1    **Genetic variation at the *Cyp6m2* putative insecticide resistance locus in**

2    ***Anopheles gambiae* and *Anopheles coluzzii***


3    **Authors**

4    Martin G. Wagah[1,*], Petra Korlević[1,2], Christopher Clarkson[1], Alistair Miles[3], The *Anopheles gambiae*

5    1000 Genomes Consortium[1], Mara K. N. Lawniczak[1], Alex Makunin[1].


6    **Contact information**

7    *Corresponding authors: mw21@sanger.ac.uk


8    **Affiliations**

9    1.  Wellcome Sanger Institute, Hinxton, Cambridgeshire, CB10 1SD, United Kingdom

10   2.  European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton,

11       Cambridgeshire, CB10 1SD, United Kingdom.

12   3.  University of Oxford, Wellcome Trust Centre for Human Genetics, Oxford, OX3 7BN, United

13       Kingdom

14

15

16

17

18

19

20

## 21 Abstract

### 22 Background

23 The emergence of insecticide resistance is a major threat to malaria control programmes in Africa,

24 with many different factors contributing to insecticide resistance in its vectors, *Anopheles* mosquitoes.

25 *CYP6M2* has previously been recognized as an important candidate in cytochrome P450-mediated

26 detoxification in *Anopheles* mosquitoes. As it has been implicated in resistance against pyrethroids,

27 organochlorines and carbamates, its broad metabolic activity makes it a potential agent in insecticide

28 cross-resistance. Currently, allelic variation within the *Cyp6m2* gene remains unknown.

### 29 Results

30 Here, we use Illumina whole-genome sequence data from Phase 2 of the *Anopheles gambiae* 1000

31 Genomes Project (Ag1000G) to examine genetic variation in the *Cyp6m2* gene across 16 populations

32 in 13 countries comprising *Anopheles gambiae* and *Anopheles coluzzii* mosquitoes. We find 15

33 missense biallelic substitutions at high frequency (defined as >5% frequency in one or more

34 populations), that fall into five distinct haplotype groups that carry the main high frequency variants:

35 A13T, D65A, E328Q, Y347F, I359V and A468S. We examine whether these alleles show evidence of

36 selection either through potentially modified enzymatic function or by being linked to variants that

37 change the transcriptional profile of the gene. Despite consistent reports of *Cyp6m2* upregulation and

38 metabolic activity in insecticide resistant Anophelines, we find no evidence of directional selection

39 occurring on these variants or on the haplotype clusters in which they are found.

### 40 Conclusion

41 Our results imply that emerging resistance associated with *Cyp6m2* is potentially driven by distant

42 regulatory loci such as transcription factors rather than by its missense variants, or that other genes

43 are playing a more significant role in conferring metabolic resistance.

## 44 Keywords

46   **Background**

47   Malaria remains a pernicious public health problem that plagues the African region, which has over

48   90% of the world's malaria cases and deaths [1]. Although concerted vector control interventions such

49   as long lasting insecticidal nets (LLINs) and indoor residual spraying (IRS) have led to the attainment

50   of key milestones, global progress has stagnated and case numbers are stable or on the rise in many

51   countries in Africa [1-3]. This is due to multiple factors, including the emergence of insecticide

52   resistance, which threaten the effectiveness of vector control interventions [4].

53   The most well understood mechanisms of insecticide resistance are classified into two main functional

54   categories depending on the underlying genes involved: target-site insensitivity and metabolic

55   sequestration and detoxification. Both types may occur concurrently within a single population or even

56   within a single mosquito [5-7]. These mechanisms have led to increasing resistance to all four

57   common insecticide classes — pyrethroids, organochlorines, carbamates and organophosphates — in

58   all major malaria vectors across Africa [7, 8].

59   Metabolic detoxification occurs mainly through the elevated activity of large and functionally diverse

60   multigene enzyme families: glutathione S-transferases (GSTs), carboxylesterases (COEs) and

61   cytochrome P450 monooxygenases (P450s) [7, 9].  Although a few candidates in these enzyme

62   families have been directly associated with resistance, our understanding of metabolic resistance has

63   lagged far behind that of pyrethroid target-site resistance, chiefly due to its complexity and the lack of

64   associated causal mutations [10]. This is despite the fact that metabolic resistance is often considered

65   a greater threat to mosquito control [9], especially since the only widely accepted occurrence of

66   malaria vector control failure was attributed to the elevated expression of resistance-associated

67   P450s in *An. funestus* [11-13]. A comprehensive understanding of metabolic resistance must

68   therefore involve disambiguating the roles that individual enzymes play and the genetic backgrounds

69   that underlie their significance in vector populations.

70   The CYP6M2 enzyme exhibits complex insecticide metabolism associated with multiple binding

71   modes for insecticides [14]. Its gene is located within a cluster of 14 Cyp6 P450 genes on

72   chromosome 3R of *An. gambiae* [15]*,* and is among the 111 known P450 genes across the *An.*

73    *gambiae* genome [16, 17]. In this genomic region, *Cyp6m2* is nested within a sub-cluster of P450s

74    containing *Cyp6m3* and *Cyp6m4* which have also been associated with xenobiotic detoxification[18].

75    *Cyp6m2* is notably one of the few specific P450s that have shown a consistent association with

76    metabolic resistance [5]. Metabolic resistance is mainly assessed through transcriptional profiling of

77    genes involved in xenobiotic detoxification. Transcriptomic experiments such as quantitative PCR and

78    microarray assays have established a link between *Cyp6m2* overexpression and the resistance

79    phenotype in field populations of *An. gambiae*, *An. coluzzii*, *An. arabiensis* and *An. sinensis,*

80    irrespective of the presence of knock-down resistance (*kdr*) mutations such as L995F or L995S in the

81    voltage gated sodium channel (VGSC) [5, 19-21]. In DDT resistant *An. gambiae* in Ghana, *Cyp6m2*

82    has been found to be overexpressed 3.2 to 5.2-fold in combination with the upregulation of additional

83    P450s like *Cyp6z2* [18]. In DDT resistant *An. coluzzii* collected in Benin, *Cyp6m2* was also found to

84    be overexpressed 1.2 to 4.6-fold in combination with *Gste2* from the *GST* gene family and in the

85    presence of fixed *kdr* alleles in the *Vgsc* gene [22]. In Nigeria, the 2.4 to 2.7-fold upregulation of

86    *Cyp6m2* was found to be associated with high levels of permethrin resistance [5] and *An. gambiae*

87    that exhibited a strong resistance to bendiocarb in In Côte d'Ivoire also had an elevated (up to 8-fold)

88    expression of the *Cyp6m2* gene [20]. In the same study, transgenic expression of *Cyp6m2* in

89    *Drosophila melanogaster* was shown to produce resistance to both DDT and bendiocarb. *In vivo*

90    functional analysis of multi-tissue overexpression induced by genetic modification has also shown

91    *Cyp6m2* to be sufficient in conferring resistance to permethrin and deltamethrin [23]. However, this

92    overexpression also increased the mosquitos' susceptibility to the organophosphate malathion.

93    Collectively, these studies indicate that *Cyp6m2* can confer metabolic resistance against insecticides

94    in 3 of the 4 known classes: both type I and type II pyrethroids [14, 18, 23], organochlorines [24], and

95    carbamates [20]. It therefore has a high potential for cross-resistance, which may make the problem

96    of malaria vector control even more intractable by limiting the options available to malaria control

97    programs for insecticide rotation or combination. The negative cross-resistance associated with

98    malathion hereby points to potential mitigating strategies [23].

99    The frequent association of *Cyp6m2* with insecticide resistance described above warrants further

100    investigation into whether there is evidence of copy number variation (CNV) or missense mutations at

101    the locus. CNVs have been implicated in augmenting gene dosage leading to increased transcription

102    of metabolic enzymes [25, 26]. A genome-wide CNV analysis conducted on the Ag1000G dataset and

103    described in detail elsewhere [25] found CNVs to be significantly enriched in metabolic resistance-

104    associated gene families and to be undergoing positive selection. These CNVs were identified across

105    P450s (such as *Cyp9k1* and at both the *Cyp6z3–Cyp6z1* and the *Cyp6aa1–Cyp6p2* gene clusters)

106    and GSTs (at the *Gstu4–Gste3* cluster). However, CNVs across the *Cyp6m2* locus were found to be

107    rare, even in populations that are known to exhibit *Cyp6m2*-mediated resistance [25]. This indicates

108    that CNVs alone are not sufficient to explain the widespread occurrence of the *Cyp6m2*-associated

109    resistance phenotype: additional factors such as allelic variation might contribute to resistance

110    associated with *Cyp6m2* activity.

111    Allelic variation can play an additional role in P450-mediated resistance by modifying either enzyme

112    catalytic activity or gene expression levels [27]. Allelic variation has been shown to be key in inducing

113    high metabolic efficiency of *Cyp6P9b* and in conferring metabolic resistance to *An. funestus* [28].

114    Allelic variants in metabolic genes have also been identified to reliably and reproducibly associate with

115    resistance, such as in *Cyp4J5* and *Coeae1d* in *An. gambiae,* and can serve as diagnostic markers of

116    phenotypic resistance [29]. However, there is still a paucity of information about allelic variation

117    associated with metabolic resistance when compared to the well-characterized target-site mutations

118    [29]. Mutations that may modulate metabolic resistance by either altering function or modifying

119    expression in *Cyp6m2* are yet to be described.

120    Following the consistent association of *Cyp6m2* with insecticide resistance in many populations, we

121    examine whole-genome Illumina sequence data from phase 2 of the *Anopheles gambiae* 1000

122    Genomes Project (Ag1000G) [30] which consists of 1,142 wild-caught mosquitoes sequenced to a

123    mean depth above 14x, and report a comprehensive analysis of genetic variation within the *Cyp6m2*

124    gene. We also examine the wider haplotypes around *Cyp6m2* spanning across the *Cyp6m* sub cluster

125    and the larger *Cyp6* supercluster for signatures of selection.

126    **Results**

127    **Cyp6m2 non-synonymous nucleotide variation**

128    Short-read whole-genome sequence data from the Ag1000G phase 2 data resource [30] were used to

129    investigate genetic variation at the *Cyp6m2* locus across 16 populations of  *An. gambiae* and *An.*

130  *coluzzii* (n = 1,142 total individuals) collected between 2000 and 2012 [*Table 1, Additional file 1*]. The

131  single nucleotide polymorphisms (SNPs) we studied here were discovered and QC'd using methods

132  described elsewhere [31]. We focused on SNPs that change the amino acid sequence of the *CYP6M2*

133  enzyme as they have a potential functional role in *Cyp6m2*-associated insecticide resistance (n = 193)

134  [*Additional file 2*]. As putative resistance variants under selection pressure from insecticides are

135  expected to increase in frequency over time, we subsequently computed allele frequencies for every

136  non-synonymous SNP in each population with reference to species and country of origin. We filtered

137  the list to focus only on those variants that were at high frequency within populations or across

138  populations (defined as >5% frequency in one or more populations). In total, this resulted in 15 non-

139  synonymous variants that we further explored [*Table 1*].

140  Analysis of the patterns of polymorphism of *Cyp6m2* from different populations showed both relative

141  homogeneity within some geographical regions and distinct variants across different regions. The

142  variants with the highest overall frequency were **I359V** (16%) and **D65A** (6%) [*Table 1*]. The most

143  widespread variant was **I359V**, which was present in West, Central and East African populations of

144  both *An. gambiae* and *An. coluzzii*. Populations with the highest frequency of **I359V** were Gabon

145  (49%) and Ghana (25%) for *An. gambiae*, and Guinea (37.5%) for *An. coluzzii*. Another mutation,

146  **E328Q**, was found across West Africa's *An. coluzzii* populations in Burkina Faso, Côte d'Ivoire,

147  Ghana, Guinea and The Gambia and ranged in frequency from 6.2 to 13.6%. Several variants were

148  found to exceed the 5% threshold only in one or two populations: **A13T** and **Y347F,** in Angola's *An.*

149  *coluzzii* (39.7%) and in Kenya (52.1%) respectively and **D65A** only in Gabon's *An. gambiae* and in

150  Kenya's populations at 42.8% and 52.1%, respectively [*Table 1*].

151

152  Table1. Allele frequencies of common *Cyp6m2* variants.

153

| Variants | | | | | Population allele frequency (%) | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Position[1] | | Ag[2] | All | AOco l[3] | GHc ol | BFc ol | CIc ol | GNc ol | G W | G M | CMga m | GHga m | BFga m | GNga m | GAga m | UGga m | GQga m | FRga m | KE |
| 69289 45 | G> A | A13T | 2.8 | 39.7 | 0 | 0.7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 69290 47 | G> C | G47R | 0.6 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 69290 50 | T> A | S48T | 2.8 | 9 | 0.9 | 4 | 2.1 | 0 | 3.3 | 2.3 | 3.2 | 4.2 | 2.7 | 2.5 | 0 | 2.2 | 0 | 0 | 0 |

| Position | | Codon | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 6929102 | A>C | D65A | 6.0 | 0 | 0 | 0 | 0 | 0 | 0.5 | 0 | 1.9 | 4.2 | 2.7 | 3.8 | 42.8 | 2.7 | 0 | 0 | 52.1 |
| 6929375 | A>T | K156I | 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1.4 | 0 | 44.4 | 0 | 0 |
| 6929506 | A>G | N200D | 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1.4 | 0 | 44.4 | 0 | 0 |
| 6929770 | T>A | S288T | 1.1 | 0 | 5.5 | 0.7 | 0.7 | 0 | 0 | 0 | 0.8 | 4.2 | 2.7 | 1.3 | 0.7 | 1.3 | 0 | 0 | 0 |
| 6929881 | G>A | E325K | 1.2 | 0 | 7.3 | 3.3 | 2.1 | 0 | 3.8 | 2.3 | 0 | 0 | 0.5 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6929890 | G>C | E328Q | 2.9 | 0.6 | 13.6 | 7.3 | 7 | 12.5 | 9.9 | 6.2 | 0.3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6929948 | A>T | Y347F | 2.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.2 | 0 | 0 | 0 | 0 | 1.8 | 0 | 0 | 52.1 |
| 6929983 | A>G | I359V | 16.0 | 0 | 4.5 | 0 | 7 | 37.5 | 16.5 | 20 | 19.4 | 25 | 17.4 | 16.3 | 48.6 | 19.6 | 0 | 0 | 0 |
| 6930206 | C>T | P407L | 0.7 | 0 | 2.7 | 0.7 | 7 | 0 | 0 | 0 | 0.2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6930242 | A>T | E419V | 0.4 | 5.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6930269 | A>G | K428R | 2.0 | 0 | 0 | 0.7 | 0 | 0 | 0.5 | 0.8 | 3.5 | 0 | 2.7 | 2.5 | 0.7 | 6.3 | 0 | 0 | 0 |
| 6930388 | G>T | A468S | 2.6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.3 | 0 | 0 | 0 | 42 | 0 | 0 | 0 | 0 |

154 [1]Position relative to the AgamP4 reference sequence, chromosome 3R.

155 [2]Codon numbering according to *Anopheles gambiae* AGAP008212-RA transcript in geneset AgamP4.12.

156 [3]AOcol=Angola *coluzzii*; GHcol=Ghana *coluzzii*; BFcol=Burkina Faso *coluzzii*; CIcol=Côte d'Ivoire *coluzzii*; GNcol=Guinea

157 *coluzzii*; GW=Guinea Bissau; GM = The Gambia; CMgam=Cameroon; GHgam = Ghana *gambiae*; BFgam = Burkina Faso

158 *gambiae*; GNgam = Guinea *gambiae*; GAgam=Gabon *gambiae*; UGgam=Uganda *gambiae*; GQgam = Equatorial Guinea

159 *gambiae;* FRgam=Mayotte *gambiae*; KE=Kenya.

160

**161 Haplotypic backgrounds of non-synonymous alleles**

162 The Ag1000G data resource contains data that not only spans across exonic regions of any given

163 gene, but also intronic and intergenic regions. This enables a comprehensive analysis of haplotypes

164 that contain putative insecticide resistance alleles, but is constrained by the fact that this resource

165 does not contain samples whose resistance status or *Cyp6m2* expression levels are known.

166 Selection pressure acting upon missense variants or linked cis regulatory variants is likely to affect the

167 haplotype structure of the gene. To study haplotype structure at *Cyp6m2*, we extracted biallelic SNPs

168 across the entire 1689bp *Cyp6m2* gene to calculate the number of SNP differences between all pairs

169 of 2,284 haplotypes derived from the mosquitoes. We identified a clustering threshold of seven SNPs

170  where the haplotype clusters corresponded to the haplotypes carrying the high frequency alleles

171  [*Table1*, *Figure 1*]. We found that these haplotypes could mostly be grouped into five distinct clusters

172  (labelled C1-C5): C1 contained haplotypes carrying **A13T**; C2 contained most haplotypes carrying

173  **D65A**, **A468S**, and some haplotypes carrying **I359V**; C3 contained most haplotypes carrying both

174  **D65A** and **Y347F**, and C5 contained haplotypes carrying **E328Q. C4** contained haplotypes with no

175  signature missense mutation [*Figure 1*].

176

177  **Figure 1. Hierarchical clustering of *Cyp6m2* haplotypes.**

178

179  Top: a dendrogram showing hierarchical clustering of haplotypes derived from wild-caught mosquitoes.

180  The colour bar indicates the population of origin for each haplotype.

181  Bottom: high frequency (>□5%) alleles identified within each haplotype (white = reference allele; black = alternative allele). The

182  lowest margin labels the major haplotype clusters.

183  Overall, haplotype cluster distribution resembled the whole genome groupings of individuals described

184  elsewhere using our dataset [30]: Cluster C5 contained haplotypes from West African *An. coluzzii*; C4

185  contained *An. gambiae* from West, Central and near-East Africa; and the rest of the clusters

186  contained haplotypes from samples from a single country and species [*Figure 2*]. The variation across

187  the haplotypes largely showed no strict or systematic difference between the two species or across

188  broad geographic regions, which is in line with recent whole genome sequencing reports [31].

189  **Figure 2. Map of haplotype cluster frequencies and distribution.**

190

191  Each pie chart indicates the haplotype group frequencies within specific sampling populations. The sizes of the wedges within

192  the pies are proportional to haplotype group frequencies within the populations. Haplotypes in group C1 carry the **A13T** allele.

193  Haplotypes in group C2 carry **D65A, I359V** and **A468S** alleles. Haplotypes in group C3 carry **D65A** and **Y347F** alleles.

194  Haplotypes in group C5 carry the **E328Q** allele. Haplotypes in group C4 had no defining non-synonymous variant, and wild type

195  (*wt*) haplotypes were all those that did not fall within the C1-C5 clusters.

196

197  We investigated patterns of association among these non-synonymous variants by computing the

198  normalized coefficient of linkage disequilibrium (D') using haplotypes from the Ag1000G phase 2

199  resource. Of the two highest frequency variants, **I359V** was found to be in perfect linkage with **A468S**

200  but this was driven only by one population (Gabon) with most backgrounds carrying **I359V** not

201  showing linkage with any other missense mutations [*Figure 1 & Supplementary Figure 1*]. **D65A** was

202  in perfect linkage with **A468S** and **Y347F**, showing that **D65A** was almost only ever found on

203  haplotypes carrying either **A468S** or **Y347F**. **I359V** and **D65A**, the highest frequency mutations across

204  all populations, were found to be only in moderate linkage disequilibrium (0.36) [*Supplementary*

205  *Figure 1*]. Other variants were found to be in weak linkage disequilibrium with the six main high

206  frequency alleles and segregated independently within their own populations. While we observed

207  some strong associations through linkage disequilibrium analysis across all populations, a deeper

208  investigation revealed that these associations were driven by population specific dynamics in

209  populations (such as Kenya) where we know bottlenecking has been an issue [31]. It is therefore

210  unlikely that the identified variants are conferring some selective advantage against existing

211  insecticide pressures.

212

213  We next explored whether the surrounding genomic region showed a similar hierarchical clustering

214  pattern to *Cyp6m2*, which might be indicative of either dominant demographic effects or selection

215  acting at other linked loci that is having a major impact on variation within *Cyp6m2*. The downstream

216  genes we selected coded for proteins that were 1-to-1 orthologs with *D. melanogaster* genes. We

217  selected *ODR2* [32], *HAM* [33] and *SH2* [34], which were 81280 bases, 457164 bases and 1198636

218  bases downstream of the *Cyp6m2* gene respectively. The distinctive haplotype clustering pattern

219  observed for *Cyp6m2* in the Kenya, Angola and Gabon populations persisted across these genes,

220  indicating that in these populations, the diversity reduction in and downstream of *Cyp6m2* is more

221  likely driven by demography rather than by a selective sweep [*Supplementary Fig. 2-4*]. We also

222  extracted biallelic SNPs across the *Cyp6m* sub cluster of 3 genes (*Cyp6m2*, *Cyp6m3* and *Cyp6m4*)

223  and across the *Cyp6* supercluster of 14 genes within which the *Cyp6m* sub cluster is located (*Cyp6s2,*

224  *Cyp6s1, Cyp6r1, Cyp6n2, Cyp6y2, Cyp6y1, Cyp6m1, Cyp6n1, Cyp6m2, Cyp6m3, Cyp6m4, Cyp6z3,*

225  *Cyp6z2* and *Cyp6z1*), and performed hierarchical clustering across these regions as described above.

226  The typical geographical stratification of haplotypes persisted, suggesting the absence of a selective

227  sweep across this region [*Supplementary Fig.5 & 6*].

228

229  We examined the genetic backgrounds carrying these alleles further by constructing median joining

230  networks (MJNs) [35] using the Ag1000G Phase 2 haplotype data. This enabled us to resolve the

231    radiation of DNA substitutions arising on haplotypes carrying the identified variants. It also allowed us

232    to reconstruct and position intermediate haplotypes while revealing the non-hierarchical relationships

233    between haplotypes that could not be resolved by hierarchical clustering alone. The MJNs were

234    constructed with reference to a maximum edge distance of two SNPs. This ensured that the

235    connected components captured only closely related haplotypes. The resulting MJNs had a close

236    correspondence with the hierarchical clustering output in assignment of haplotypes to clusters (88%

237    overall concordance across all clusters).

238    The median joining networks showed more clearly the distinctive demographic stratification of the high

239    frequency variants that was highlighted by the hierarchical clustering networks [*Figure 3*]. Most nodes

240    containing secondary variants arising from the main nodes were small, which is inconsistent with

241    directional selection where larger nodes are expected. Only one of the **I359V** nodes contained

242    haplotypes from mosquitoes of both species, however the secondary nodes did not contain

243    haplotypes from more than one species. This indicates that although **I359V** is shared by both *An.*

244    *gambiae* and *An. coluzzii*, it is unlikely that this is because of an introgression event across the

245    *Cyp6m2* gene.

246    **Figure 3. Haplotype networks.**

247    Median joining network for haplotypes carrying **A13T**, **D65A**, **E328Q**, **Y347F, I359V** and **A468S,** with a maximum edge distance

248    of two SNPs. Node size indicates haplotype counts and node colour indicates the population/species of haplotypes.

249    AO=Angola; GH=; BF=Burkina Faso; CI=Côte d'Ivoire; GN=Guinea; CM=Cameroon; GW=Guinea Bissau; GM = The Gambia;

250    GA=Gabon; UG=Uganda; FR=Mayotte; GQ=Equatorial Guinea; KE=Kenya.

251

252    **Positive selection of non-synonymous alleles**

253    Extended Haplotype Heterozygosity (EHH) decay [36] was calculated  to explore evidence for

254    directional selection on the haplotypes carrying high frequency non-synonymous variants. It is

255    expected that the presence of ongoing or recent directional selection pressure would lead to the

256    increase in frequency of haplotypes, which on average will have longer regions of haplotype

257    homozygosity relative to haplotypes that are not under selection. This diversity reduction would

258    produce signatures of selection that would be conspicuous across a large genomic region. EHH

259     analysis would therefore be able to detect diversity reduction caused by ongoing directional selection

260     being driven either by amino acid substitutions identified within the gene or by mutations within *cis-*

261     acting elements next to the gene that may be under selection.

262     To perform the EHH decay analysis, we defined a core region of 1689 bases that spans across the

263     entire gene. This was identical to what was used to differentiate the identified haplotype groups

264     though hierarchical clustering. This region contained multiple distinct haplotypes above 1% frequency

265     within the cohort, including haplotypes corresponding to the C1-C5 haplotype clusters. All haplotypes

266     that did not correspond to C1-C5 were considered to be wild type (wt). Although there were several

267     different haplotypes in each population that fit this description, we do not distinguish between them

268     and call all these wild type, as *Cyp6m2* has no known resistance alleles and a true wild type remains

269     to be discovered. EHH decay was then computed for each core haplotype up to 200 kilobases

270     upstream and downstream [*Supplementary Fig. 7*]: beyond 200 kb, the EHH had decayed to zero.

271     We noted that haplotype clusters containing high frequency variants (C1-C5) did not exhibit a

272     significantly slower EHH decay relative to the wild types, showing no evidence of positive selection.

273     However, one Kenyan wild type haplotype group had a dramatically slower EHH decay relative to wild

274     type haplotypes from other populations. In order to account for this difference within wild type groups

275     across multiple populations and to reveal potential signs of selection that would be obscured by a

276     collective analysis across all populations, we separated the haplotypes by population and species and

277     recomputed EHH decay for each core haplotype as above.

278     **Figure 4. Extended haplotype homozygosity per population.**

279     No evidence for drastic difference in linkage disequilibrium within populations around core haplotypes across *Cyp6m2*.

280     Extended Haplotype Heterozygosity (EHH) decay was calculated around cluster (C1 to C5) and non-cluster (wt) haplotypes

281     using SNPs across and flanking the *Cyp6m2* region. KE=Kenya, GAgam=Gabon *An. gambiae*, AOcol= Angola *An. coluzzi*i,

282     GW=Guinea Bissau, CIcol=Côte d'Ivoire *An. coluzzii*, GHcol=Ghana *An. coluzzii*.

283     Kenyan mosquito populations are known to have an extreme demographic history, as they have

284     experienced a severe recent bottleneck, and the Angola and Gabon populations are known to be

285     geographically unique populations which are strongly differentiated from all other populations[31].

286     Hence, their haplotypes exhibited a considerably slower decay than West African haplotypes [*first*

287     *three panels: Figure 4*]. However, the putative resistance haplotypes C1-C5 did not experience a

288   slower EHH decay relative to their wild type haplotypes, showing no evidence of positive selection

289   acting upon those haplotypes in those populations.


290   As expected, the West African *An. coluzzii* haplotypes exhibited a much faster decay of EHH than

291   specimens from Kenya, Angola, or Gabon, highlighting the demographic differences previously

292   observed for these collections [31] [*last three panels, Figure 4*]. The C5 haplotype was a promising

293   candidate for potential selection as it occurred within a more diverse population, and it was interesting

294   to note that some wild type haplotypes in Côte d'Ivoire's *An. coluzzii* had a slightly slower decay than

295   others within West Africa [*fifth panel, Figure 4*]. However, these haplotypes were not part of the C5

296   cluster, and did not carry the widespread **E328Q** mutation. The C5 haplotype did not exhibit a

297   dramatically slower decay of EHH than wild type haplotypes in the populations in which it was found,

298   suggesting that it is not under positive selection.


299   **Discussion**


300   *Cyp6m2* has been implicated in many *Anopheles* populations as a key P450 that contributes to the

301   insecticide resistance phenotype [5, 14, 20, 24]. It has been reported that allelic variants across some

302   P450s can affect enzyme conformational dynamics and substrate binding affinity [28], offering

303   potential mechanisms that may modulate enzyme activity and efficiency, and thus account for

304   additional *Cyp6m2* resistance where CNVs alone may not suffice. However, little is also known about

305   *Cyp6m2* allelic variation across Africa.


306   In this study, we report a comprehensive account of the distribution of amino acid substitutions

307   occurring within the *Cyp6m2* gene. We also examine the haplotype structure of the gene to probe for

308   selective sweeps by performing hierarchical clustering of haplotypes. We also examine the genetic

309   background upon which the missense variants are found by plotting both median joining networks and

310   decay of extended haplotype homozygosity, which are useful for revealing signatures of selection. We

311   note that the distinct haplotype groups therein are stratified demographically and largely correspond to

312   signature missense variants found in specific populations. This is in contrast to the strong signals of

313   recent positive selection at other cytochrome P450 gene loci such as at *Cyp6p3* [31] which is often

314   upregulated in tandem with *Cyp6m2* in multiple pyrethroid resistant populations [5, 37, 38].

315   It is still unclear how the identified non-synonymous variants may modulate *Cyp6m2* binding activity,

316   in either the presence or absence of multiple competitive substrates and metabolites. The two

317   aromatic residues (Phe 108 and Phe 121) that have been previously identified to be vital in

318   deltamethrin orientation in the *Cyp6m2* active site[14] were not found to contain high frequency

319   variants in our dataset.

320   None of the haplotype groups identified that carried missense variants were found to be under

321   directional selection. This is despite the existence of a widespread variant (**E328Q**) linked to a

322   geographic region (West Africa) where *Cyp6m2* upregulation has been associated with emerging

323   metabolic resistance [20, 37]. In *An. coluzzii* originating from both Côte d'Ivoire and Ghana, the C5

324   haplotype that carried **E328Q** was shown to have an even faster decay of EHH than the wild type

325   haplotypes, further indicating an absence of directional selection. The stratification of other main

326   haplotype clusters from Angola (C1), Gabon (C2) and Kenya (C3) was also consistent with the strong

327   demographic differentiation and overall reduced heterozygosity of these populations described

328   elsewhere [31].

329   While the genomic data quality across the *Cyp6m2* gene and its putative promoter region was

330   satisfactory, there was a ~10,000 base region of inaccessibility upstream of *Cyp6m2* that cut across

331   the intergenic region into *Cyp6n1* [39]. A similar inaccessible region was also present 1 kb

332   downstream of the gene in the intergenic region between *Cyp6m2* and *Cyp6m3*, which is likely

333   caused by the presence of repeats that inhibit read mapping. Although it is possible that the upstream

334   region of inaccessibility could contain a regulatory variant that is susceptible to selection, it is unlikely

335   to obscure signatures of selection.

336   It has been shown in multiple studies that target-site resistance (i.e. VGSC-*kdr*) provides a strong

337   persistent baseline of resistance as it rises towards fixation within populations [40]. In the presence of

338   insecticide selection pressure, target-site mutations and metabolic resistance have also been shown

339   to act synergistically to confer a stronger resistance phenotype to pyrethroids [29, 41]. While

340   signatures of selection have previously been identified in some metabolic gene clusters within

341   populations that have a high *kdr* frequency[31], further studies need to examine whether directional

342   selection occurring on one locus can obscure selection on another locus. To resolve this conundrum,

343   genomic analysis must be performed on populations sampled across generations and whose

344    transcriptomic and phenotypic characteristics are known, in order to tease out the individual

345    contributions of specific sources of resistance.

346    Independent studies employing different experimental designs have also shown that metabolic

347    resistance manifests as a cascade of multiple upregulated genes [42]. These genes, like *Cyp6m2*, are

348    part of the normal cellular mechanism for xenobiotic detoxification that involves a linked, coordinated

349    response of large multi-gene enzyme families in complicated pathways. Therefore, it is likely that

350    identifying signatures of selection due to insecticide pressure will involve thorough analysis across this

351    vast network. The Cap 'n' Collar isoform-C (CncC) transcription factor sub-family has been shown to

352    work in tandem with other transcription factors to regulate the transcription of phase I, II and III

353    detoxification loci of multiple insects such as *Culex quinquefasciatus* and *D. melanogaster* [43, 44].

354    *CncC* knockdown or upregulation has been shown to directly affect phenotypic resistance in

355    *Anopheles gambiae* as well, modulating the expression of key P450s enzymes such as *Cyp6z2*,

356    *Cyp6z3* and *Cyp6m2* that are located in the same genomic region[43, 44]. Given that we have

357    detected no evidence of selection on amino acid variants in the *Cyp6m2* gene, it is possible that the

358    emergence of *Cyp6m2* associated resistance is being driven by selection pressures acting upon

359    genes coding for distant regulatory proteins such as transcription factors. These transcription factors

360    can regulate downstream gene expression across large genomic distances. These transcription

361    factors have also been implicated in the differential expression of other detoxification enzyme families

362    also associated with insecticide resistance (GSTs, COEs, UDP-glucuronosyltransferases (UGTs) and

363    ABC transporters). It is therefore likely that the centre of selection leading to the *Cyp6m2* associated

364    resistance phenotype will be identified through whole genome selection scans of susceptible and

365    resistant populations rather than by single loci analysis. Further research on Anopheline epigenomics,

366    transcriptomics, proteomics and systems biology will also be game changers in mapping the complex

367    regulatory network of insecticide resistance, aiding the identification of critical targets and the

368    development of new strategies to control the spread of metabolic insecticide resistance.

### Conclusion

370    The scale up of insecticide-based interventions has caused increased selection pressure and higher

371    levels of insecticide resistance across Africa. While the *CYP6M2* enzyme has been associated with

372    emerging metabolic resistance in Africa, our data indicates that allelic variation within the *Cyp6m2*

373    gene itself or across its Cyp6 supercluster has not been subject to recent positive selection in any of

374    the populations sampled. This is in contrast to other Cytochrome P450 genes where CNV alleles are

375    clearly under strong selection. Our results do not rule out a role for *Cyp6m2* in insecticide resistance

376    in natural populations, but highlight the need for a deeper understanding of the regulatory networks

377    affecting Cytochrome P450 gene expression in malaria vectors. This will require large-scale, holistic

378    experimental work that collects genomic, transcriptomic and phenotypic datasets which when

379    juxtaposed can resolve the complexities of metabolic resistance.

380    **Methods**

381    **Data collection and analysis**

382    In this study, we followed the species nomenclature of Coetzee *et al* [45] where *An. gambiae* refers to

383    *An. gambiae sensu stricto* (S form) and *An. coluzzii* refers to *An. gambiae sensu stricto* (M form). A

384    detailed description of the Ag1000G sample collection, DNA extraction, sequencing, variant calling,

385    quality control and phasing can be found here [31]. Briefly, Anopheline samples were collected from

386    33 sampling sites across 16 populations in 13 countries in sub-Saharan Africa [*Table 1 & Additional*

387    *file 1*]. The sampling procedure covered different ecosystems and aimed at collecting a minimum of

388    30 specimens per country. The specimens consisted of *An. gambiae* and *An. coluzzii*: only An.

389    coluzzii were sampled from Angola, both *An. gambiae* and *An. coluzzii* were sampled from Burkina

390    Faso, while all other populations consisted of *An gambiae,* except Kenya and Guinea Bissau where

391    the species identity was indeterminate.

392    Whole genome sequencing of all mosquitoes was performed on the Illumina HiSeq 2000 platform.

393    The generated 100 base paired-end reads were aligned to the *An. gambiae* AgamP3 reference

394    genome assembly [46] and variants were called using GATK UnifiedGenotyper. Samples with mean

395    coverage ≲14× and variants with attributes that correlated with Mendelian error in genetic crosses

396    were removed during quality control.

397    The SnpEff v4.1b software was used for the functional annotation of Ag1000G variant data [47] using

398    locations from geneset AgamP4.12. All variants in transcript AGAP008212-RA with a SnpEff

399  annotation of "missense" were regarded as nonsynonymous variants. The *Cyp6m2* gene has not

400  been shown to exhibit alternative splicing, and no alternative transcripts have been reported.

**Haplotype clustering, linkage disequilibrium and mapping of haplotype clusters**

402  To reveal the haplotype structure at *Cyp6m2, Cyp6m* sub-cluster, *Cyp6* supercluster, *HAM, ODR-2*

403  and *SH2*, we computed the Hamming distance between all haplotype pairs and performed

404  hierarchical clustering of haplotypes. We worked through arbitrary clustering threshold values to cut

405  the dendrograms at genetic distances that would best highlight the most relevant clusters. We used

406  Lewontin's $D'$ [48] to compute the linkage disequilibrium (LD) between all pairs of missense *Cyp6m2*

407  mutations. Image rendering for the haplotype clustering, linkage disequilibrium and haplotype cluster

408  frequencies map was performed using the matplotlib Python package [49].  Geography handling for

409  the haplotype cluster frequencies map was done using cartopy [50].

**Haplotype Networks**

411  We constructed haplotype networks using the median-joining algorithm [35] implemented in Python

412  [51]. Haplotypes carrying the main high frequency mutations were analysed with a maximum edge

413  distance of two SNPs. The Graphviz library was used to render the networks and the composite figure

414  was constructed in Inkscape [52].

**Extended haplotype homozygosity**

416  We defined the core haplotype on a 1689 base region spanning the *Cyp6m2*, from chromosome arm

417  3R, starting at position 6928858 and ending at position 6930547. We selected this region to ensure a

418  1:1 haplotype correspondence with that used in the hierarchical clustering analysis. We computed

419  extended haplotype homozygosity (EHH) across all core haplotypes in all populations as described in

420  Sabeti et al. [36] using scikit-allel version 1.1.9 [53]. EHH composite plots were made using the

421  matplotlib Python package [49].

**List of abbreviations**

423  CncC:  Cap 'n' Collar isoform-C

424  CNV:   Copy Number Variation

425    COEs:  Carboxylesterases

426    DDT:    Dichlorodiphenyltrichloroethane

427    EHH:    Extended Haplotype Heterozygosity

428    GSTs:   Glutathione S-Transferases

429    HAM:    Transcription Factor Hamlet

430    IRS:      Indoor Residual Spraying

431    *kdr*:      Knock-Down Resistance

432    LLINs:   Long Lasting Insecticidal Nets

433    MJNs:   Median-Joining Networks

434    ODR2:  Odd-Skipped Related

435    SH2:    SRC Homology 2

436    SNPs:   Single Nucleotide Polymorphisms

437    P450s:  Cytochrome P450 Monooxygenases

438    UGTs:   UDP-glucuronosyltransferases

439    VGSC:  Voltage Gated Sodium Channel

440    *wt*:       wild type

441

## Declarations

443    **Acknowledgements**

446

447    **Availability of data and materials**

448    Jupyter Notebooks and scripts containing all analyses, tables and figures can be found in the GitHub

449    repository [51]. Variant calls and phased haplotype data from the Ag1000G Phase 2 AR3 data release

450    were used, and can be found here [54].

451

452 **Authors contribution**

453 AM and MKNL designed the study. AM and CC developed the base code. MGW and AM performed

454 all analyses. MGW drafted the manuscript. All authors read and approved the final manuscript.

455

456 **Competing interests statement**

457 The authors declare no competing interests.

458

459 **Consent for publication**

460 Not applicable

461

462 **Ethics approval and consent to participate**

463 Not applicable.

464

465 **Funding**

466 The Wellcome Sanger Institute is funded by the Wellcome Trust (grant 206194/Z/17/Z), which

467 supports M.K.N.L. and part of the sequencing, analysis, informatics, and management of the

468 *Anopheles gambiae* 1000 Genomes Project.

469

470

471 **Supplementary Figures**

472 **Supplementary Figure 1. Linkage disequilibrium (*D*′) between non-synonymous variants**.

473

474 A value of 1 shows perfect linkage between the alleles. A value of −1 shows that the alleles are never found conjointly. The bar

475 plot indicates allele frequencies within the Ag1000G phase 2 cohort.

476

477 **Supplementary Figure 2. Hierarchical clustering and missense mutations for *ODR2*.**

478 Top: a dendrogram showing hierarchical clustering of haplotypes across the ODR2 gene. The gene is located at position

479 7,059,422 to 7,119,244: 128,875 bases downstream of *Cyp6m2*.

480 The colour bar indicates the population of origin for each haplotype.

481 Bottom: high frequency (>◻5%) alleles identified within each haplotype (white = reference allele; black = alternative allele).

482

483 **Supplementary Figure 3. Hierarchical clustering and missense mutations for *HAM.***

484 Top: a dendrogram showing hierarchical clustering of haplotypes across the HAM gene. The gene is located at position

485 7,435,306 to 7,485,012: 504,759 bases downstream of *Cyp6m2*.

486 The colour bar indicates the population of origin for each haplotype.

487 Bottom: high frequency (>◻5%) alleles identified within each haplotype (white = reference allele; black = alternative allele).

488

489 **Supplementary Figure 4. Hierarchical clustering and missense mutations for *SH2.***

490 Top: a dendrogram showing hierarchical clustering of haplotypes across the SH2 gene. The gene is located at position

491 8,176,778 to 8,183,084: 1,246,231 bases downstream of *Cyp6m2*.

492 The colour bar indicates the population of origin for each haplotype.

493 Bottom: high frequency (>◻5%) alleles identified within each haplotype (white = reference allele; black = alternative allele).

494

495 **Supplementary Figure 5. Hierarchical clustering and missense mutations for *Cyp6m* sub**

496 **cluster**.

497   Top: a dendrogram showing hierarchical clustering of haplotypes across the Cyp6m sub cluster of genes containing Cyp6m2,

498   Cyp6m3 and Cyp6m4. The genes are located at position 6928858 to 6935721.

499   The colour bar indicates the population of origin for each haplotype.

500   Bottom: high frequency (>5%) alleles identified within each haplotype (white = reference allele; black = alternative allele).

501

502   **Supplementary Figure 6. Hierarchical clustering and missense mutations for *Cyp6***

503   **supercluster.**

504   Top: a dendrogram showing hierarchical clustering of haplotypes across the Cyp6 supercluster of 14 P450 genes containing

505   *Cyp6s2, Cyp6s1, Cyp6r1, Cyp6n2, Cyp6y2, Cyp6y1, Cyp6m1, Cyp6n1, Cyp6m2, Cyp6m3, Cyp6m4, Cyp6z3, Cyp6z2* and

506   *Cyp6z1*. The genes are located at position 6903106 to 6978142.

507   The colour bar indicates the population of origin for each haplotype.

508   Bottom: high frequency (>70%) alleles identified within each haplotype (white = reference allele; black = alternative allele).

509

510   **Supplementary Figure 7. Extended haplotype homozygosity across all populations.**

511   A rapid decay of EHH in comparison to other haplotypes implies absence of positive selection.

512

# References

514   1.   WHO: **World Malaria Report 2019**. 2019.

515   2.   Bhatt S, Weiss DJ, Cameron E, Bisanzio D, Mappin B, Dalrymple U, Battle K, Moyes
516        CL, Henry A, Eckhoff PA *et al*: **The effect of malaria control on Plasmodium**
517        **falciparum in Africa between 2000 and 2015**. *Nature* 2015, **526**(7572):207-211.

518   3.   Cibulskis RE, Alonso P, Aponte J, Aregawi M, Barrette A, Bergeron L, Fergus CA,
519        Knox T, Lynch M, Patouillard E *et al*: **Malaria: Global progress 2000 - 2015 and**
520        **future challenges**. *Infect Dis Poverty* 2016, **5**(1):61.

521   4.   Ranson H, Lissenden N: **Insecticide Resistance in African Anopheles**
522        **Mosquitoes: A Worsening Situation that Needs Urgent Action to Maintain**
523        **Malaria Control**. *Trends Parasitol* 2016, **32**(3):187-196.

524   5.   Djouaka RF, Bakare AA, Coulibaly ON, Akogbeto MC, Ranson H, Hemingway J,
525        Strode C: **Expression of the cytochrome P450s, CYP6P3 and CYP6M2 are**
526        **significantly elevated in multiple pyrethroid resistant populations of Anopheles**
527        **gambiae s.s. from Southern Benin and Nigeria**. *BMC Genomics* 2008, **9**:538.

528   6.   Djouaka R, Riveron JM, Yessoufou A, Tchigossou G, Akoton R, Irving H, Djegbe I,
529        Moutairou K, Adeoti R, Tamò M *et al*: **Multiple insecticide resistance in an**
530        **infected population of the malaria vector Anopheles funestus in Benin**. *Parasit*
531        *Vectors* 2016, **9**:453.

532   7.   WHO: **Global report on insecticide resistance in malaria vectors: 2010--2016**.
533        2018.

534  8.    **WHO Malaria Threats Map**
535        [https://apps.who.int/malaria/maps/threats/?theme=prevention&mapType=prevention
536        %3A0&bounds=%5B%5B-54.61667525407141%2C-
537        26.993804332606665%5D%2C%5B66.07511128112793%2C35.549094294064915
538        %5D%5D&insecticideClass=PYRETHROIDS&insecticideTypes=&assayTypes=MOL
539        ECULAR_ASSAY%2CBIOCHEMICAL_ASSAY%2CSYNERGIST-
540        INSECTICIDE_BIOASSAY&synergistTypes=&species=&vectorSpecies=&surveyTyp
541        es=&deletionType=HRP2_PROPORTION_DELETION&plasmodiumSpecies=P._FAL
542        CIPARUM&drug=DRUG_AL&mmType=1&endemicity=false&countryMode=false&sto
543        ryMode=false&storyModeStep=0&filterOpen=false&filtersMode=filters&years=2010%
544        2C2018]
545  9.    Ranson H, N'Guessan R, Lines J, Moiroux N, Nkuni Z, Corbel V: **Pyrethroid**
546        **resistance in African anopheline mosquitoes: what are the implications for**
547        **malaria control?** *Trends Parasitol* 2011, **27**(2):91-98.
548  10.   Wilding CS, Weetman D, Steen K, Donnelly MJ: **High, clustered, nucleotide**
549        **diversity in the genome of Anopheles gambiae revealed through pooled-**
550        **template sequencing: implications for high-throughput genotyping protocols**.
551        *BMC Genomics* 2009, **10**:320.
552  11.   Hargreaves K, Koekemoer LL, Brooke BD, Hunt RH, Mthembu J, Coetzee M:
553        **Anopheles funestus resistant to pyrethroid insecticides in South Africa**. *Med*
554        *Vet Entomol* 2000, **14**(2):181-189.
555  12.   Wondji CS, Morgan J, Coetzee M, Hunt RH, Steen K, Black WCt, Hemingway J,
556        Ranson H: **Mapping a quantitative trait locus (QTL) conferring pyrethroid**
557        **resistance in the African malaria vector Anopheles funestus**. *BMC Genomics*
558        2007, **8**:34.
559  13.   Wondji CS, Irving H, Morgan J, Lobo NF, Collins FH, Hunt RH, Coetzee M,
560        Hemingway J, Ranson H: **Two duplicated P450 genes are associated with**
561        **pyrethroid resistance in Anopheles funestus, a major malaria vector**. *Genome*
562        *Res* 2009, **19**(3):452-459.
563  14.   Stevenson BJ, Bibby J, Pignatelli P, Muangnoicharoen S, O'Neill PM, Lian L-Y,
564        Müller P, Nikou D, Steven A, Hemingway J *et al*: **Cytochrome P450 6M2 from the**
565        **malaria vector Anopheles gambiae metabolizes pyrethroids: Sequential**
566        **metabolism of deltamethrin revealed**. *Insect Biochem Mol Biol* 2011, **41**(7):492-
567        502.
568  15.   **Chromosome 3R: 6,928,825-6,930,580 - Region in detail - Anopheles gambiae -**
569        **VectorBase**
570        [https://www.vectorbase.org/Anopheles_gambiae/Location/View?db=core;g=AGAP00
571        8212;r=3R:6928825-6930580;t=AGAP008212-RA]
572  16.   Ranson H, Claudianos C, Ortelli F, Abgrall C, Hemingway J, Sharakhova MV, Unger
573        MF, Collins FH, Feyereisen R: **Evolution of supergene families associated with**
574        **insecticide resistance**. *Science* 2002, **298**(5591):179-181.
575  17.   Ranson H, Paton MG, Jensen B, McCarroll L, Vaughan A, Hogan JR, Hemingway J,
576        Collins FH: **Genetic mapping of genes conferring permethrin resistance in the**
577        **malaria vector, Anopheles gambiae**. *Insect Mol Biol* 2004, **13**(4):379-386.
578  18.   Müller P, Donnelly MJ, Ranson H: **Transcription profiling of a recently colonised**
579        **pyrethroid resistant Anopheles gambiae strain from Ghana**. *BMC Genomics*
580        2007, **8**:36.
581  19.   Nardini L, Christian RN, Coetzer N, Ranson H, Coetzee M, Koekemoer LL:
582        **Detoxification enzymes associated with insecticide resistance in laboratory**
583        **strains of Anopheles arabiensis of different geographic origin**. *Parasit Vectors*
584        2012, **5**:113.
585  20.   Edi CV, Djogbénou L, Jenkins AM, Regna K, Muskavitch MAT, Poupardin R, Jones
586        CM, Essandoh J, Kétoh GK, Paine MJI *et al*: **CYP6 P450 enzymes and ACE-1**
587        **duplication produce extreme and multiple insecticide resistance in the malaria**
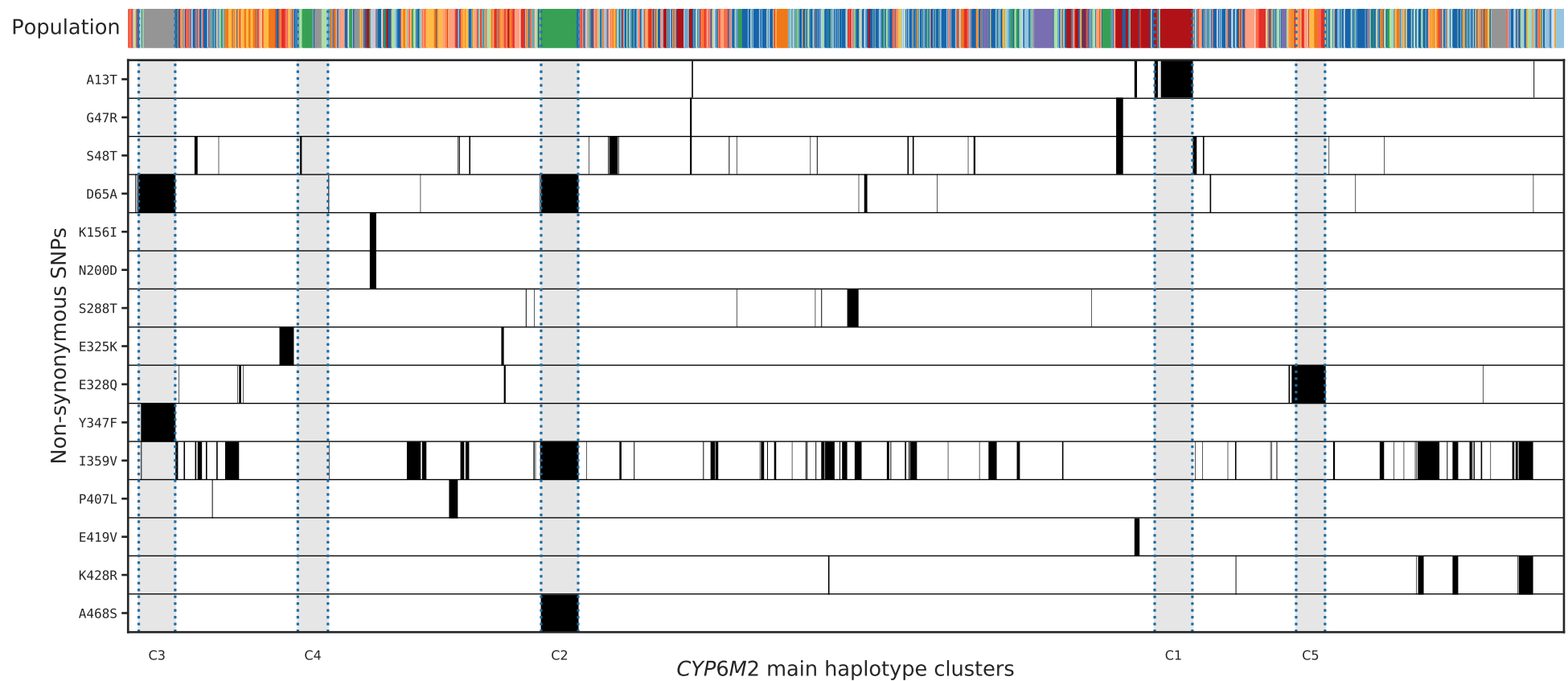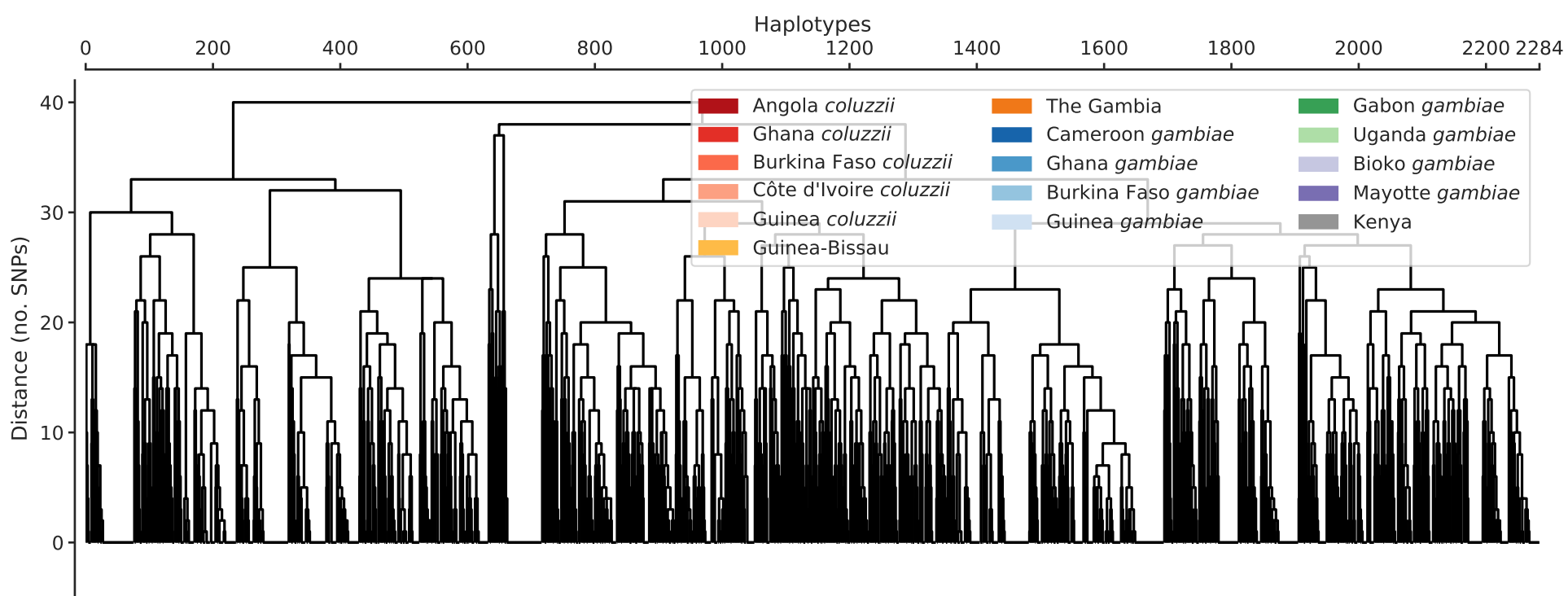588        **mosquito Anopheles gambiae**. *PLoS Genet* 2014, **10**(3):e1004236.

589   21.   Yan Z-W, He Z-B, Yan Z-T, Si F-L, Zhou Y, Chen B: **Genome-wide and**
590         **expression-profiling analyses suggest the main cytochrome P450 genes**
591         **related to pyrethroid resistance in the malaria vector, Anopheles sinensis**
592         **(Diptera Culicidae)**. *Pest Manag Sci* 2018, **74**(8):1810-1820.
593   22.   Djègbè I, Agossa FR, Jones CM, Poupardin R, Cornelie S, Akogbéto M, Ranson H,
594         Corbel V: **Molecular characterization of DDT resistance in Anopheles gambiae**
595         **from Benin**. *Parasit Vectors* 2014, **7**:409.
596   23.   Adolfi A, Poulton B, Anthousi A, Macilwee S, Ranson H, Lycett GJ: **Functional**
597         **genetic validation of key genes conferring insecticide resistance in the major**
598         **African malaria vector, Anopheles gambiae**. *Proc Natl Acad Sci U S A* 2019,
599         **116**(51):25764-25772.
600   24.   Mitchell SN, Stevenson BJ, Müller P, Wilding CS, Egyir-Yawson A, Field SG,
601         Hemingway J, Paine MJI, Ranson H, Donnelly MJ: **Identification and validation of**
602         **a gene causing cross-resistance between insecticide classes in Anopheles**
603         **gambiae from Ghana**. *Proc Natl Acad Sci U S A* 2012, **109**(16):6147-6152.
604   25.   Lucas ER, Miles A, Harding NJ, Clarkson CS, Lawniczak MKN, Kwiatkowski DP,
605         Weetman D, Donnelly MJ, Anopheles gambiae Genomes C: **Whole-genome**
606         **sequencing reveals high complexity of copy number variation at insecticide**
607         **resistance loci in malaria mosquitoes**. *Genome Res* 2019, **29**(8):1250-1261.
608   26.   Weetman D, Djogbenou LS, Lucas E: **Copy number variation (CNV) and**
609         **insecticide resistance in mosquitoes: evolving knowledge or an evolving**
610         **problem?** *Curr Opin Insect Sci* 2018, **27**:82-88.
611   27.   Schuler MA, Berenbaum MR: **Structure and function of cytochrome P450S in**
612         **insect adaptation to natural and synthetic toxins: insights gained from**
613         **molecular modeling**. *J Chem Ecol* 2013, **39**(9):1232-1245.
614   28.   Ibrahim SS, Riveron JM, Bibby J, Irving H, Yunta C, Paine MJI, Wondji CS: **Allelic**
615         **Variation of Cytochrome P450s Drives Resistance to Bednet Insecticides in a**
616         **Major Malaria Vector**. *PLoS Genet* 2015, **11**(10):e1005618.
617   29.   Weetman D, Wilding CS, Neafsey DE, Müller P, Ochomo E, Isaacs AT, Steen K,
618         Rippon EJ, Morgan JC, Mawejje HD *et al*: **Candidate-gene based GWAS identifies**
619         **reproducible DNA markers for metabolic pyrethroid resistance from standing**
620         **genetic variation in East African Anopheles gambiae**. *Sci Rep* 2018, **8**(1):2920.
621   30.   Clarkson CS, Miles A, Harding NJ, Lucas ER, Battey CJ, Amaya-Romero JE, Cano
622         J, Diabate A, Constant E, Nwakanma DC *et al*: **Genome variation and population**
623         **structure among 1,142 mosquitoes of the African malaria vector species**
624         **<em>Anopheles gambiae</em> and <em>Anopheles coluzzii</em>**. *bioRxiv*
625         2019:864314.
626   31.   Consortium TAgG: **Genetic diversity of the African malaria vector Anopheles**
627         **gambiae**. *Nature* 2017, **552**(7683):96-100.
628   32.   **Chromosome 3R: 7,059,422 - 7,119,244 - Region in detail - Anopheles gambiae -**
629         **VectorBase** [https://vectorbase.org/vectorbase/app/record/gene/AGAP008222]
630   33.   **Chromosome 3R:7,435,306 - 7,485,012 - Anopheles gambiae - VectorBase**
631         [https://vectorbase.org/vectorbase/app/record/gene/AGAP008232]
632   34.   **Chromosome 3R: 8,176,778 - 8,183,084 - Region in detail - Anopheles gambiae -**
633         **VectorBase**
634    [https://vectorbase.org/vectorbase/app/record/gene/AGAP008273]
635   35.   Bandelt HJ, Forster P, Röhl A: **Median-joining networks for inferring intraspecific**
636         **phylogenies**. *Mol Biol Evol* 1999, **16**(1):37-48.
637   36.   Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, Schaffner SF, Gabriel
638         SB, Platko JV, Patterson NJ, McDonald GJ *et al*: **Detecting recent positive**
639         **selection in the human genome from haplotype structure**. *Nature* 2002,
640         **419**(6909):832-837.
641   37.   Müller P, Warr E, Stevenson BJ, Pignatelli PM, Morgan JC, Steven A, Yawson AE,
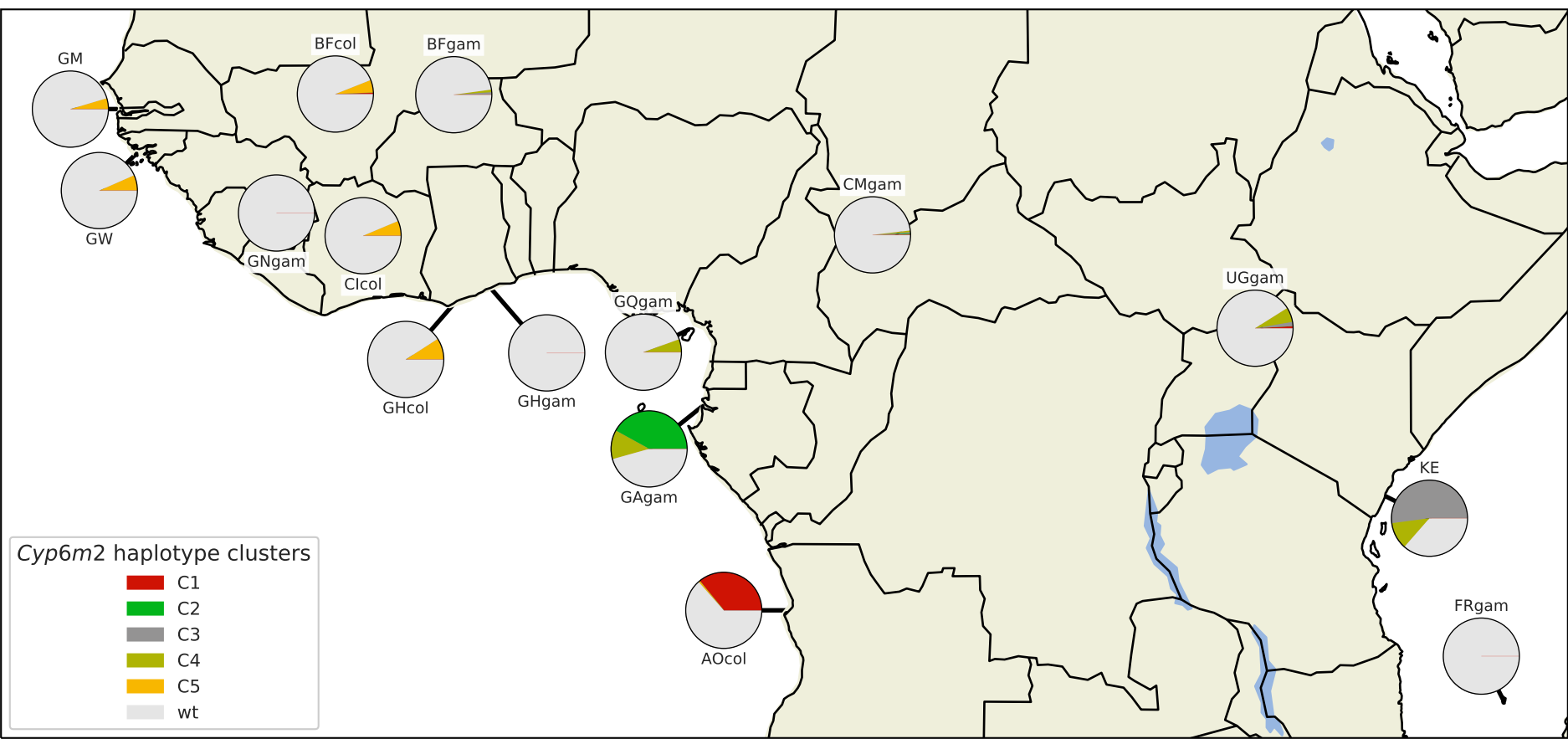642         Mitchell SN, Ranson H, Hemingway J *et al*: **Field-caught permethrin-resistant**

643    **Anopheles gambiae overexpress CYP6P3, a P450 that metabolises pyrethroids**.
644    *PLoS Genet* 2008, **4**(11):e1000286.

645  38.  Stica C, Jeffries CL, Irish SR, Barry Y, Camara D, Yansane I, Kristan M, Walker T,
646    Messenger LA: **Characterizing the molecular and metabolic mechanisms of**
647    **insecticide resistance in Anopheles gambiae in Faranah, Guinea**. *Malar J* 2019,
648    **18**(1):244.

649  39.  **Ag1000G - AR3 Panoptes genome browser**
650    [https://www.malariagen.net/apps/ag1000g/phase1-
651    AR3/index.html?dataset=Ag1000G&workspace=workspace_1&view=f6c6c7c8-23c9-
652    11eb-a4f3-22000a6287ed&state=genomebrowser]

653  40.  Clarkson CS, Miles A, Harding NJ, Weetman D, Kwiatkowski D, Donnelly M, The
654    Anopheles gambiae Genomes C: **The genetic architecture of target-site**
655    **resistance to pyrethroid insecticides in the African malaria vectors Anopheles**
656    **gambiae and Anopheles coluzzii**. 2018.

657  41.  Hemingway J: **The role of vector control in stopping the transmission of**
658    **malaria: threats and opportunities**. *Philos Trans R Soc Lond B Biol Sci* 2014,
659    **369**(1645):20130431.

660  42.  Liu N: **Insecticide resistance in mosquitoes: impact, mechanisms, and research**
661    **directions**. *Annu Rev Entomol* 2015, **60**:537-559.

662  43.  Ingham VA, Pignatelli P, Moore JD, Wagstaff S, Ranson H: **The transcription factor**
663    **Maf-S regulates metabolic resistance to insecticides in the malaria vector**
664    **Anopheles gambiae**. *BMC Genomics* 2017, **18**(1):669.

665  44.  Wilding CS: **Regulating resistance: CncC:Maf, antioxidant response elements**
666    **and the overexpression of detoxification genes in insecticide resistance**. *Curr*
667    *Opin Insect Sci* 2018, **27**:89-96.

668  45.  Coetzee M, Hunt RH, Wilkerson R, Della Torre A, Coulibaly MB, Besansky NJ:
669    **Anopheles coluzzii and Anopheles amharicus, new members of the Anopheles**
670    **gambiae complex**. *Zootaxa* 2013, **3619**:246-274.

671  46.  Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, Nusskern DR,
672    Wincker P, Clark AG, Ribeiro JMC, Wides R *et al*: **The genome sequence of the**
673    **malaria mosquito Anopheles gambiae**. *Science* 2002, **298**(5591):129-149.

674  47.  Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden
675    DM: **A program for annotating and predicting the effects of single nucleotide**
676    **polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster**
677    **strain w1118; iso-2; iso-3**. *Fly* 2012, **6**(2):80-92.

678  48.  Lewontin RC: **The Interaction of Selection and Linkage. I. General**
679    **Considerations; Heterotic Models**. *Genetics* 1964, **49**(1):49-67.

680  49.  Hunter JD: **Matplotlib: A 2D Graphics Environment**. *Comput Sci Eng* 2007,
681    **9**(3):90-95.

682  50.  Cartopy: **Using cartopy with matplotlib — cartopy 0.18.0 documentation**. In.,
683    0.17.0 edn. https://scitools.org.uk/; 2020.

684  51.  Wagah MG: **ag1000g-phase2-cyp6m2**. In., 9/11/2020 edn. https://github.com/;
685    2020.

686  52.  Harrington B: **Inkscape**. In., 1.0.1 edn; 2005.

687  53.  Miles A: **scikit-allel - Explore and analyse genetic variation — scikit-allel 1.3.2**
688    **documentation**. In. https://github.com; 2018.

689  54.  Consortium TAgG: **Ag1000G phase 2 AR1 data release**. In., 1 edn. MalariaGen
690    Genomic Epidemiology Network; 2017.

691

692  **Additional files**

693    1.  Additional file 1.

694         a.   File name = Additional file 1

695         b.   Title = List of *An. gambiae* and *An. coluzzii* genome samples and haplotypes from

696            Ag1000G Phase 2-AR3.

697         c.   Format = csv

698         d.   Description = Table showing Ag1000G Phase 2-AR3 sample properties such as

699            population, country, region, sex, species identity and haplotype cluster.

700   2.   Additional file 2.

701         a.   File name = Additional file 2

702         b.   Title = List of synonymous and non-synonymous genetic variants in *Cyp6m2*.

703         c.   Format = csv

704         d.   Description = Table showing Ag1000G Phase 2-AR3 *Cyp6m2* variant calls and

705            variant properties stratified by population and effect.

706

707

Cyp6m2 haplotype clusters

- C1
- C2
- C3
- C4
- C5
- wt

KE

Core haplotype
— C3
— C4
— *wt*

GAgam

Core haplotype
— C2
— C4
— *wt*

AOcol

Core haplotype
— C1
— *wt*

GW

Core haplotype
— C5
— *wt*

CIcol

Core haplotype
— C5
— *wt*

GHcol

Core haplotype
— C5
— *wt*

*Cyp6m2*

Genes +

Chromosome 3R position (Mbp)