

1 **The evolution and spread of target-site**
2 **resistance to pyrethroid insecticides in the**
3 **African malaria vectors *Anopheles gambiae***
4 **and *Anopheles coluzzii***

5 Chris S. Clarkson¹, Alistair Miles^{2,1}, Nicholas J. Harding², Dominic
6 Kwiatkowski^{1,2}, Martin Donnelly^{3,1}, and The *Anopheles gambiae*
7 1000 Genomes Consortium⁴

8 ¹Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA

9 ²Big Data Institute, Old Road, Oxford OX3 7FZ

10 ³Liverpool School of Tropical Medicine, Pembroke Place, Liverpool
11 L3 5QA

12 ⁴Ag1000g Consortium, MalariaGEN, Big Data Institute, Old Road,
13 Oxford OX3 7FZ

14 Work in progress

15 **Abstract**

16 Resistance to pyrethroid insecticides is a major concern for malaria vector control,
17 because these are the only compounds approved for use in insecticide-treated bed-nets
18 (ITNs). Pyrethroids target the voltage-gated sodium channel (VGSC), an essential

component of the mosquito nervous system, but substitutions in the amino acid sequence can disrupt the activity of these insecticides, inducing a resistance phenotype. Here we use Illumina whole-genome sequence data from phase 1 of the *Anopheles gambiae* 1000 Genomes Project (Ag1000G) to provide a comprehensive account of genetic variation at the *Vgsc* locus in mosquito populations from 8 African countries. In addition to three known resistance variants, we describe 18 non-synonymous variants at appreciable frequency in one or more populations that are previously unknown in mosquitoes. For each variant we predict a resistance phenotype based on genetic evidence for positive selection, patterns of linkage between variants, and functional evidence from other species. We then analyse the genetic backgrounds on which resistance variants are found, to refine our understanding of the origins and spread of resistance between species and geographical locations. We identify ten distinct outbreaks of resistance, of which five appear to be localised to a single geographical location, and five have spread between two or more countries. The most successful and widespread outbreak (F1) originates in West Africa and has subsequently spread to countries in Central and Southern Africa. Our results demonstrate that the molecular basis of pyrethroid resistance in African malaria vectors is more complex than previously appreciated, and provide a foundation for the design of new genetic tools for outbreak surveillance to inform insecticide resistance management and track the further spread of resistance

Introduction

An estimated 663 million cases of malaria were averted in Africa between 2000 and 2015 due to public health interventions, of which 68% were prevented by insecticide-treated bed-nets (ITNs) and 10% through indoor residual spraying of insecticides (IRS) [1]. However, over this same period, insecticide resistance has become increasingly prevalent in malaria vector populations [2, 3]. Four chemical classes of insecticides – organophosphates, carbamates, pyrethroids and organochlorines – are licensed for use in public health, but only pyrethroids are approved by the World Health Organisation (WHO) for use in ITNs. Pyrethroids are also commonly used for IRS and in agriculture, and mosquito populations are under pressure to evolve molecular mechanisms of pyrethroid resistance. There is evidence that pyrethroid resistance has a direct impact on the effectiveness of ITNs and

IRS [4, 5], although assessing the impact on disease prevalence is difficult and has been hampered by the fact that pyrethroid resistance is now so pervasive that it is difficult to find mosquito populations with pyrethroid susceptibility at baseline to serve as controls [6]. Nevertheless, the position of WHO remains that insecticide resistance poses a grave threat to the substantial gains made in reducing malaria across Africa [7, 6]. Improvements are needed in our ability to monitor resistance, and gaps must be filled in our knowledge of the molecular basis of resistance.

The voltage-gated sodium channel (VGSC) is the physiological target of pyrethroids and of the organochlorine DDT and is integral to the insect nervous system, involved in the transmission of nerve impulses. Both pyrethroids and DDT have a similar mode of action, binding to sites within the protein channel and preventing normal nerve function, causing paralysis (“knock-down”) and then death. However, amino acid substitutions at key positions within the channel can alter the interaction with the insecticide molecule, thereby substantially increasing the dosage of insecticide required for knock-down (hence described as knock-down resistance or *kdr* [8]. If this tolerance exceeds the dosage present in ITNs or on indoor surfaces following IRS, these interventions may be rendered ineffective [4, 5]. In the African malaria vectors *Anopheles gambiae* and *An. coluzzii*, three substitutions have been found in natural populations and shown experimentally to cause pyrethroid and DDT resistance. Two of these substitutions occur in codon 995, with the Leucine → Phenylalanine (L995F) substitution prevalent in West and Central Africa [9, 3], and the Leucine → Serine (L995S) substitution found in Central and East Africa [10, 3]. A third variant N1570Y has been found in association with L995F in Central Africa and shown to increase resistance above L995F alone [11]. Codon numbering is given here relative to transcript AGAP004707-RA as defined in the AgamP4.4 gene annotations. A mapping of codon numbers from AGAP004707-RA to *Musca domestica*, the system in which the *kdr* mutations were first discovered [12], is given in Table 1 and in @@Supplementary data.

Target-site resistance to pyrethroids and DDT has also been studied in a range of other insect species, including disease vectors as well as domestic and crop pests. Because of its essential function, the VGSC protein is highly conserved across insect species [13], and therefore knowledge gained from one species is relevant to another. Many resistance-associated variants have been described in these other species, and thus there are many

possible amino acid substitutions that could induce a resistance phenotype in malaria vectors other than the known variants in codons 995 and 1570 [14],[15]. Some of these variants fall within transmembrane domains, and thus may directly interact with insecticide molecules [16, 13]. However, functional studies have also demonstrated that variants within internal linker domains can substantially enhance the the level of resistance, when present in combination with channel modifications [11]. Most previous studies of *An. gambiae* and/or *An. coluzzii* have performed targeted sequencing of small regions within the gene [17, 18, 19], and there has been no comprehensive survey of variation across the entire gene in multiple mosquito populations.

Insecticide resistance monitoring in malaria vector populations now often incorporates some form of genetic assay to detect the allele present at *Vgsc* codon 995 (e.g. [20]). Both alleles are present at high frequency in multiple geographical locations, and the L995F allele is present in both *An. gambiae* and *An. coluzzii* [3]. The extent of mosquito migration remains an open question, however mosquitoes do travel between different locations and have the potential to spread resistance alleles from one population to another (adaptive gene flow) [21]. Hybridization between mosquito species also occurs and has the potential to transfer resistance alleles between species (adaptive introgression); studies in West Africa have shown that the L995F allele has been transferred from *An. gambiae* into *An. coluzzii* populations [22]. A resistance allele may also arise independently in multiple populations in the absence of gene flow [23], either because of multiple mutational events occurring after insecticides are introduced (selection on new mutations), or because resistance alleles were already present at low frequency in mosquito populations prior to insecticide use (selection on standing variation). Previous studies have found evidence that the L995F allele occurs on several different genetic backgrounds, suggesting multiple origins of resistance [17, 18]. However, these studies have used information from only a small region of the gene, and have limited resolution to make inferences about geographical origins or history of spread. Better information about the origins and spread of resistance could improve insecticide resistance monitoring and inform strategies for insecticide resistance management.

Here we report an in-depth analysis of the *Vgsc* locus using genotype and haplotype data from phase 1 of the *Anopheles gambiae* 1000 Genomes Project (Ag1000G) [24]. These

data are derived from whole-genome Illumina sequencing of 765 individual mosquitoes collected from natural populations in 8 African countries. We presented some initial results regarding *Vgsc* from analyses of these data in Ag1000g Consortium *et al.* [24], and here we describe a number of new analyses to confirm, extend and elaborate on our initial findings. Our aim is to provide a comprehensive account of genetic variation at the *Vgsc* locus and its implications for the management of pyrethroid resistance in natural vector populations.

Results

Functional variation

To identify single nucleotide polymorphisms (SNPs) with a potentially functional role in pyrethroid resistance, we extracted SNPs from the Ag1000G phase 1 data resource that alter the amino acid sequence of the VGSC protein, and computed their allele frequencies among 9 populations defined by species and country of origin. SNPs that confer resistance are expected to increase in frequency under selective pressure, and we refined the list of potentially functional SNPs to retain only those at an appreciable frequency ($>5\%$) in one or more populations (Table 1). The resulting list comprises 20 SNPs, including the known L995F, L995S and N1570Y variants, and a further 17 SNPs not previously described in these species. We reported 15 of these novel SNPs in our initial analysis of the Ag1000G phase 1 data (@@REF Ag1000G), and we extend the analyses here to incorporate two tri-allelic SNPs affecting codons 402 and 410.

The two alleles in codon 995 are clearly the main drivers of resistance at this locus. The L995F allele at high frequency in populations of both species from West, Central and Southern Africa, and the L995S allele at high frequency among *An. gambiae* populations from Central and East Africa (Table 1; @@REF Ag1000G). All haplotypes carrying L995F or L995S have evidence for strong recent positive selection (@@REF Ag1000G). Both alleles were present in populations sampled from Cameroon and Gabon, including some individuals with a hybrid L995F/S genotype. Within these populations, the L995F and L995S alleles were (@@TODO were not?) in Hardy-Weinberg equilibrium ($P=@@$), thus there does not (@@does?) appear to be selection against hybrids.

Table 1. Non-synonymous nucleotide variation in the voltage-gated sodium channel gene. AO=Angola; BF=Burkina Faso; GN=Guinea; CM=Cameroon; GA=Gabon; UG=Uganda; KE=Kenya; GW=Guinea-Bissau; *Ac*=*An. coluzzii*; *Ag*=*An. gambiae*. All variants are at 5% frequency or above in one or more of the 9 Ag1000G phase 1 populations, with the exception of 2,400,071 G>T which is only found in the CMAg population at 0.4% frequency but is included because another mutation (2,400,071 G>A) is found at the same position causing the same amino acid substitution (M490I); and 2,431,019 T>C (F1920S) which is at 4% frequency in GA*Ag* but also found in CMAg and linked to L995F.

Variant			Population allele frequency (%)										Function	
Position ¹	<i>Ag</i> ²	<i>Md</i> ³	AOAc	BFAC	GNAg	BFAG	CMAg	GAAg	UGAg	KE	GW	Domain ⁴	Resistance phenotype ⁵	
2,390,177 G>A	R254K	R261	0	0	0	0	32	21	0	0	0	IN (I.S4-I.S5)	L995F enhancer (predicted)	
2,391,228 G>C	V402L	V410	0	7	0	0	0	0	0	0	0	TM (I.S6)	I1527T enhancer (predicted)	
2,391,228 G>T	V402L	V410	0	7	0	0	0	0	0	0	0	TM (I.S6)	I1527T enhancer (predicted)	
2,399,997 G>C	D466H	-	0	0	0	0	7	0	0	0	0	IN (I.S6-II.S1)	L995F enhancer (predicted)	
2,400,071 G>A	M490I	M508	0	0	0	0	0	0	0	18	0	IN (I.S6-II.S1)	none (predicted)	
2,400,071 G>T	M490I	M508	0	0	0	0	0	0	0	0	0	IN (I.S6-II.S1)	none (predicted)	
2,416,980 C>T	T791M	T810	0	1	13	14	0	0	0	0	0	TM (II.S1)	L995F enhancer (predicted)	
2,422,651 T>C	L995S	L1014	0	0	0	0	15	64	100	76	0	TM (II.S6)	driver	
2,422,652 A>T	L995F	L1014	86	85	100	100	53	36	0	0	0	TM (II.S6)	driver	
2,424,384 C>T	A1125V	K1133	9	0	0	0	0	0	0	0	0	IN (II.S6-III.S1)	none (predicted)	
2,425,077 G>A	V1254I	I1262	0	0	0	0	0	0	0	0	5	IN (II.S6-III.S1)	none (predicted)	
2,429,617 T>C	I1527T	I1532	0	14	0	0	0	0	0	0	0	TM (III.S6)	driver (predicted)	
2,429,745 A>T*	N1570Y	N1575	0	26	10	22	6	0	0	0	0	IN (III.S6-IV.S1)	L995F enhancer	
2,429,897 A>G	E1597G	E1602	0	0	6	4	0	0	0	0	0	IN (III.S6-IV.S1)	L995F enhancer (predicted)	
2,429,915 A>C	K1603T	K1608	0	5	0	0	0	0	0	0	0	TM (IV.S1)	L995F enhancer (predicted)	
2,430,424 G>T	A1746S	A1751	0	0	11	13	0	0	0	0	0	TM (IV.S5)	L995F enhancer (predicted)	
2,430,817 G>A	V1853I	V1858	0	0	8	5	0	0	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	
2,430,863 T>C	I1868T	I1873	0	0	18	25	0	0	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	
2,430,880 C>T	P1874S	P1879	0	21	0	0	0	0	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	
2,430,881 C>T	P1874L	P1879	0	7	45	26	0	0	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	
2,431,019 T>C	F1920S	Y1925	0	0	0	0	1	4	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	
2,431,061 C>T	A1934V	A1939	0	12	0	0	0	0	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	
2,431,079 T>C	I1940T	I1945	0	4	0	0	7	0	0	0	0	IN (IV.S6-)	L995F enhancer (predicted)	

¹ Position relative to the AgamP3 reference sequence, chromosome arm 2L. Variants marked with an asterisk (*) failed conservative variant filters applied genome-wide in the Ag1000G phase 1 AR3 callset, but appeared sound on manual inspection of read alignments.

² Codon numbering according to *Anopheles gambiae* transcript AGAP004707-RA in geneset AgamP4.4.

³ Codon numbering according to *Musca domestica* EMBL accession X96668 [12].

⁴ Position of the variant within the protein. IN=internal domain; TM=trans-membrane domain. The protein contains four homologous repeats (I-IV), each having six transmembrane segments (1-6). Codes in parentheses identify the specific domain, e.g., “I.S4” refers to trans-membrane segment 4 in repeat I, and “IS4-IS5” refers to the linker segment between I.S4 and I.S5.

⁵ Phenotype predictions are based on population genetic evidence and have not been confirmed experimentally.

141 The I1527T allele is present in *An. coluzzii* from Burkina Faso at 14% frequency, and
 142 there is evidence that haplotypes carrying this allele have been positively selected (@@REF
 143 Ag1000G). Codon 1527 occurs within trans-membrane domain segment III.S6, immedi-
 144 ately adjacent to a second predicted binding pocket for pyrethroid molecules, thus it is
 145 plausible that I1527T could alter insecticide binding (@@REF Dong). We also found that
 146 the two variant alleles affecting codon 402, both of which induce a V402L substitution,
 147 were in strong linkage with I1527T (D'>@@N; Figure 1), and almost all haplotypes car-
 148 rying I1527T also carried a V402L substitution. The most parsimonious explanation for
 149 this pattern of linkage is that the I1527T mutation occurred first, and mutations in codon
 150 402 subsequently arose on this genetic background. Codon 402 also occurs within a trans-
 151 membrane segment (I.S6), and the V402L substitution has by itself been shown experi-
 152 mentally to increase pyrethroid resistance in @@species and *Xenopus* oocytes (@@REFs).
 153 However, because V402L appears secondary to I1527T in our cohort, we classify I1527T
 154 as a putative resistance driver and V402L as a putative enhancer. Because of the limited
 155 geographical distribution of these alleles, we hypothesize that the I1527T+V402L combi-
 156 nation represents a pyrethroid resistance allele that arose in West African *An. coluzzii*
 157 populations; however, the L995F allele is at higher frequency (85%) in our Burkina Faso
 158 *An. coluzzii* population, and is known to be increasing in frequency (@@REFs), there-
 159 fore L995F may provide a stronger resistance phenotype and is replacing I1527T+V402L
 160 in these populations.

161 Of the other 16 SNPs, 13 occurred almost exclusively in combination with L995F (Figure
 162 @@; @@REF Ag1000G). These include the N1570Y allele, known to enhance pyrethroid
 163 resistance in *An. gambiae* in combination with L995F. These also include two variants
 164 in codon 1874 (P1874S, P1874L). P1874S has previously been found in a colony of the
 165 crop pest *Plutoblah blahdiblah* with a pyrethroid resistance phenotype, but has not been
 166 shown to confer resistance experimentally. 10 of these variants, including N1570Y and
 167 P1874S/L, occur within internal linker domains of the protein, and so fit the model of
 168 variants that may enhance or compensate for the driver phenotype by modifying channel
 169 gating behaviour (@@CHECK; @@REFs). The remaining 3 variants are within trans-
 170 membrane domains, and so may enhance resistance by @@TODO how. Because of the
 171 tight linkage between these 13 SNPs and the L995F allele, we classify all as putative L995F

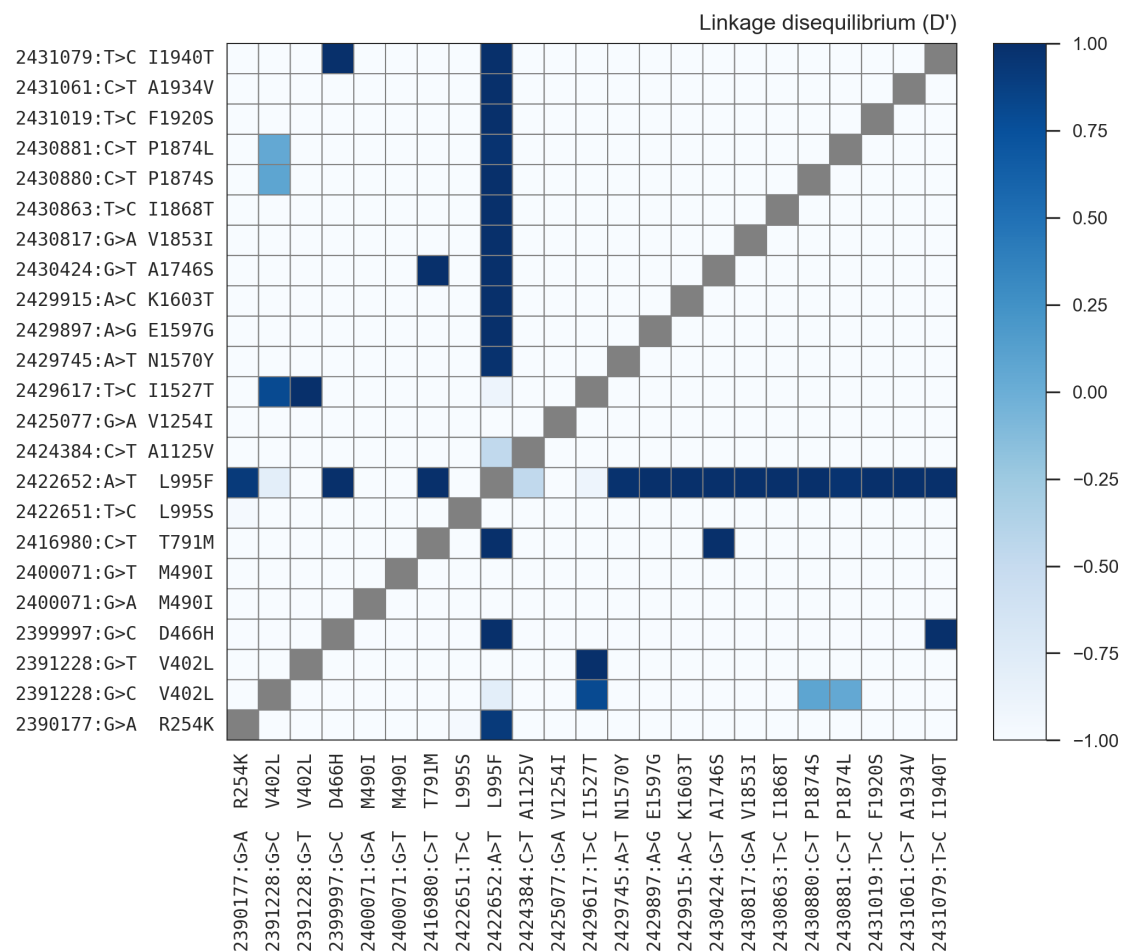


Figure 1. Linkage disequilibrium between non-synonymous variants. A value of 1 indicates that the two variants always occur in combination, and conversely a value of -1 indicates that the two variants never occur in combination. @TODO nuance this?

enhancers, although experimental work is required to confirm a resistance phenotype.

The remaining 3 variants (M490I, A1125V, V1254I) do not occur in combination with any known resistance allele, and do not appear to be associated with haplotypes under selection (@@REF Ag1000G). A possible exception is the M490I allele found at 18% frequency in the Kenyan population, although the fact that this population has experienced a recent population crash makes it difficult to test for evidence of selection at this locus. All 3 variants occur in internal linker domains, and so do not fit the model of a resistance driver, although experimental work is required to rule out a resistance phenotype.

180 Haplotype structure

181 Although it is known that pyrethroid resistance is increasing in prevalence in malaria
182 vector populations across Africa, it has not been clear whether this is being driven by the
183 spread of resistance alleles via gene flow, or by resistance alleles emerging independently in
184 multiple locations, or by some combination of both processes. The Ag1000G data resource
185 provides a potentially rich source of information about the evolutionary and demographic
186 history of insecticide resistance in any given gene, because data are available not only for
187 SNPs in gene coding regions, but also SNPs in introns and flanking intergenic regions,
188 and in neighbouring genes. These additional variants can be used to analyse the genetic
189 backgrounds (haplotypes) on which resistance alleles are found. In sexually reproducing
190 species, DNA sequences are transmitted from parents to progeny in chunks, rearranged via
191 recombination at each generation, and haplotypes convey information about this history
192 of transmission and recombination, especially when haplotypes from many individuals can
193 be compared.

194 In our initial analysis of the *Vgsc* (@@REF Ag1000G), we used 1710 biallelic SNPs
195 from within the @@70 kbp *Vgsc* gene (@@N exonic, @@N intronic) to compute the num-
196 ber of SNP differences between all pairs of 1530 haplotypes derived from 765 wild-caught
197 mosquitoes. This genetic distance measurement is a rough proxy for the degree of re-
198 latedness between haplotypes, in the sense that two haplotypes with a small number of
199 SNP differences must be closely related and share a common ancestor in the recent past.
200 This measurement cannot be used to directly estimate the time to most recent common
201 ancestor (TMRCA) for any pair of haplotypes, however, because it does not account for
202 the possibility of recombination events within the gene, which is increasingly likely for
203 pairs of haplotypes that are more distantly related. Nevertheless, it provides a useful tool
204 for exploring patterns of similarity and dissimilarity within the data. To visualise these
205 patterns, we used the pairwise genetic distances to perform hierarchical clustering, which
206 groups similar haplotypes together into clusters. We found that haplotypes carrying resis-
207 tance alleles were grouped into 10 distinct clusters. Five of these clusters carried the L995F
208 allele (labelled F1-F5), and a further five clusters carried L995S (labelled S1-S5). Within
209 each cluster, haplotypes were nearly identical across all 1710 SNPs (spanning @@70 kbp),

and therefore each cluster represents a collection of haplotypes with a very recent common ancestor. Within some of these clusters, we found haplotypes from mosquitoes collected from different locations. Specifically, cluster F1 contained haplotypes from Guinea, Burkina Faso, Cameroon and Angola; clusters @@ each contained haplotypes from Cameroon and Gabon; and cluster @@ contained haplotypes from Uganda and Kenya. The F1 cluster also contained haplotypes from both *An. gambiae* and *An. coluzzii* individuals. If we assume that haplotypes within each cluster share a common ancestor since the introduction of insecticides, which is reasonable given the high degree of similarity, then each of these clusters provides evidence that resistance alleles have been spreading between geographical locations and species via adaptive gene flow. Here we present several new analyses of these haplotype data, to confirm our initial inferences regarding gene flow, and provide further details regarding the origins and movement of resistance alleles.

To provide an alternative view of the genetic similarity between haplotypes carrying resistance alleles, we used haplotype data from within the *Vgsc* gene region to construct median-joining networks (Figure 2). This analysis is very similar to hierarchical clustering, except that it allows for the reconstruction and placement of intermediate haplotypes that may not be observed in the data. We constructed these networks up to a maximum distance of @@2 SNP differences, to ensure that each connected component in the resulting networks represents a collection of haplotypes with a recent common ancestor, and thus which is also likely to be minimally affected by recombination within the gene. For haplotypes carrying L995F, the resulting network confirms the presence of five distinct clusters, with close correspondance to the clusters F1-F5 identified previously. The L995S network also confirms five distinct clusters, in concordance with our previous analysis.

The haplotype networks bring into sharp relief the explosive evolution of amino acid substitutions secondary to the L995F allele. Within the F1 network, nodes carrying non-synonymous variants radiate out from a central node carrying only L995F, indicating that the central node represents the ancestral haplotype carrying L995F alone which initially came under selection, and these secondary variants have arisen subsequently as new mutations. Many of the nodes carrying secondary variants are large, consistent with positive selection and a functional role for these secondary variants as enhancers of the L995F resistance phenotype. The F1 network also allows us to infer multiple introgression events

241 between the two species. The central (ancestral) node comprises haplotypes from both
 242 species, as do nodes carrying the N1570Y, P1874L, and @@TODO one more variant@@.
 243 This structure is consistent with an initial introgression of the ancestral F1 haplotype, fol-
 244 lowed by introgression of haplotypes carrying secondary mutations. The contrast between
 245 the haplotype networks for the L995F and L995S alleles is striking because of the near-
 246 total absence of non-synonymous variation within the L995S networks. As we reported
 247 previously, this difference is highly significant – the ratio of non-synonymous to synony-
 248 mous nucleotide diversity (π_N/π_S) is @N times higher among haplotypes carrying
 249 L995F relative to haplotypes carrying L995S (@Test; $P=@@$) (@REF Ag1000G). Some
 250 secondary variants are present within the L995S networks, but all are at low frequency,

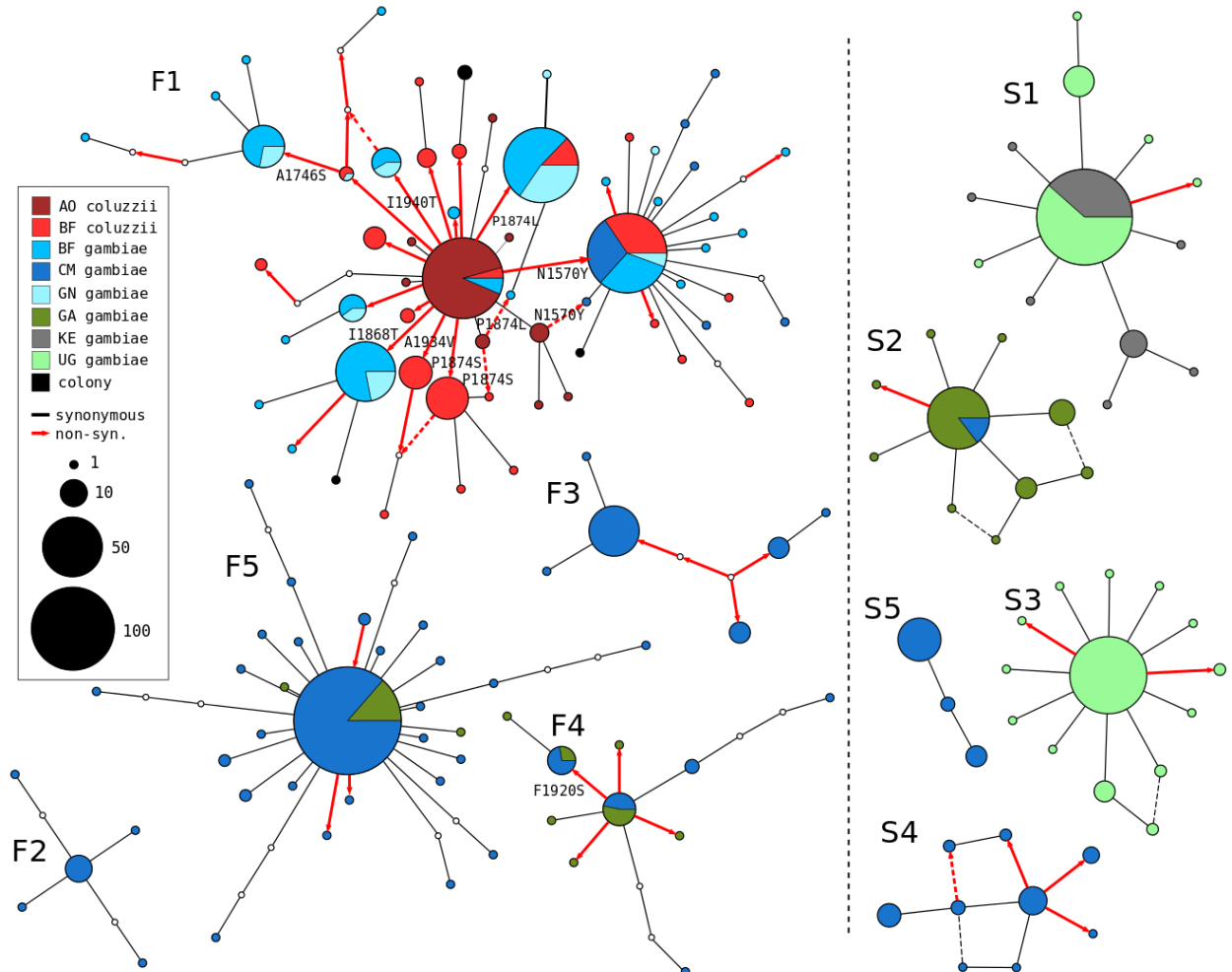


Figure 2. Haplotype networks. @@TODO redo the figure. @@TODO annotate non-syn edges in cluster F3. @@TODO mention if any clusters fixed for non-syn variants so not shown. @@TODO annotate other non-syn edges, e.g., in S4?

251 and thus may be neutral or mildly deleterious variants that are hitch-hiking on selective
252 sweeps for the L995S allele.

253 While the haplotype clustering and network analyses provide evidence for the spread
254 of resistance alleles via adaptive gene flow, and for the secondary evolution of L995F
255 enhancer alleles, they have several limitations. Within haplotype clusters where gene flow
256 has occurred, they have poor resolution to infer the origin and direction of gene flow. This
257 is because the analyses only leverage information about genetic distance within the *Vgsc*
258 gene, and for very recent events, insufficient time has elapsed for informative mutations
259 to accumulate within this relatively small genome region. Also, the fact that we observe
260 five distinct clusters for each of the codon 995 alleles suggests that each cluster is in some
261 sense independent from the others, and thus gene flow is not required for resistance to
262 emerge in multiple geographical locations. However, the threshold for the genetic distance
263 at which we have chosen to divide haplotypes into different networks or clusters is to
264 a certain extent arbitrary, and based on an intuitive sense of how much variation could
265 have accumulated among the descendants of a single resistant ancestor since the onset of
266 selective pressure. We also need to clarify what we mean by “independent”, as there are
267 several possible scenarios under which resistance could evolve in multiple populations in
268 the absence of gene flow. Finally, analyses of genetic distance within a fixed genome region
269 can be confounded by recombination events occurring within that region. For example,
270 a recombination event within the *Vgsc* gene upstream of codon 995 could cause us to
271 split a collection of haplotypes into two clusters, even though they are ancestrally related
272 within the region downstream of the recombination event. In the next sub-sections we
273 provide some conceptual foundations to help clarify these ambiguities, and use analyses
274 of haplotype sharing from the genome regions flanking the *Vgsc* gene to provide finer
275 resolution to diagnose recent gene flow events.

276 **Insecticide resistance outbreaks**

277 To provide an aid to further interpretation of the genetic data, and relating them to the
278 challenges of insecticide resistance management, we introduce the concept of an **insec-**
279 **ticide resistance outbreak**. Informally, we define a resistance outbreak by analogy
280 with the epidemiological concept of an outbreak, as a rapid increase in the prevalence

281 of insecticide resistance among mosquitoes at a particular place and time. Note that
 282 this does not imply that the overall abundance of mosquitoes is increase, just that the
 283 relative frequency of resistance within mosquito populations is increasing. We also re-
 284 quire that all occurrences of insecticide resistance within the same outbreak are connected
 285 by a chain of transmission of resistance alleles from parent to progeny mosquitoes, and
 286 thus can be traced back to a single resistant common ancestor. A resistance outbreak
 287 can be **localised**, meaning that it affects a small group of mosquitoes of a single species
 288 from a limited geographical area. Alternatively, a resistance outbreak may be **spreading**,
 289 meaning that resistance alleles have been transmitted since the introduction of insecti-
 290 cides by interbreeding of mosquitoes of different species and/or originating from different
 291 geographical locations.

292 Our goal for the *Vgsc* gene can now be restated, which is to perform an insecticide
 293 resistance outbreak analysis. We would like to diagnose how many separate outbreaks have
 294 occurred, which outbreaks are localised, and which are spreading. For spreading outbreaks,
 295 we would like to reconstruct the path of transmission of resistance alleles between mosquito

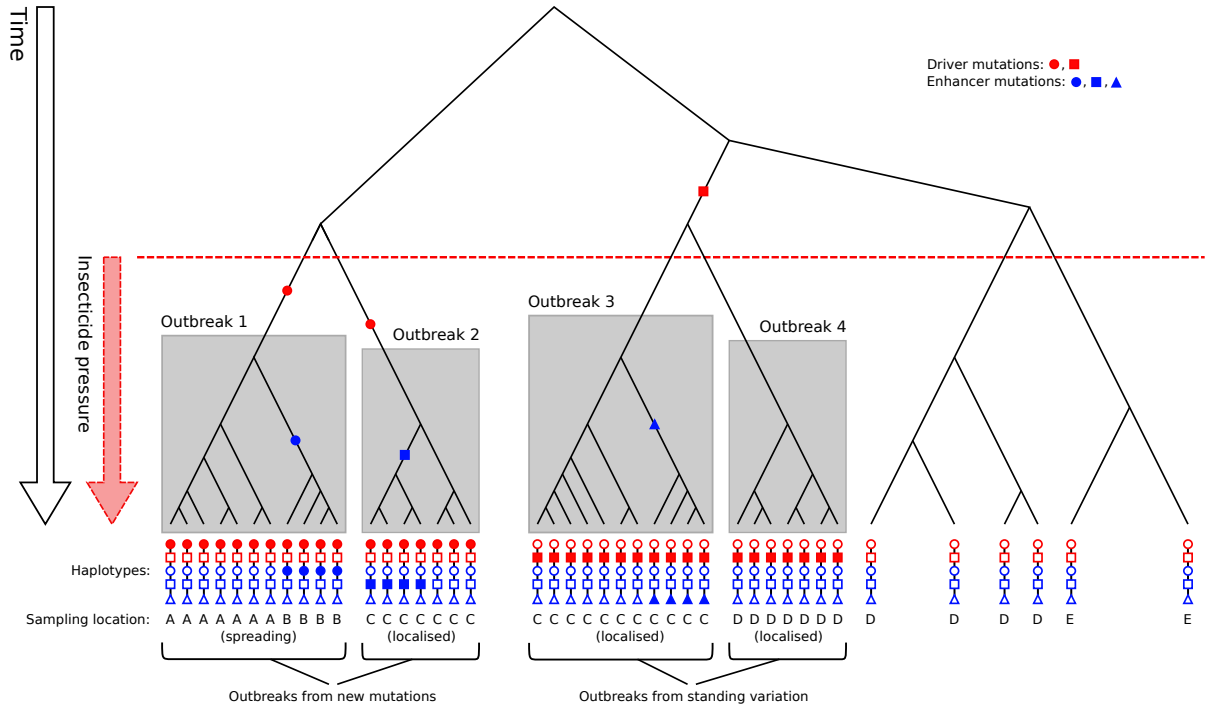


Figure 3. Illustration of insecticide resistance outbreaks. @@TODO explanation.

populations, and to provide information on the probable source. We would, of course, also like to identify the primary and secondary genetic factors that are driving each outbreak. Stated in this way, it is easier to discuss how this information is potentially relevant to insecticide resistance management, and to frame key epidemiological questions. For example, we would like to begin to build a picture of where and when local conditions have favoured the evolution of insecticide resistance, and whether those conditions are relatively patchy (and hence outbreaks are mainly localised) or whether conditions are consistent over broad areas (and hence can support a spreading outbreak). We would also like to know which mosquito populations are sufficiently connected to enable outbreak spread, and if there is any consistent pattern to the direction of spread. This information could be relevant to discussions about how resources for insecticide resistance management might be targeted, what strategies are appropriate in which settings, and where and when insecticide resistance management needs to be coordinated between different countries and/or at different levels of administration.

For clarity, we also define the concept of an insecticide resistance outbreak formally in terms of coalescent theory, as a collection of lineages (1) sharing a resistance driver allele by descent, (2) coalescing more recently than the onset of insecticide pressure, and (3) having increased in frequency because of positive selection due to insecticides. This definition is illustrated for four hypothetical outbreaks in Figure 3. Because mosquitoes are sexually recombining, genealogical trees vary along the genome, and so we define resistance outbreaks with respect to a specific gene locus, which for the present study is codon 995 within the *Vgsc* gene. Note that separate outbreaks may be driven by the same resistance allele, and this can occur if multiple mutational events occur after the introduction of insecticides (Figure 3, outbreaks 1 and 2), or if a resistance allele is present in mosquito populations as standing variation prior to insecticide use (Figure 3, outbreaks 3 and 4). Here we are primarily concerned with whether outbreaks are localised or spreading, because this has immediate epidemiological relevance. We do not attempt to infer whether separate outbreaks with the same driver allele arose via standing variation or new mutations, however this is an interesting biological question to address in future studies. As a technical note, there is a simple correspondance with terminology conventionally used in the population genetics literature to describe selective sweeps. At

327 a given gene locus, a hard selective sweep gives rise to a single resistance outbreak, and a
 328 soft selective sweep gives rise to multiple resistance outbreaks.

329 **Outbreak analysis from haplotype age**

330 As described above, haplotype data from genome regions both within and flanking the
 331 *Vgsc* gene provide a higher resolution for reconstructing recent historical events. To lever-
 332 age this information, we used a heuristic approach to estimate the time to most recent
 333 common ancestor (TMRCA) or “age” for each pair of haplotypes in our dataset, centering
 334 the analysis on *Vgsc* codon 995. For each pair of haplotypes, we estimated the length
 335 of the region shared identical by descent (IBD), and the number of mutations that have
 336 accumulated since the most recent common ancestor. We then combined these two pieces
 337 of information to produce a point estimate for the haplotype age (Methods). We studied
 338 the overall distribution of pairwise haplotype ages (Figure 4), and used hierarchical clus-
 339 tering to construct a dendrogram and visualise the overall age structure (Figure 5). We
 340 caution that although the estimated ages are in units of generations, these estimates have
 341 not been calibrated, and there is substantial uncertainty regarding both the mutation and
 342 recombination rate parameters. The ages therefore should not be interpreted as reliable
 343 absolute values, but they can be compared to each other to investigate the relative age of
 344 different events.

345 A key feature of the overall age distribution is that it is bimodal, with a minor mode of
 346 haplotypes coalescing recently, and a major mode coalescing further in the past (Figure

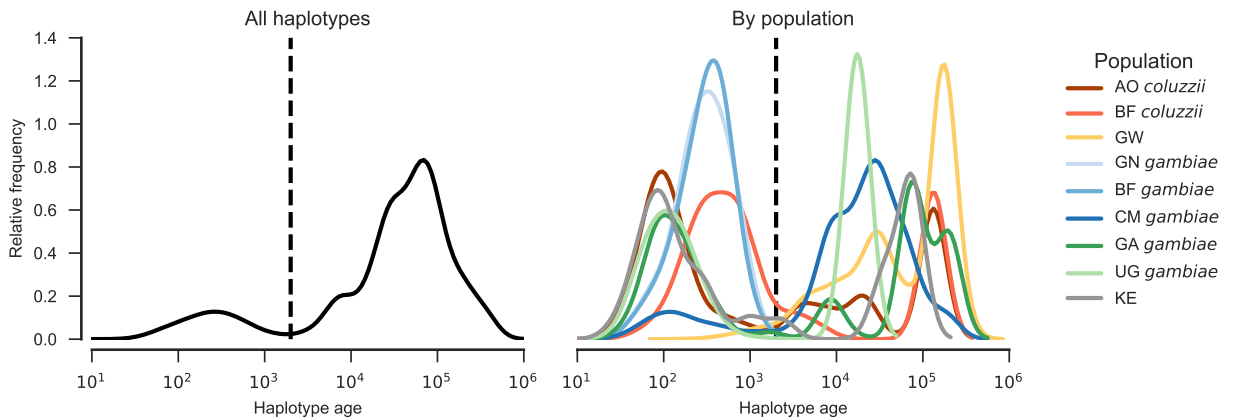


Figure 4. Haplotype age distribution. @@TODO caption.

4). This is expected at an insecticide resistance locus experiencing one or more resistance outbreaks. Within each outbreak, all haplotypes share a very recent common ancestor, but between outbreaks and among haplotypes without any resistance allele, haplotypes are more distantly related, and the distribution of ages is influenced by mosquito population size and other demographic factors. In particular, mosquito populations generally have a large effective population size (@@REF Ag1000G), and so in the absence of selection, haplotypes are expected to coalesce slowly. The bimodal age distribution is not due to geographical population structure, because the same bimodality is observed within several populations. We take the midpoint between these two modes as an estimate for the earliest time of onset of selective pressure due to insecticides, and thus for the maximum age of a resistance outbreak. To identify haplotype clusters representing putative resistance outbreaks, we then cut the haplotype dendrogram at this maximum outbreak age (Figure 5). Comparing this to previous analyses of haplotype structure based on genetic distance, we find clusters F1-F5 and S1-S3 recapitulated with close correspondence, and S4 and S5 merged into a single cluster. We label a new cluster “L@@” representing an outbreak driven by the I1527T allele in combination with one or the other V402L allele. We also label

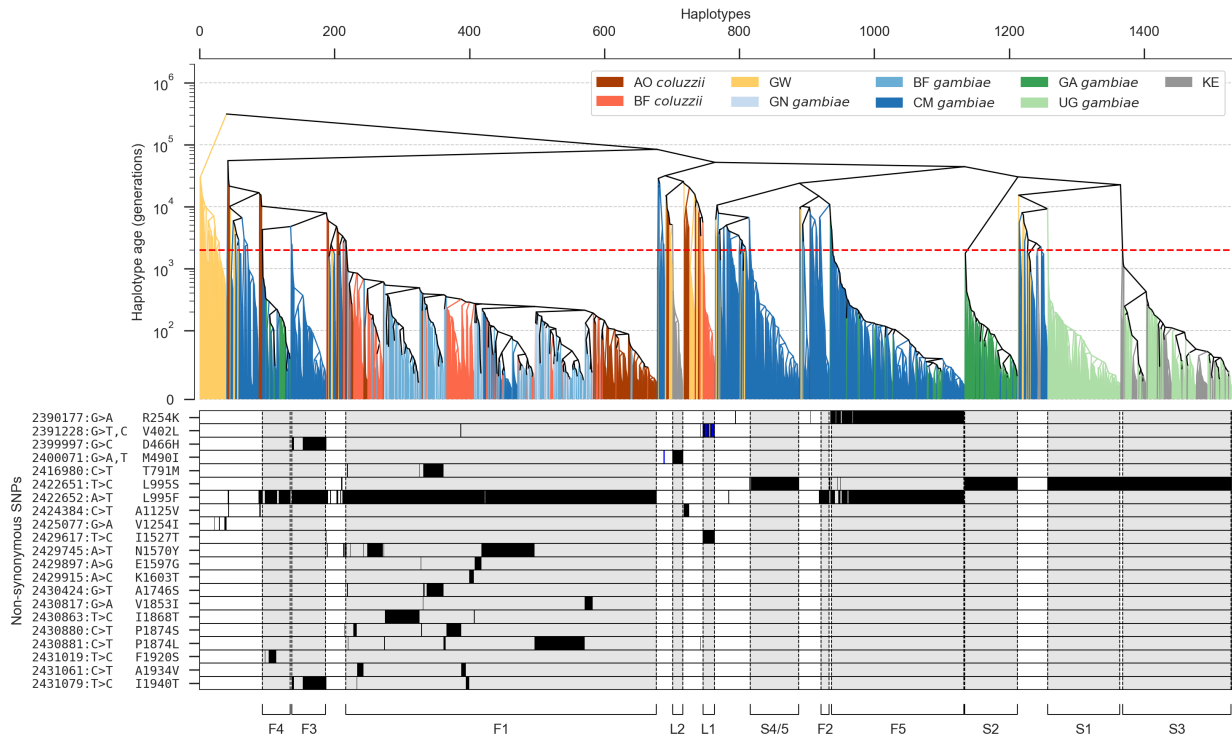


Figure 5. Clustering of haplotypes by age. @@TODO caption.

363 a cluster “L@@” capturing a set of haplotypes from Kenya carrying the M490I variant,
 364 although the fact that these haplotypes all share a recent common ancestor may be a
 365 reflection of the unusual demography of the Kenyan population which has experienced
 366 a severe population crash (@@REF) and not be due to recent selection for insecticide
 367 resistance. As in earlier analyses, clusters F1, F4, F5 and S3 all include haplotypes
 368 sampled from multiple geographical locations, and thus represent spreading outbreaks.
 369 Clusters F2, F3, S1, S2, S4/5 and L1 include only haplotypes from a single sampling
 370 location, and thus appear to represent localised outbreaks.

371 We then studied the distribution of haplotype ages within each spreading outbreak, to
 372 attempt to reconstruct information about the historical path of transmission of resistance
 373 alleles between locations. To do this, we grouped the haplotypes within each spreading
 374 outbreak by sampling location, and compared the distribution of haplotype ages both
 375 within and between locations. To aid in interpreting these data, we define three possi-
 376 ble spreading scenarios, being: (1) a directional spread from one population to another;
 377 (2) spread from an unsampled population into the sampled populations; and (3) a com-
 378 plex scenario involving multiple gene flow events. In Figure 6 we illustrate the expected
 379 genealogy and haplotype age distribution under each of these scenarios.

380 The clearest result was obtained for outbreak F1 (Figure 7). Within this outbreak,

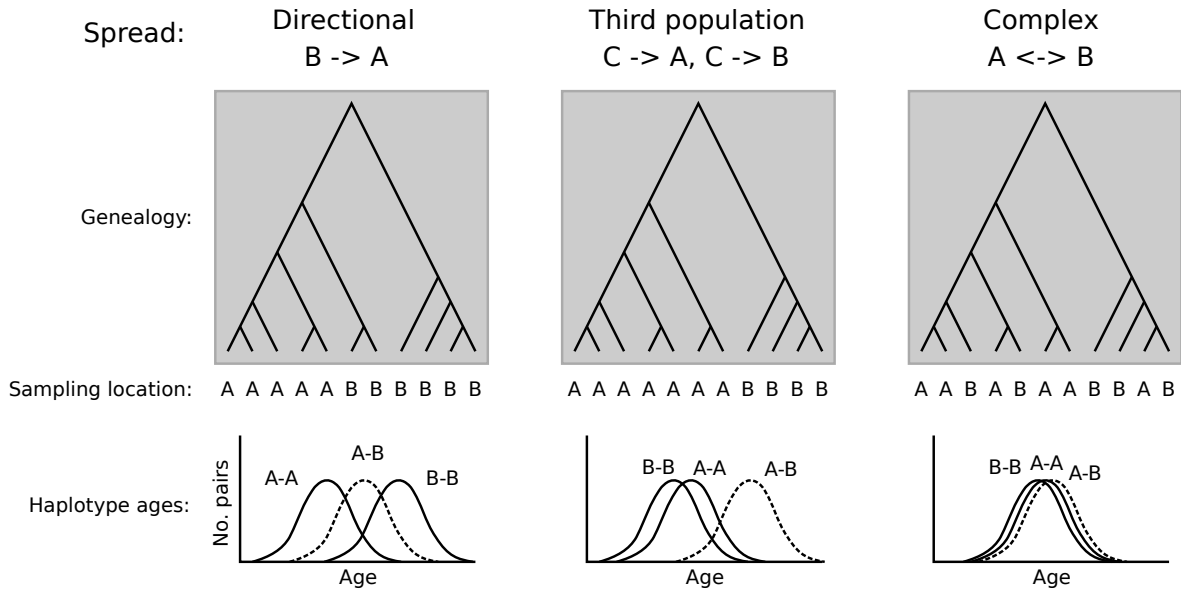


Figure 6. Inferring history of spread from haplotype ages. @@TODO explain.

haplotypes from Cameroon and Angola are significantly younger than haplotypes from Burkina Faso and Guinea. The age distributions are consistent with an outbreak originating in West Africa and subsequently spreading towards Cameroon and separately towards Angola. We were surprised that the age distributions for *An. gambiae* and *An. coluzzii* from Burkina Faso are very similar, despite the fact that previous studies have shown that introgression has occurred from *An. gambiae* into *An. coluzzii*. This may indicate that the initial introgression event happened during the early phases of the outbreak, but is also consistent with a complex history of multiple gene flow events between the species.

Outbreaks F4, F5 and S2 each involve haplotypes from both Cameroon and Gabon. Interpreting the age distributions for these outbreaks is difficult, because mosquitoes from Gabon were collected at a much earlier time point (2000) than mosquitoes from Cameroon (20@@). If our haplotype age estimates were well-calibrated, and we also had reliable estimates for the number of mosquito generations per year, then we might be able to adjust for this time difference, however we are not able to do so presently. An interesting feature of these outbreaks, however, is that we would expect haplotypes from Gabon to appear older due to the time of sampling, which is observed for outbreak S2 but not for F4 or F5. Indeed, S2 is at a high frequency among all Gabon haplotypes and a low frequency among Cameroon haplotypes, whereas the reverse is true for F4 and F5. These data suggest that F4 and F5 have spread from Cameroon towards Gabon, while S2 has

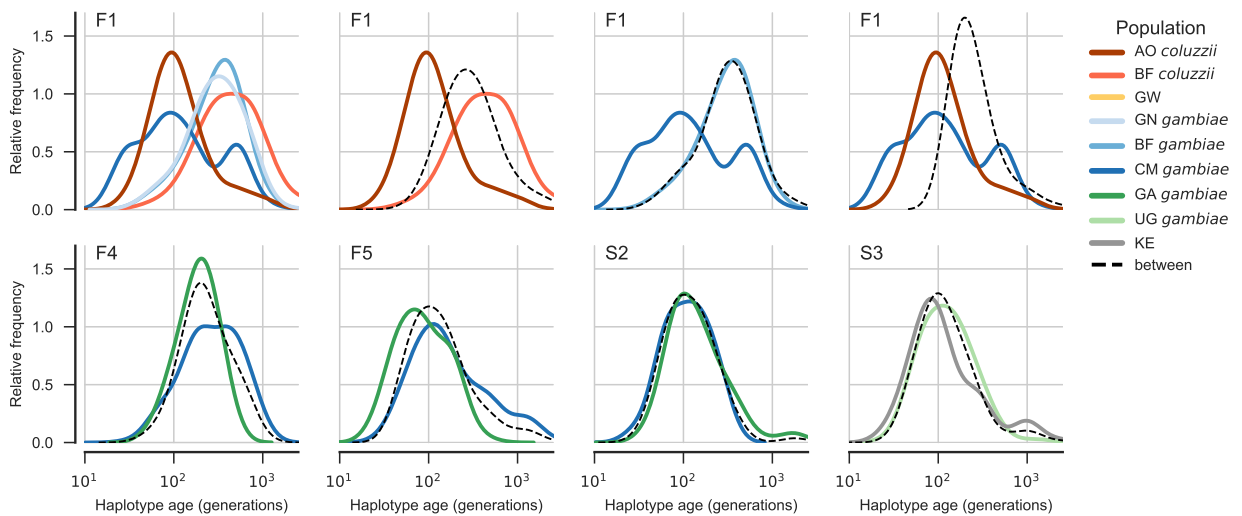


Figure 7. Haplotype age distributions within spreading outbreaks. @@TODO caption.

spread in the opposite direction. A lot can happen in mosquito populations in @N years, however, and these conclusions remain highly speculative pending further sampling from both locations.

For outbreak S3 involving haplotypes from Uganda and Kenya, the age distributions do not suggest any clear direction of gene flow. This could reflect multiple gene flow events in either or both directions. However, another outbreak (S1) is localised in Uganda and represented within the Ugandan population at roughly equal frequency with S3. If transmission was occurring from Uganda towards Kenya, we might expect both outbreaks to have spread to Kenya. Thus the localisation of S1 suggests S3 has spread into Uganda from Kenya or another location. Again, this conclusion remains tentative and requires confirmation via further sampling.

To summarise these conclusions in a concise way, we have depicted the distribution and spread of resistance outbreaks via the map shown in Figure 8. We have plotted haplotypes from each sampling location as a pie chart. The overall size of each pie chart represents the number of haplotypes sampled, and coloured wedges within each pie represent the frequency of each resistance outbreak within the population. Coloured arrows are used

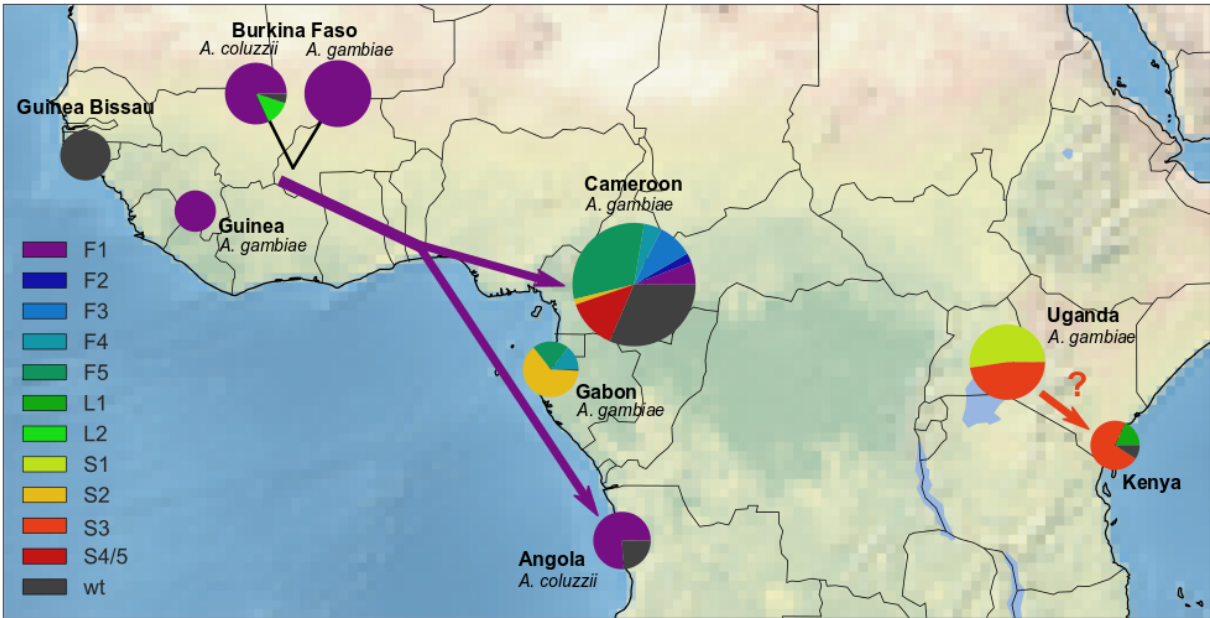


Figure 8. Geographical distribution of resistance outbreaks. @TODO arrows for Gabon <-> Cameroon. @TODO change arrow for Kenya -> Uganda. @TODO add source area for F1.

to depict our inferences regarding the transmission paths for spreading outbreaks. Our conclusions regarding direction of spread for outbreaks F4, F5, S2 and S3 are tentative, and we indicate this with a question mark. Because of the relatively sparse geographical representation within the Ag1000G phase 1 dataset, and the fact that collections were not synchronized but span several years, we cannot be precise about the geographical origins of these resistance outbreaks. Even for outbreak F1 where we have clear evidence of spread from West Africa towards Central and Southern Africa, we have only sampled mosquitoes from Guinea and Burkina Faso, and the true source of the outbreak may not be either of these countries. We indicate this uncertainty regarding the outbreak source as a coloured area with a dashed border. This representation is imperfect, as is our knowledge regarding the sources and transmission paths of these outbreaks, but we hope this depiction may at least serve to stimulate further sampling, analysis and discussion, with the aim of improving our knowledge of resistance outbreaks for *Vgsc* as well as other insecticide resistance genes.

Design of genetic assays for outbreak surveillance

The insecticide resistance outbreaks we have identified here are undoubtedly ongoing, affecting many more mosquito populations than we have sampled in Ag1000G phase 1, and continuing to spread. In addition, other outbreaks may be occurring in populations that we have not sampled, or in populations we have sampled but since the sampling date. Whole genome sequencing of individual mosquitoes clearly provides data of sufficient resolution to detect resistance outbreaks and provide ongoing outbreak surveillance. The cost of whole genome sequencing continues to fall, with the present cost being approximately 100 GBP to obtain 30X coverage of an individual genome. Mobile sequencing technology is also developing rapidly, and may be a realistic prospect for mosquito population surveillance within a few years. There is an interim period, however, during which it may be more practical to develop targeted genetic assays for outbreak surveillance that could scale to tens of thousands of mosquitoes at a fraction of the cost of whole genome sequencing. For example, SNP genotyping using mass spectrometry and amplicon sequencing are two available technologies that could be applied now at scale and at modest cost.

To facilitate the development of targeted genetic assays for *Vgsc* insecticide resistance

446 outbreak surveillance, we have produced two supplementary data tables. In Supplemen-
 447 tary Table 1 we provide a list of all biallelic SNPs discovered with high confidence in this
 448 study within the *Vgsc* gene and in the 100 kbp upstream and downstream flanking re-
 449 gions. Both amplicon sequencing and genotyping by mass spectrometry require the design
 450 of PCR primers to amplify the targeted genome region. To aid in primer design, for each
 451 SNP we provide the flanking sequence for 250 bp upstream and downstream of the SNP
 452 position, including information about polymorphisms within these flanking regions. Not
 453 all SNPs are informative for detecting whether an individual mosquito carries a haplotype
 454 from a resistance outbreak, and we provide some summary statistics for each SNP to aid in
 455 the selection of the most informative SNPs. For each SNP we report the allele frequencies
 456 within each of the outbreaks identified here, as well as for populations of susceptible hap-

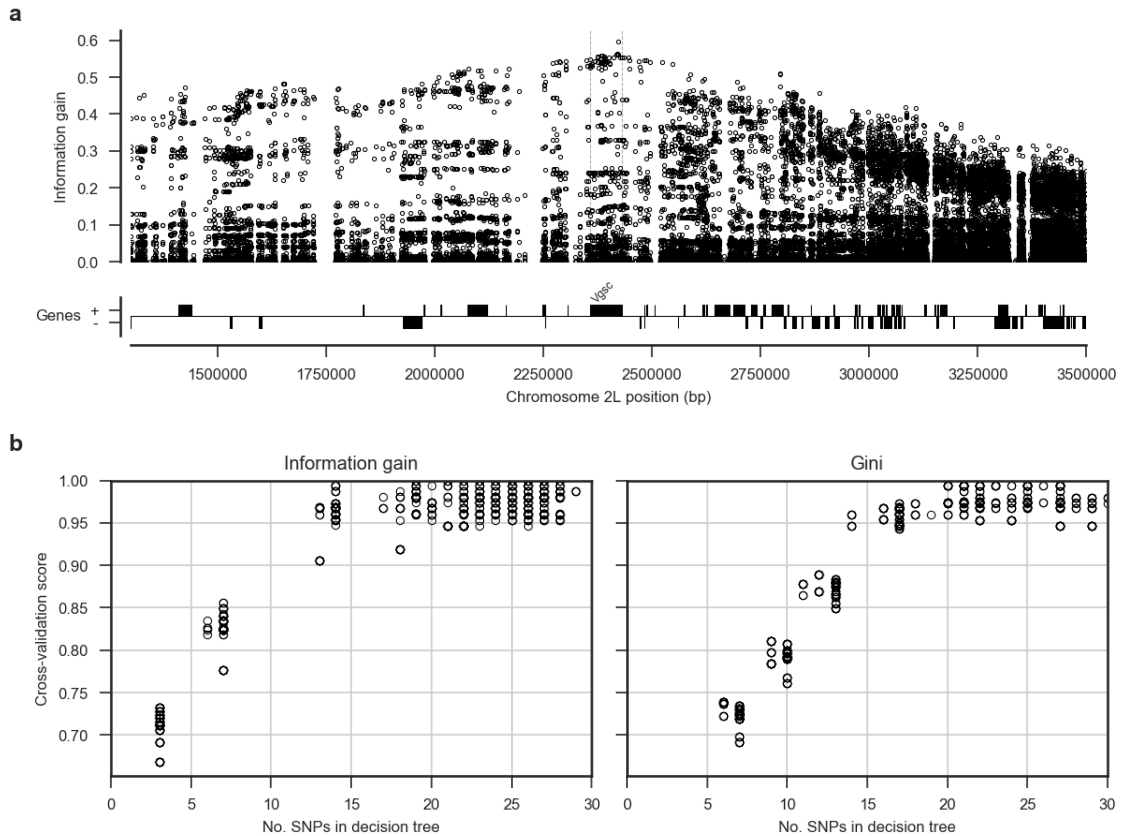


Figure 9. Informative SNPs for outbreak surveillance. **a**, Each data point represents a single SNP. The information gain value for each SNP provides an indication of how informative the SNP is likely to be if used as part of a genetic assay for testing whether a mosquito carries a resistance haplotype, and if so, which resistance outbreak it derives from. **b**, Number of SNPs required to accurately classify which outbreak a haplotype derives from. Decision trees using either information gain or the Gini impurity index as the decision criterion. Accuracy was evaluated using 10-fold stratified cross-validation.

lotypes. We also provide the overall variance in allele frequencies, the information gain, and the Gini impurity for each SNP. Note that recombination events are more likely at increasing distances upstream and downstream of the resistance variants under selection, and thus the most informative SNPs are found closest to the resistance variants within the gene (Figure 9). However, SNPs with some information gain are available throughout the gene and in flanking regions.

We suggest that the design of a genetic assay proceed by (1) performing an initial round of filtering to remove SNPs which are not informative (e.g., low information gain); (2) performing a round of primer design to remove SNPs for which primers are unlikely to be successful; (3) performing a full analysis of the remaining SNPs to select a subset that is sufficient to classify all outbreaks identified here; (4) finalise primer designs for the chosen panel of SNPs. A possible methodology for step 3 would be to use an algorithm such as ID3 to build a decision tree. To aid in the development of a classification algorithm, in Supplementary Table 2 we provide our classification for each of the 1530 haplotypes sampled here, along with the alleles carried by each haplotype for each of the SNPs included in Supplementary Table 1. To test the methodology, we constructed decision trees using either information gain (LD3) or the Gini impurity index as a decision criterion, and using all available SNPs from within the *Vgsc* plus 20 kbp flanking regions as input features (i.e., assuming primers could be designed in all cases). Figure 9b shows the cross-validation scores obtained for trees constructed allowing increasing numbers of SNPs. This analysis suggests that it should be possible to construct a tree able to classify haplotypes from all 10 resistance outbreaks with >95% accuracy using 15 SNPs or less.

Recombination

As mentioned earlier, analyses of haplotype structure based on genetic distance within the fixed window of the *Vgsc* gene could be affected if recombination events occurred within the gene. Our analyses of haplotype age should be less affected by recombination, because they explicitly take recombination into account, estimating the positions at which recombination events have occurred to interrupt regions shared IBD between pairs of haplotypes. However, these analyses were based on a heuristic method for estimating recombination breakpoints, and there are several potential sources of error. To study

the evidence for recombination within the genome region spanning the *Vgsc* gene, and provide some additional confirmation that our inferences regarding insecticide resistance outbreaks have not been affected by recombination or other sources of error, we performed an additional analysis of genetic distance between haplotypes. We first constructed a putative ancestral haplotype for each of the outbreaks we identified, by starting from the codon 995 position and separately moving upstream and downstream, assuming the major allele at each SNP bifurcation point represents the ancestral haplotype. We then computed the genetic distance (D_{XY}) between each of our sampled haplotypes and each of the inferred ancestral outbreak haplotypes, computing the distance in overlapping windows of bp across a 2 Mbp region spanning the *Vgsc* gene. The results for outbreaks F1-F5 are plotted in Figure 10, and outbreaks S1-S4/5 are shown in Figure ???. In these plots we expect that all haplotypes from a given outbreak should share very close genetic similarity ($D_{XY} \approx 0$) with each other and with the ancestral haplotype for that outbreak within the *Vgsc* gene itself, with an increasing number of haplotypes recombining away from the ancestral outbreak haplotype as we move away from the gene in either the upstream or downstream direction. Conversely, haplotypes from one outbreak should not share any close genetic similarity ($D_{XY} > 0$) with the inferred ancestral haplotype from a different outbreak, either within the *Vgsc* gene or in flanking regions.

The results for all outbreaks are largely consistent with this expectation. For this analysis we treated S4/5 as a single outbreak, as indicated by the haplotype age analysis, and we can gain some insight into why these two were split into separate clusters in earlier analyses. All haplotypes in the S4/5 outbreak share close similarity with the ancestral haplotype on both flanks of the *Vgsc* gene, but there is a short region of within the gene where a subset of haplotypes are diverged. This region of divergence accounts for the S4/S5 split in earlier analyses. @TODO explain @TODO also note relatively low divergence among F2, F3, F4 on upstream flank and explain

Discussion

@TODO Discuss accessibility, have we missed any functional variation?

@TODO Discuss weaknesses, caveats and potential improvements to method for esti-

516 mating haplotype age.

517 @@TODO What are the implications for insecticide resistance management? Realisti-
518 cally how could this information be used?

519 @@TODO What about DDT? If prior selection for DDT resistance, how might this
520 complicate the picture? Do we see any evidence for multiple phases of selection?

521 @@TODO Speculate on why L995F but not L995S has evolved secondary variation.

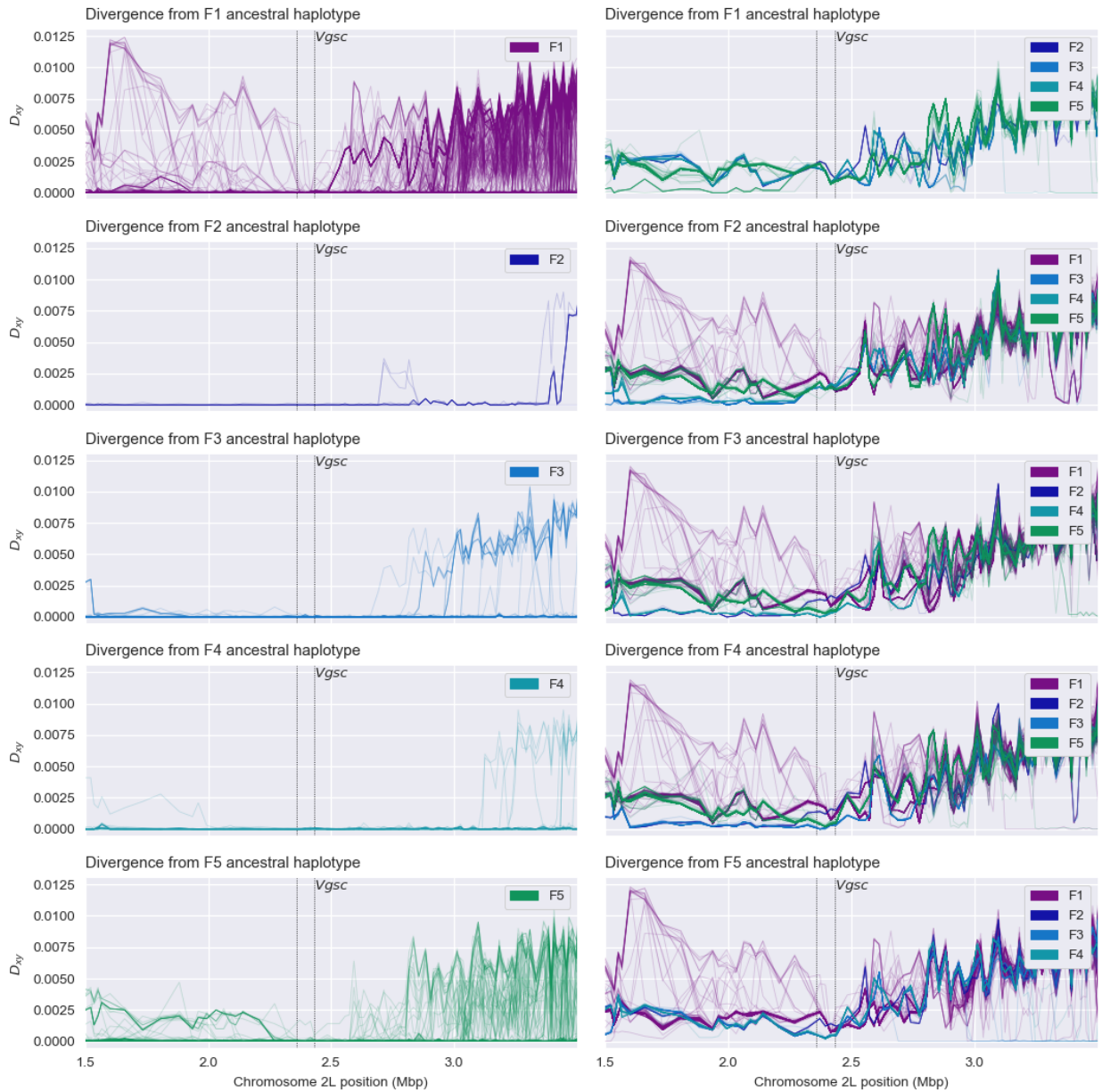


Figure 10. Recombination and ancestral haplotypes for L995F. @@TODO legend

522 **Methods**

523 @@TODO

524 **References**

- 525 [1] S. Bhatt et al. ‘The effect of malaria control on Plasmodium falciparum in Africa
526 between 2000 and 2015’. In: *Nature* 526.7572 (2015), pp. 207–211. ISSN: 0028-0836.
527 arXiv: [arXiv:1011.1669v3](#).
- 528 [2] Janet Hemingway et al. ‘Averting a malaria disaster: Will insecticide resistance derail
529 malaria control?’ In: *The Lancet* 387.10029 (2016), pp. 1785–1788. ISSN: 1474547X.
- 530 [3] Ana Paula B Silva et al. ‘Mutations in the voltage-gated sodium channel gene of
531 anophelines and their association with resistance to pyrethroids: a review’. In: *Par-*
532 *asites & Vectors* 7.1 (2014), p. 450. ISSN: 1756-3305.
- 533 [4] N’Guessan R, Corbel V, Akogbéto M, Rowland M. ‘Reduced efficacy of insecticide-
534 treated nets and indoor residual spraying for Malaria control in area of pyrethroid
535 resistance, Benin.’ In: *Emerging Infectious Diseases* 13 (2007), pp. 199–206. ISSN:
536 10806040.
- 537 [5] Kobié H. Toé et al. ‘Increased pyrethroid resistance in malaria vectors and decreased
538 bed net effectiveness Burkina Faso’. In: *Emerging Infectious Diseases* 20.10 (2014),
539 pp. 1691–1696. ISSN: 10806059.
- 540 [6] World Health Organization. *Implications of insecticide resistance for malaria vector*
541 *control*. Tech. rep. Geneva, 2016.
- 542 [7] World Health Organization. *Global Plan for Insecticide Resistance Management*
543 *(GPIRM)*. Tech. rep. Geneva, 2012.
- 544 [8] T. G.E. Davies et al. ‘A comparative study of voltage-gated sodium channels in the
545 Insecta: Implications for pyrethroid resistance in Anopheline and other Neopteran
546 species’. In: *Insect Molecular Biology* 16.3 (2007), pp. 361–375. ISSN: 09621075.

- [9] D. Martinez-Torres et al. ‘Molecular characterization of pyrethroid knockdown resistance (kdr) in the major malaria vector *Anopheles gambiae* s.s.’ In: *Insect Molecular Biology* 7.2 (1998), pp. 179–184. ISSN: 09621075.
- [10] H. Ranson et al. ‘Identification of a point mutation in the voltage-gated sodium channel gene of Kenyan *Anopheles gambiae* associated with resistance to DDT and pyrethroids’. In: *Insect Molecular Biology* 9.5 (2000), pp. 491–497. ISSN: 09621075.
- [11] Christopher M Jones et al. ‘Footprints of positive selection associated with a mutation (N1575Y) in the voltage-gated sodium channel of *Anopheles gambiae*.’ In: *Proceedings of the National Academy of Sciences of the United States of America* 109.17 (2012), pp. 6614–9. ISSN: 1091-6490.
- [12] Martin S. Williamson et al. ‘Identification of mutations in the housefly para-type sodium channel gene associated with knockdown resistance (kdr) to pyrethroid insecticides’. In: *Molecular and General Genetics* 252.1-2 (1996), pp. 51–60. ISSN: 00268925.
- [13] T. G. E. Davies et al. ‘DDT, pyrethrins, pyrethroids and insect sodium channels’. In: *IUBMB Life* 59.3 (2007), pp. 151–162. ISSN: 1521-6543.
- [14] Frank D. Rinkevich, Yuzhe Du and Ke Dong. ‘Diversity and convergence of sodium channel mutations involved in resistance to pyrethroids’. In: *Pesticide Biochemistry and Physiology* 106.3 (2013), pp. 93–100. ISSN: 00483575. arXiv: NIHMS150003.
- [15] Ke Dong et al. *Molecular biology of insect sodium channels and pyrethroid resistance*. 2014. arXiv: 15334406.
- [16] Andrias O. O’Reilly et al. ‘Modelling insecticide-binding sites in the voltage-gated sodium channel’. In: *Biochemical Journal* 396.2 (2006), pp. 255–263. ISSN: 0264-6021.
- [17] J Pinto et al. ‘Multiple origins of knockdown resistance mutations in the Afrotropical mosquito vector *Anopheles gambiae*’. In: *PLoS One* 2 (2007), e1243. ISSN: 19326203.
- [18] Josiane Etang et al. ‘Polymorphism of intron-1 in the voltage-gated sodium channel gene of *Anopheles gambiae* s.s. populations from cameroon with emphasis on insecticide knockdown resistance mutations’. In: *Molecular Ecology* 18.14 (2009), pp. 3076–3086. ISSN: 09621083.

- 576 [19] Federica Santolamazza et al. ‘Remarkable diversity of intron-1 of the para voltage-
577 gated sodium channel gene in an *Anopheles gambiae*/*Anopheles coluzzii* hybrid
578 zone.’ In: *Malaria journal* 14.1 (2015), p. 9. ISSN: 1475-2875.
- 579 [20] Chris Bass et al. ‘Detection of knockdown resistance (<i>kdr</i>) mutations in
580 <i>Anopheles gambiae</i>: a comparison of two new high-throughput assays with
581 existing methods.’ In: *Malaria journal* 6 (2007), p. 111. ISSN: 1475-2875.
- 582 [21] A. Dao et al. ‘Signatures of aestivation and migration in Sahelian malaria mosquito
583 populations’. In: *Nature* 516.7531 (2014), pp. 387–390. ISSN: 0028-0836.
- 584 [22] Chris S. Clarkson et al. ‘Adaptive introgression between *Anopheles* sibling species
585 eliminates a major genomic island but not reproductive isolation’. In: *Nature Com-*
586 *munications* 5 (2014). ISSN: 2041-1723.
- 587 [23] L. J. Reimer et al. ‘An unusual distribution of the *kdr* gene among populations of
588 *Anopheles gambiae* on the island of Bioko, Equatorial Guinea’. In: *Insect Molecular*
589 *Biology* 14.6 (2005), pp. 683–688. ISSN: 09621075.
- 590 [24] Ag1000g Consortium. ‘Natural diversity of the malaria vector *Anopheles gambiae*’.
591 In: *Nature* ?? (2017), ?