analyzed the data. DW and MJD wrote the manuscript. All  
752 authors read, corrected and approved the manuscript.

753

754 **Additional Information**

755 The author(s) declare no competing financial interests.

756 **Figure Legends**

757 **Figure 1.** Analysis of sources of within-population structure in individually genotyped adult females.  
758 Plots show  $F_{ST}$  values for each SNP compared among each of the three clusters identified using BAPS.  
759 The upper bar shows the location of the 2La inversion region.

760 **Figure 2.** Family pool association analysis for permethrin resistance (A) corrected for humidity as a  
761 covariate; (B) corrected for humidity and PC1 (a proxy for 2La polymorphism; see Fig 1). Test  
762 probabilities are shown for each SNP arranged on a physical scale across chromosomes, centromere  
763 positions are shown by solid vertical lines; chromosome breaks by dashed lines. The purple bar  
764 indicates the 2La inversion region.

765 **Figure 3.** Repeatability of SNP associations in independent field samples from Uganda and Kenya  
766 with three different pyrethroid insecticides (see Table S4 for full results). The dashed line indicates  
767  $P=0.05$ . Uga = Uganda; insecticides used for bioassays:  $\lambda$ -cy = lambda cyhalothrin; delta =  
768 deltamethrin; perm=permethrin. The year of collection is shown.

769 **Figure 4.** Evolutionary relationships of the resistant (Phenylalanine) and susceptible (Leucine) alleles  
770 of *CYP4J5* marker L43F in Ugandan samples. The optimal tree with sum of branch lengths = 0.062 is  
771 shown. Percentage bootstrap values (1000 replicates) are shown next to branches. The sample prefix  
772 indicates mosquito origin (Jinja or Oyam, Uganda)

773

774 **Tables**

775 **Table 1.** Stepwise regression of pyrethroid resistance by the environmental correlate humidity and  
 776 candidate SNPs among families of the Tororo discovery population. The number in the predictor  
 777 column indicates the order of entry of the SNP into the model.

| chromosome | gene type                   | minor allele frequency | predictor         | model $r^2$ | P-value |
|------------|-----------------------------|------------------------|-------------------|-------------|---------|
|            | environmental correlate     |                        | 1. Humidity       | 0.167       | 0.00002 |
| 2L         | metabolic (P450)            | 0.37                   | 2. <i>Cyp4j5</i>  | 0.284       | 0.00009 |
| 2L         | target site ( <i>Vgsc</i> ) | 0.06                   | 3. L1014S         | 0.335       | 0.005   |
| 2L         | metabolic (COE)             | 0.49                   | 4. <i>Coae1d</i>  | 0.363       | 0.026   |
| 2L         | metabolic (P450)            | 0.45                   | 5. <i>Cyp4j10</i> | 0.393       | 0.021   |

778

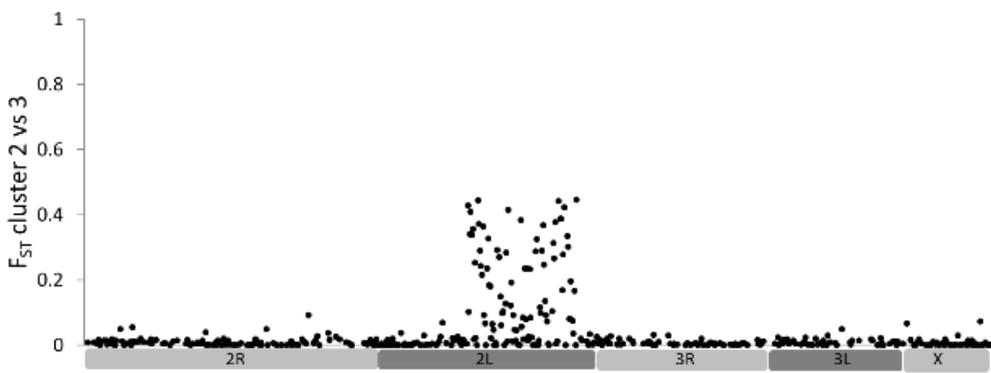
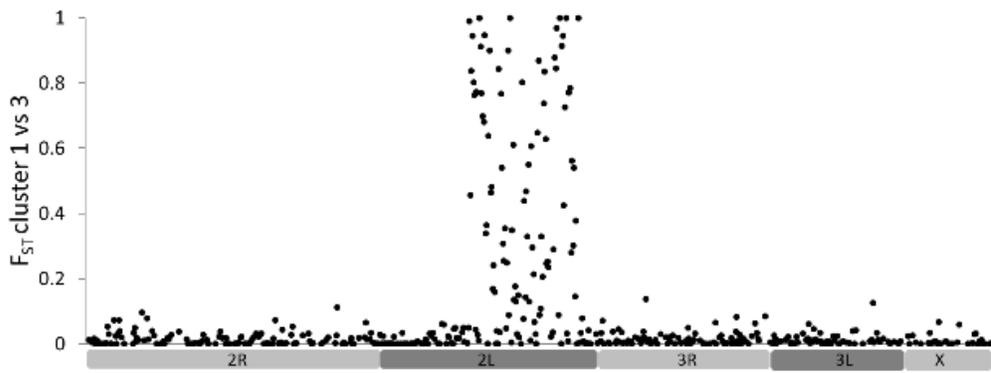
779 **Table 2.** Allele frequencies and association with resistance measured by odds ratios for the four  
 780 candidate SNPs in populations from Kenya and Uganda in which reproducibility was tested.  
 781 Standard deviations are measured across populations.

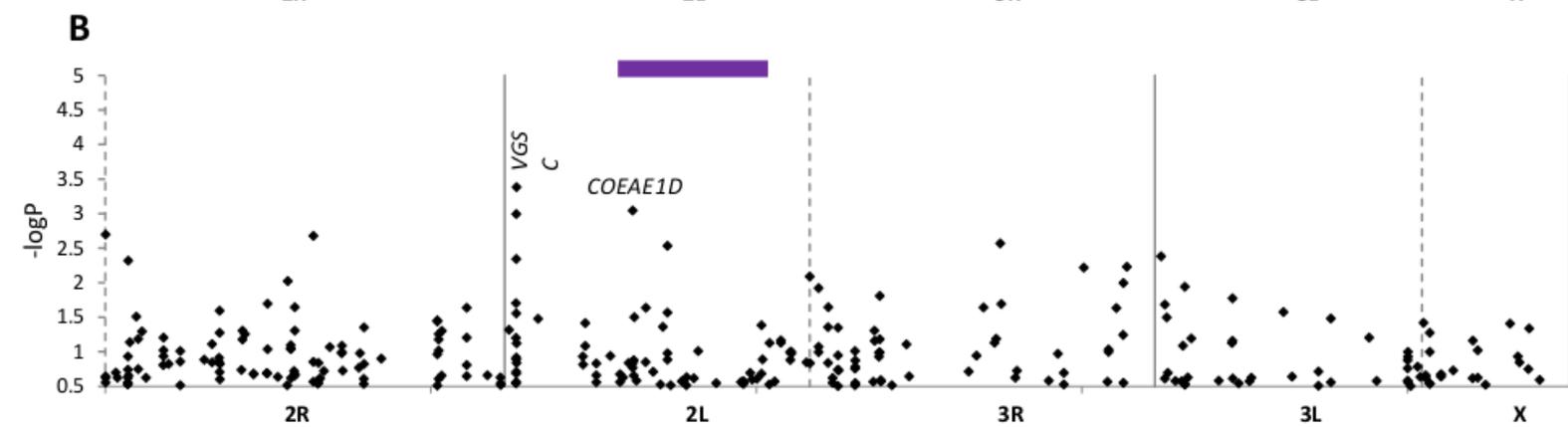
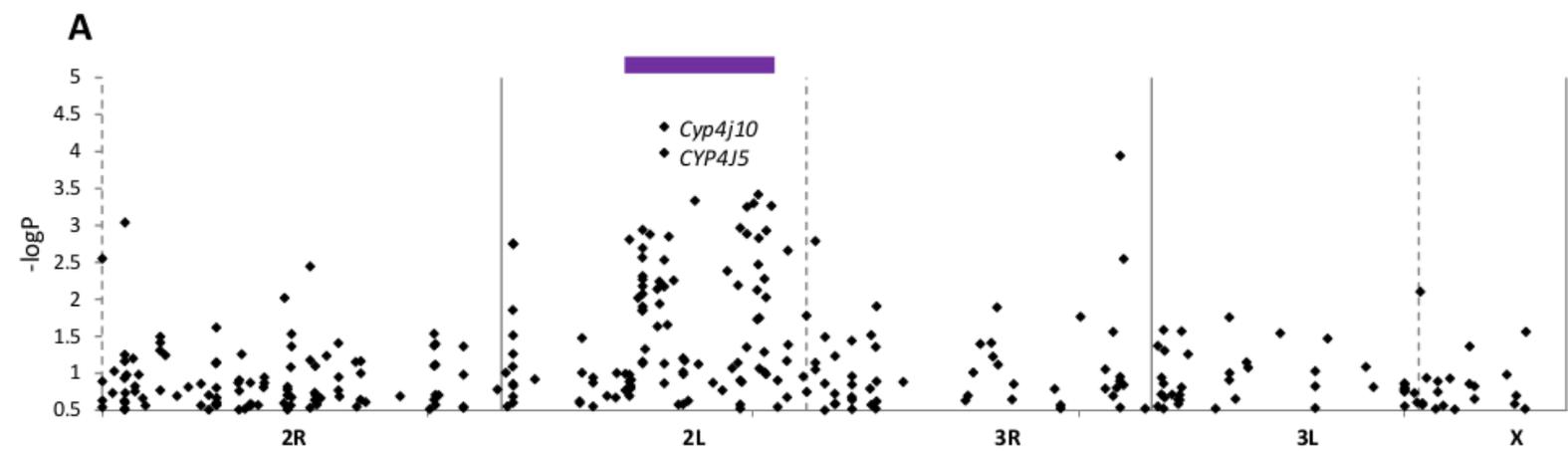
| SNP in gene    | Populations genotyped | Total allele count | Mean resistant allele frequency (st. dev) | Mean odds ratio (st. dev) |
|----------------|-----------------------|--------------------|---|---------------------------|
| <i>Cyp4j5</i>  | 7                     | 1370               | 0.61 (0.13)                               | 2.94 (1.10)               |
| <i>Vgsc</i>    | 7                     | 1068               | 0.84 (0.22)                               | 3.59 (3.71)               |
| <i>Coae1d</i>  | 5                     | 738                | 0.53 (0.04)                               | 1.93 (0.87)               |
| <i>Cyp4j10</i> | 3                     | 502                | 0.47 (0.05)                               | 1.04 (0.24)               |

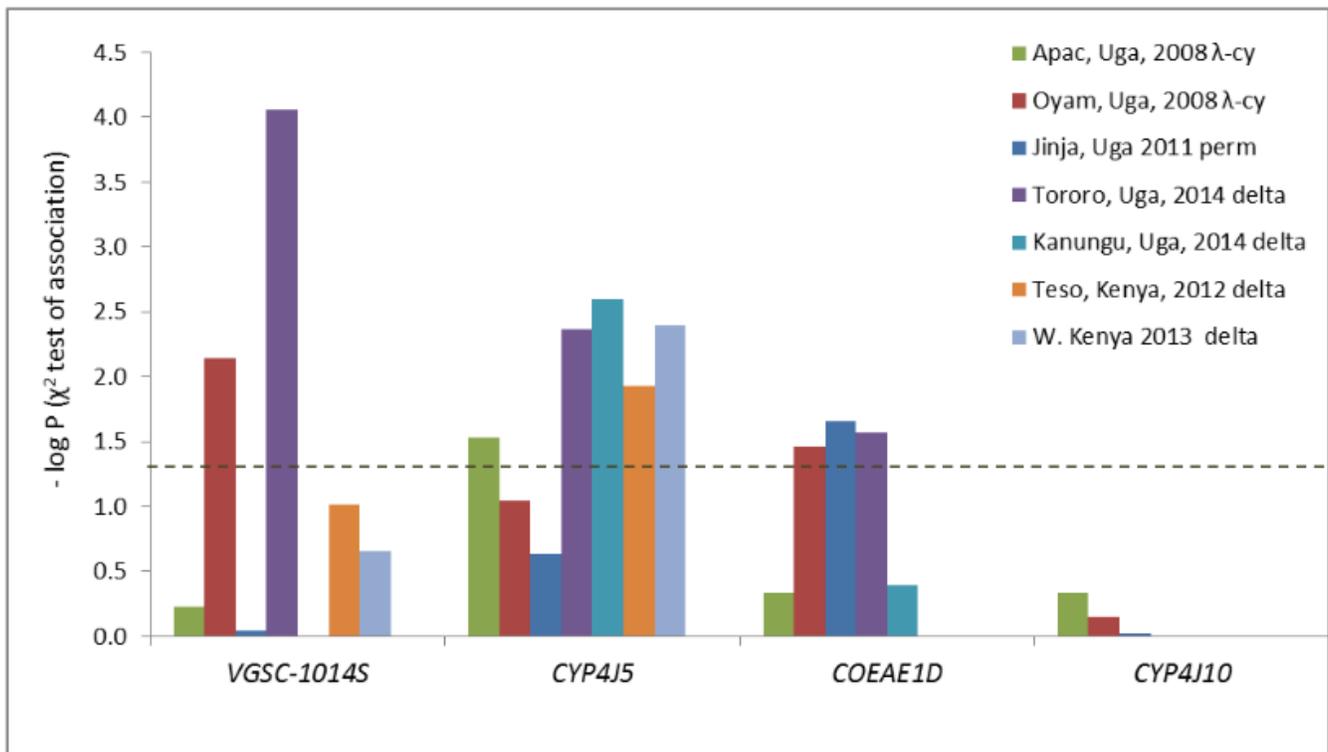
782

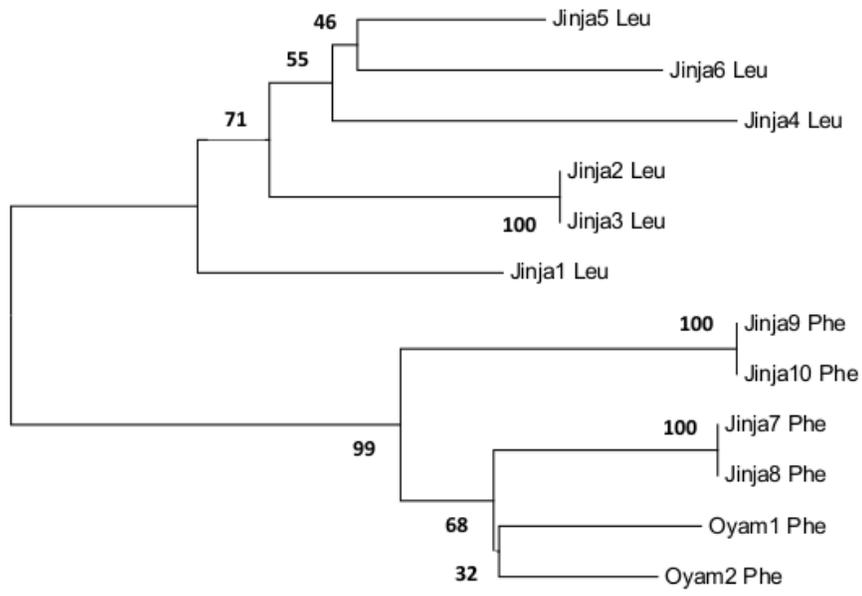
783

784









0.002