

1 **Genome variation and population structure**
2 **among 1,142 mosquitoes of the African**
3 **malaria vector species *Anopheles gambiae***
4 **and *Anopheles coluzzii***

5 The *Anopheles gambiae* 1000 Genomes Consortium¹

6 ¹A list of consortium members appears at the end of the paper

7 18th February 2020

8 **Abstract**

9 Mosquito control remains a central pillar of efforts to reduce malaria burden in
10 sub-Saharan Africa. However, insecticide resistance is entrenched in malaria vector
11 populations, and countries with high malaria burden face a daunting challenge to
12 sustain malaria control with a limited set of surveillance and intervention tools. Here
13 we report on the second phase of a project to build an open resource of high quality
14 data on genome variation among natural populations of the major African malaria
15 vector species *Anopheles gambiae* and *Anopheles coluzzii*. We analysed whole genomes
16 of 1,142 individual mosquitoes sampled from the wild in 13 African countries, and
17 a further 234 individuals comprising parents and progeny of 11 lab crosses. The
18 data resource includes high confidence single nucleotide polymorphism (SNP) calls
19 at 57 million variable sites, genome-wide copy number variation (CNV) calls, and
20 haplotypes phased at biallelic SNPs. We used these data to analyse genetic population
21 structure, and characterise genetic diversity within and between populations. We also
22 illustrate the utility of these data by investigating species differences in isolation by
23 distance, genetic variation within proposed gene drive target sequences, and patterns
24 of resistance to pyrethroid insecticides. This data resource provides a foundation for

25 developing new operational systems for molecular surveillance, and for accelerating
26 research and development of new vector control tools.

27 **Introduction**

28 The 10 countries with the highest malaria burden in Africa account for 65% of all malaria
29 cases globally, and attempts to reduce that burden further are stalling in the face of
30 significant challenges [1]. Not least among these, resistance to pyrethroid insecticides is
31 widespread throughout African malaria mosquito populations, potentially compromising
32 the efficacy of mosquito control interventions which remain a core tenet of global malaria
33 strategy [2, 3]. There is a broad consensus that further progress cannot be made if in-
34 terventions are applied blindly, but must instead be guided by data from epidemiological
35 and entomological surveillance [4]. Genome sequencing technologies are considered to be a
36 key component of future malaria surveillance systems, providing insights into evolutionary
37 and demographic events in mosquito and parasite populations that are otherwise difficult
38 to obtain [5]. Genomic surveillance systems will not work in isolation, but will depend
39 on high quality open genomic data resources, including baseline data on genome variation
40 from multiple mosquito species and geographical locations, against which comparisons can
41 be made and inferences regarding new events can be drawn.

42 Better surveillance can increase the impact and longevity of available mosquito control
43 tools, but sustaining malaria control will also require the development and deployment of
44 new tools [4]. This includes repurposing existing insecticides not previously used in public
45 health [6, 7], developing entirely new insecticide classes, and developing tools that don't
46 rely on insecticides, such as genetic modification of mosquito populations [8]. Research and
47 development of new mosquito control tools has been greatly facilitated by the availability of
48 open genomic data resources, including high quality genome assemblies [9, 10], annotations
49 [11], and more recently by high quality resources on genetic variation among natural
50 mosquito populations [12]. Further expansion of these open data resources to incorporate
51 unsampled mosquito populations and new types of genetic variation can provide new
52 insights into a range of biological and ecological processes, and help to accelerate scientific
53 discovery from basic biology through to operational research.

54 The *Anopheles gambiae* 1000 Genomes (Ag1000G) project¹ was established in 2013 to
55 build a large scale open data resource on natural genetic variation in malaria mosquito
56 populations. The Ag1000G project forms part of the Malaria Genomic Epidemiology Net-
57 work² (MalariaGEN), a data-sharing community of researchers investigating how genetic
58 variation in humans, mosquitoes and malaria parasites can inform the biology, epidemi-
59 ology and control of malaria. The first phase of the Ag1000G project released data from
60 whole genome Illumina deep sequencing of the major Afrotropical malaria vector species
61 *Anopheles gambiae* and *Anopheles coluzzii* [12], two closely related siblings within the
62 *Anopheles gambiae* species complex [13]. Mosquitoes were sampled in 8 African countries
63 from a broad geographical range, spanning Guinea-Bissau in West Africa to Kenya in East
64 Africa. Genetic diversity was found to be high in most populations, but there were marked
65 patterns of population structure, and clear differences between populations in the mag-
66 nitude and architecture of genetic diversity, indicating complex and varied demographic
67 histories. However, both of these species have a large geographical range [14], and many
68 countries and ecological settings are not represented in the Ag1000G phase 1 resource.
69 Also, only SNPs were studied in Ag1000G phase 1, but other types of genetic variation
70 are known to be important. In particular, copy number variation has long been suspected
71 to play a key role in insecticide resistance [15, 16, 17], but no previous attempts to call
72 genome-wide CNVs have been made in these species.

73 This paper describes the data resource produced by the second phase of the Ag1000G
74 project. Within this phase, sampling and sequencing was expanded to include additional
75 wild-caught mosquitoes collected from five countries not represented in phase 1. This
76 includes three new locations with *An. coluzzii*, providing greater power for genetic com-
77 parisons with *An. gambiae*, and two island populations, providing a useful reference point
78 to compare against mainland populations. Seven new lab crosses are also included, provid-
79 ing a substantial resource for studying genome variation and recombination within known
80 pedigrees. In this phase we studied both SNPs and CNVs, and rebuilt a haplotype ref-
81 erence panel using all wild-caught specimens. Here we describe the data resource, and
82 use it to re-evaluate major population divisions and characterise genetic diversity. We

¹<https://www.malariagen.net/projects/ag1000g>

²<https://www.malariagen.net>

83 also illustrate the broad utility of the data by comparing geographical population struc-
84 ture between the two mosquito species to investigate evidence for differences in dispersal
85 behaviour; analyse genetic diversity within a gene in the sex-determination pathway cur-
86 rently targeted for gene drive development; and provide some preliminary insights into
87 the prevalence of different molecular mechanisms of pyrethroid resistance.

88 **Results**

89 **Population sampling and sequencing**

90 We performed whole genome sequencing of 377 individual wild-caught mosquitoes, includ-
91 ing individuals collected from 3 countries (The Gambia, Côte d'Ivoire, Ghana) and two
92 oceanic islands (Bioko, Mayotte) not represented in the previous project phase. We also
93 sequenced 152 individuals comprising parents and progeny from seven lab crosses, where
94 parents were drawn from the Ghana, Kisumu, Pimperena, Mali and Akron colonies. We
95 then combined these data with the sequencing data previously generated during phase 1
96 of the project, to create a total resource of data from 1,142 wild-caught mosquitoes (1,058
97 female, 84 male) from 13 countries (Figure 1; Table S1) and 234 mosquitoes from 11 lab
98 crosses (Table S2). As in the previous project phase, all mosquitoes were sequenced indi-
99 vidually on Illumina technology using 100 bp paired-end reads to a target depth of 30X,
100 and all 1,142 mosquitoes in the final resource had a mean depth above 14X.

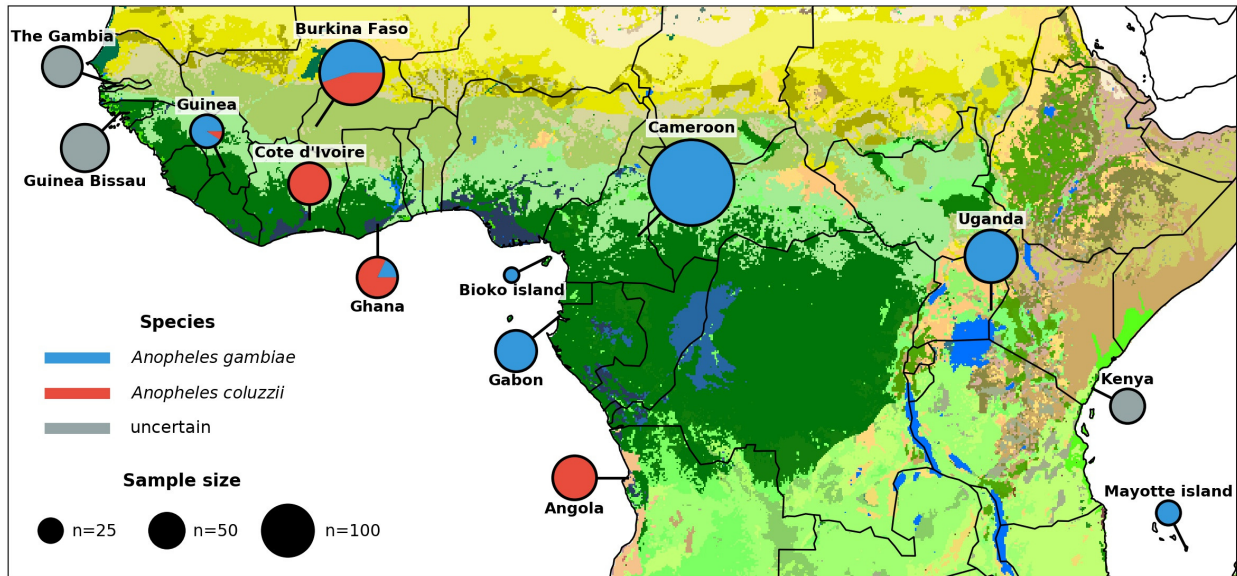


Figure 1. Ag1000G phase 2 sampling locations. Colour of circle denotes species and area represents sample size. Species assignment is labelled as uncertain for samples from Guinea-Bissau, The Gambia and Kenya, because all individuals from those locations carry a mixture of *An. gambiae* and *An. coluzzii* ancestry informative markers, see main text and Figure S1 for details. Map colours represent ecosystem classes, dark green designates forest ecosystems; see Figure 9 in [18] for a complete colour legend.

101 Genome variation

102 Sequence reads from all individuals were aligned to the AgamP3 reference genome [9, 10]
103 and SNPs were discovered using methods described previously [12]. In total, we discovered
104 57,837,885 SNPs passing all variant quality filters, 11% of which were newly discovered in
105 this project phase. Of these high quality SNPs, 24% were found to be multiallelic (three or
106 more alleles). We also analysed genome accessibility to identify all genomic positions where
107 read alignments were of sufficient quality and consistency to support accurate discovery
108 and genotyping of nucleotide variation. Similar to the previous project phase, we found
109 that 61% (140 Mbp) of genome positions were accessible, including 91% (18 Mbp) of the
110 exome and 58% (121 Mbp) of non-coding positions. Overall we discovered an average
111 of one variant allele every 1.9 bases of the accessible genome. We then used high quality
112 biallelic SNPs to construct a new haplotype reference panel including all 1,142 wild-caught
113 individuals, via a combination of read-backed phasing and statistical phasing as described
114 previously [12].

115 In this project phase we also performed a genome-wide CNV analysis, described in detail

116 elsewhere [19]. In brief, for each individual mosquito, we called CNVs by fitting a hidden
117 Markov model to windowed data on depth of sequence read coverage, then compared
118 calls between individuals to identify shared CNVs. The CNV callset comprises 31,335
119 distinct CNVs, of which 7,086 were found in more than one individual, and 1,557 were
120 present at at least 5% frequency in one or more populations. CNVs spanned more than
121 68 Mbp in total and overlapped 7,190 genes. CNVs were significantly enriched in gene
122 families associated with metabolic resistance to insecticides, with three loci in particular
123 (two clusters of cytochrome P450 genes *Cyp6p/aa*, *Cyp9k1* and a cluster of glutathione
124 S-transferase genes *Gste*) having a large number of distinct CNV alleles, multiple alleles
125 at high population frequency, and evidence that CNVs are under positive selection [19].
126 CNVs at these loci are thus likely to be playing an important role in adaptation to mosquito
127 control interventions.

128 **Species assignment**

129 The conventional and most widely used molecular assays for differentiating *An. gambiae*
130 from *An. coluzzii* are based on fixed differences in the centromeric region of the X chro-
131 mosome [20, 21]. In the first phase of the Ag1000G project, we compared the results
132 of these assays with genotypes at 506 ancestry-informative SNPs distributed across all
133 chromosome arms, and found that in some cases the conventional assays were not con-
134 cordant with species ancestry at other genome locations. In particular, all individuals
135 from two sampling locations (Kenya, Guinea-Bissau) carried a mixture of *An. gambiae*
136 and *An. coluzzii* alleles, creating uncertainty regarding the appropriate species assignment
137 [12]. Applying the same analysis to the new samples in Ag1000G phase 2, we found that
138 mosquitoes from The Gambia also carried a mixture of alleles from both species, in similar
139 proportions to mosquitoes from Guinea-Bissau (Figure S1). In all other locations, alleles
140 at ancestry-informative SNPs were concordant with conventional diagnostics [20, 21], ex-
141 cept on chromosome arm 2L where there has been a known introgression event carrying
142 an insecticide resistance allele from *An. gambiae* into *An. coluzzii* [22, 23, 24, 25]. We
143 observed this introgression in *An. coluzzii* from both Burkina Faso and Angola in the
144 phase 1 cohort, and it was also present among *An. coluzzii* from Côte d’Ivoire, Ghana and
145 Guinea in the phase 2 cohort.

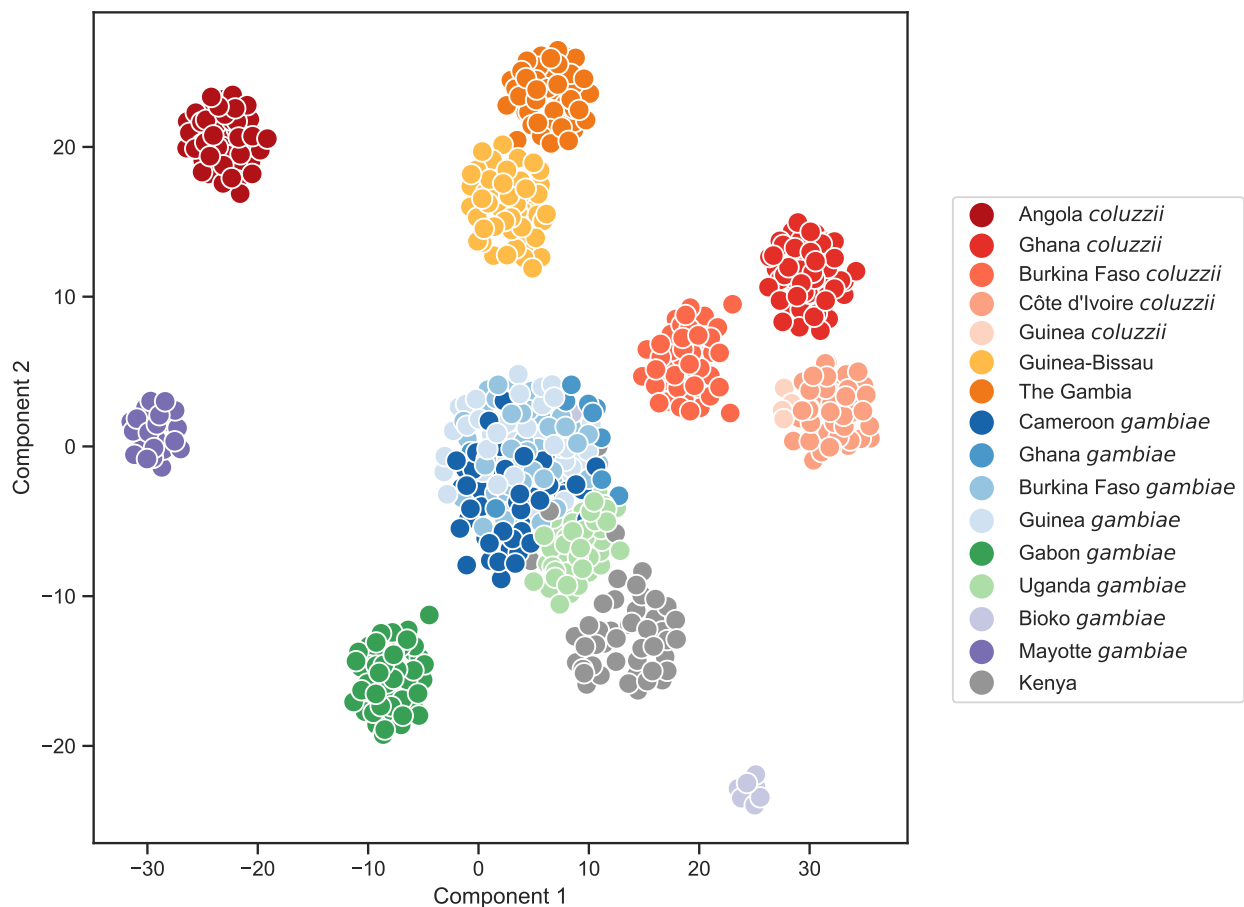


Figure 2. Population structure analysis of the wild-caught mosquitoes using UMAP [26]. Genotype data at biallelic SNPs from euchromatic regions of Chromosome 3 were projected onto two components. Each marker represents an individual mosquito. Mosquitoes from each country and species were randomly downsampled to at most 50 individuals.

146 Population structure

147 We investigated genetic population structure within the cohort of wild-caught mosquitoes
148 by performing dimensionality reduction analyses on the genome variation data, including
149 UMAP [26] and PCA [27] of biallelic SNPs from euchromatic regions of Chromosome
150 3 (Figure 2; Figure S2), and PCA of CNVs from the whole genome (Figure S3). To
151 complement these analyses, we fitted models of population structure and admixture [28]
152 to the SNP data (Figure S4). We also used SNPs to compute two measures of genetic
153 differentiation, average F_{ST} and rates of rare variant sharing, between pairs of populations
154 defined by country of origin and species (Figure 3). From these analyses, three major
155 groupings of individuals from multiple countries were evident: *An. coluzzii* from West

156 Africa (Burkina Faso, Ghana, Côte d'Ivoire, Guinea); *An. gambiae* from West, Central
157 and near-East Africa (Burkina Faso, Ghana, Guinea, Cameroon, Uganda); individuals
158 with uncertain species status from far-West Africa (Guinea-Bissau, The Gambia). Within
159 each of these groupings, samples clustered closely in all PCA and UMAP components
160 and in admixture models for up to $K = 5$ ancestral populations, and differentiation
161 between countries was weak, consistent with relatively unrestricted gene flow between
162 countries. Each of the remaining PCA clusters comprised samples from a single country
163 and species (Angola *An. coluzzii*; Gabon *An. gambiae*, Mayotte *An. gambiae*; Bioko *An.*
164 *gambiae*; individuals with uncertain species status from Kenya), and in general each of
165 these populations was more strongly differentiated from all other populations, consistent
166 with a role for geographical factors limiting gene flow. The admixture analyses for Mayotte
167 and Kenya modelled individuals from both populations as a mixture of multiple ancestral
168 populations. This could represent some true admixture in these populations' histories, but
169 could also be an artefact due to strong genetic drift [29], and requires further investigation.
170 A comparison of the two *An. gambiae* island populations is interesting because Mayotte
171 was highly differentiated from all other populations, but individuals from Bioko were
172 more closely related to other West African *An. gambiae*, suggesting that Bioko may not
173 be isolated from continental populations despite a physical separation of more than 30
174 km.

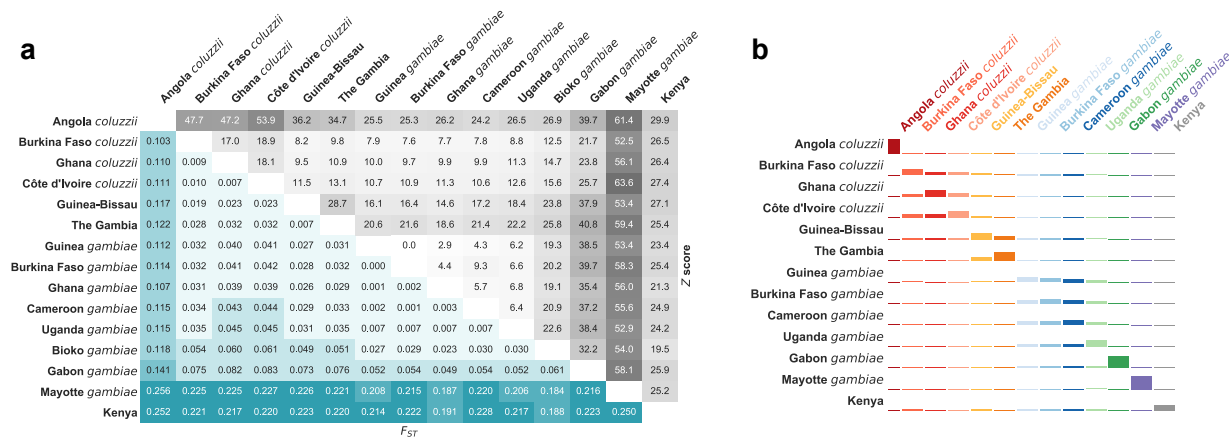


Figure 3. Genetic differentiation between populations, computed using using biallelic SNPs from euchromatic regions of Chromosome 3. **(a)** Average allele frequency differentiation (F_{ST}) between pairs of populations. The bottom left triangle shows average F_{ST} values between each population pair. The top right triangle shows the Z score for each F_{ST} value estimated via a block-jackknife procedure. **(b)** Allele sharing in doubleton (f_2) variants. For each population, we identified the set of doubletons with at least one allele originating from an individual in that population. We then computed the fraction of those doubletons shared with each other population and the fraction shared only within itself. The height of the coloured bars represent the probability of sharing a doubleton allele between or within populations. Heights are normalized row-wise for each population so that the sum of coloured bars in each row equals 1.

175 The new locations sampled in this project phase allow more comparisons to be made
 176 between *An. gambiae* and *An. coluzzii*, and there are many open questions regarding
 177 their behaviour, ecology and evolutionary history. For example, it would be valuable to
 178 know whether there are any differences in long-range dispersal behaviour between the
 179 two species [30] as have been suggested by recent studies in Sahelian regions [31, 32].
 180 Providing a comprehensive answer to this question is beyond the scope of this study, but
 181 we performed a preliminary analysis by estimating Wright’s neighbourhood size for each
 182 species [33]. This statistic is an approximation for the effective number of potential mates
 183 for an individual, and can be viewed as a measurement of how genetic differentiation
 184 between populations correlates with the geographical distance between them (isolation
 185 by distance). We used Rousset’s method for estimating neighbourhood size based on a
 186 regression of normalised F_{ST} against the logarithm of geographical distance [34]. To avoid
 187 any confounding effect of major ecological discontinuities, we used only populations from
 188 West Africa and Central Africa north of the equatorial rainforest. We found that average
 189 neighbourhood sizes are significantly lower in *An. coluzzii* than in *An. gambiae* (Wilcoxon,
 190 $W = 1320$, $P < 2.2e - 16$) (Figure 4), indicating stronger isolation by distance among

191 *An. coluzzii* populations and suggesting a lower rate and/or range of dispersal. However,
192 we do not have representation of both species at all sampling locations, and so further
193 sampling will be needed to confirm this result.

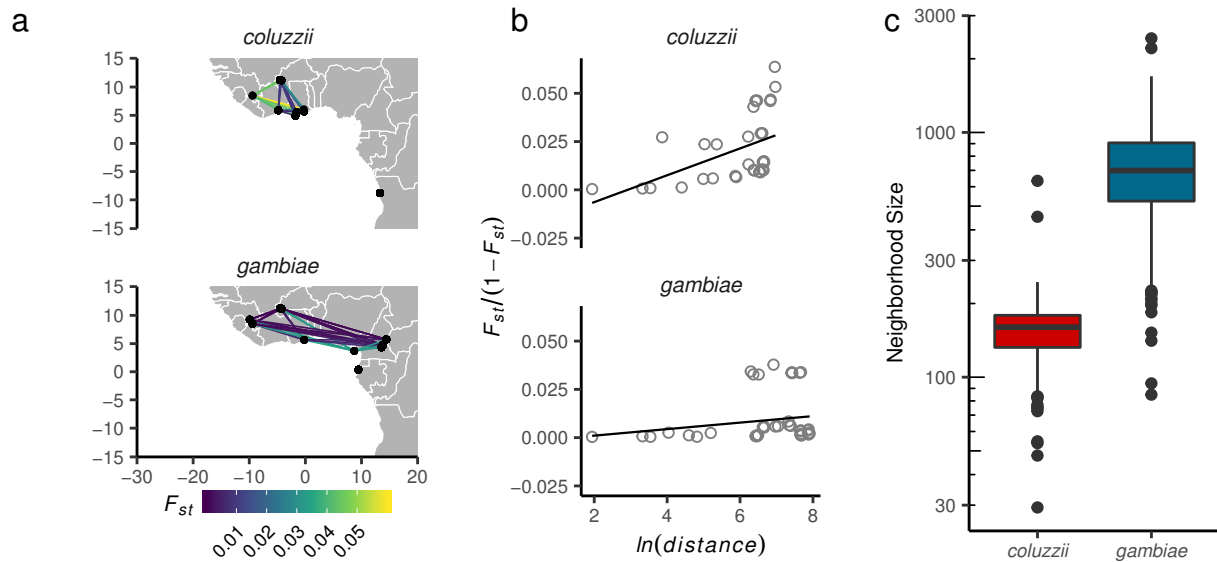


Figure 4. Comparison of isolation by distance between *An. coluzzii* and *An. gambiae* populations from locations in West and Central Africa north of the equatorial rainforest. **(a)** Study region and pairwise F_{ST} . **(b)** Regressions of average genome-wide F_{ST} against geographic distance, following Rousset [34]. Neighbourhood size is estimated as the inverse slope of the regression line. **(c)** Difference in neighbourhood size estimates by species. Box plots show medians and 95% confidence intervals of the distribution of estimates calculated in 200 kbp windows across the euchromatic regions of Chromosome 3.

194 Genetic diversity

195 The populations represented in the Ag1000G phase 2 cohort can serve as a reference point
196 for comparisons with populations sampled by other studies at other times and locations.
197 To facilitate population comparisons, we characterised genetic diversity within each of 16
198 populations in our cohort defined by country of origin and species by computing a variety
199 of summary statistics using SNP data from the whole genome. These statistics included
200 nucleotide diversity (θ_{π} ; Figure 5a), the density of segregating sites (θ_W ; Figure S5),
201 Tajima's D (Figure 5b) and site frequency spectra (SFS; Figure S6). We also estimated
202 runs of homozygosity (ROH; Figure 5c) within each individual and runs of identity by de-
203 scent (IBD; Figure 5d) between individuals, both of which provide additional information
204 about haplotype sharing and patterns of relatedness within populations.

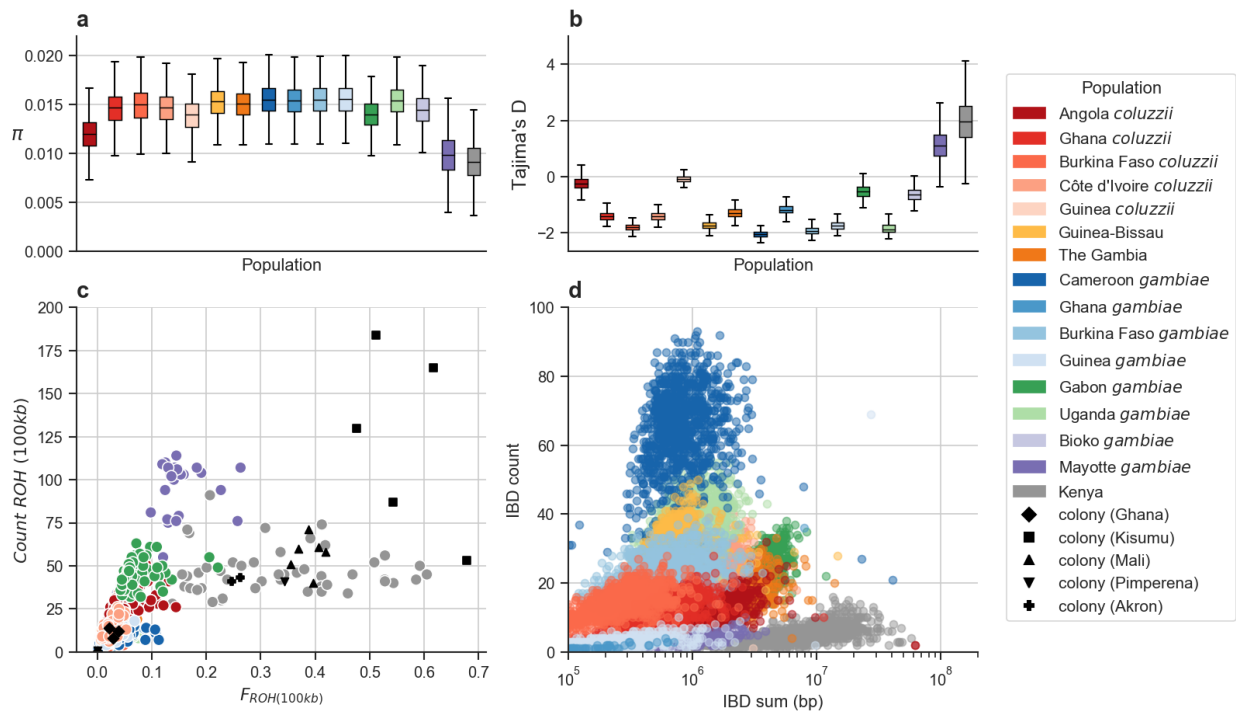


Figure 5. Genetic diversity within populations. **(a)** Nucleotide diversity (θ_{π}) calculated in non-overlapping 20 kbp genomic windows using SNPs from euchromatic regions of Chromosome 3. **(b)** Tajima's D calculated in non-overlapping 20 kbp genomic windows using SNPs from euchromatic regions of Chromosome 3. **(c)** Runs of homozygosity (ROH) in individual mosquitoes. Each marker represents an individual mosquito. **(d)** Runs of identity by descent between individuals. Each marker represents a pair of individuals drawn from the same population.

205 The two easternmost populations (Kenya, Mayotte) were outliers in all statistics calcu-
 206 lated, with lower diversity, a deficit of rare variants relative to neutral expectation, and a
 207 higher degree of haplotype sharing within and between individuals. The Kenyan popula-
 208 tion was represented in Ag1000G phase 1, and we previously described how the patterns of
 209 diversity in this population were consistent with a severe and recent population bottleneck
 210 [12]. The similarities between Kenya and Mayotte suggest that the Mayotte population
 211 has also experienced a population bottleneck, which would be expected given that May-
 212 otte is an oceanic island 310 km from Madagascar and 500 km from continental Africa,
 213 and may have been colonised by *An. gambiae* via small numbers of individuals. Although
 214 ROH and IBD were elevated in both populations, Mayotte individuals had a larger num-
 215 ber of shorter tracts than Kenyan individuals, which may reflect differences in the timing
 216 and/or strength of a bottleneck. In contrast, the *An. gambiae* individuals from Bioko
 217 Island had similar patterns of diversity to *An. gambiae* populations from West and Cen-

218 tral Africa, supporting other analyses which suggest that this population is not strongly
219 isolated from continental populations (Figures S2, 3). The additional *An. coluzzii* popu-
220 lations (Ghana, Côte d’Ivoire) were similar to the previously sampled Burkina Faso *An.*
221 *coluzzii* population, and the newly sampled Gambian population with uncertain species
222 status was similar to the previously sampled Guinea-Bissau population, consistent with
223 evidence from population structure analyses that these populations form groupings with
224 shared demographic histories and ongoing gene flow.

225 **Design of Cas9 gene drives**

226 Nucleotide variation data from this resource is being used to inform the development of
227 gene drives, a novel mosquito control technology using engineered selfish genetic elements
228 to cause mosquito population suppression or modification [35, 36, 37, 38, 8]. Promising
229 results have been obtained with a Cas9 homing endonuclease gene drive targeting a locus in
230 the doublesex gene (*dsx*), which is a critical component of the sex determination pathway
231 [8]. This locus was chosen in part because it has extremely low genetic diversity both
232 within and between species in the *An. gambiae* complex [12]. Low diversity is required
233 because any natural variation within the target sequence could inhibit association with
234 the Cas9 guide RNA and cause resistance to the gene drive [39]. We reviewed nucleotide
235 variation within *dsx* using the expanded cohort of wild-caught samples in the phase 2
236 cohort, and found no new nucleotide variants within the sequence targeted for Cas9 gene
237 drive, other than the previously known SNP at 2R:48,714,641, which has been shown not
238 to interfere with the gene drive process in lab populations [8]. To facilitate the search for
239 other potential gene drive targets in *dsx* and other genes, we computed allele frequencies
240 for all SNPs in all populations and included those data in the resource. We also compiled
241 a table of all potential Cas9 target sites (23 bp regions with a protospacer-adjacent motif)
242 in the genome that overlap a gene exon. This table includes a total of 20 Cas9 targets that
243 overlap *dsx* exon 5 and that contain at most one SNP within the Ag1000G phase 2 cohort
244 (Figure 6). Thus there may be multiple viable targets for gene drives disrupting the sex
245 determination pathway, providing opportunities to mitigate the impact of resistance due
246 to variation within any single target.

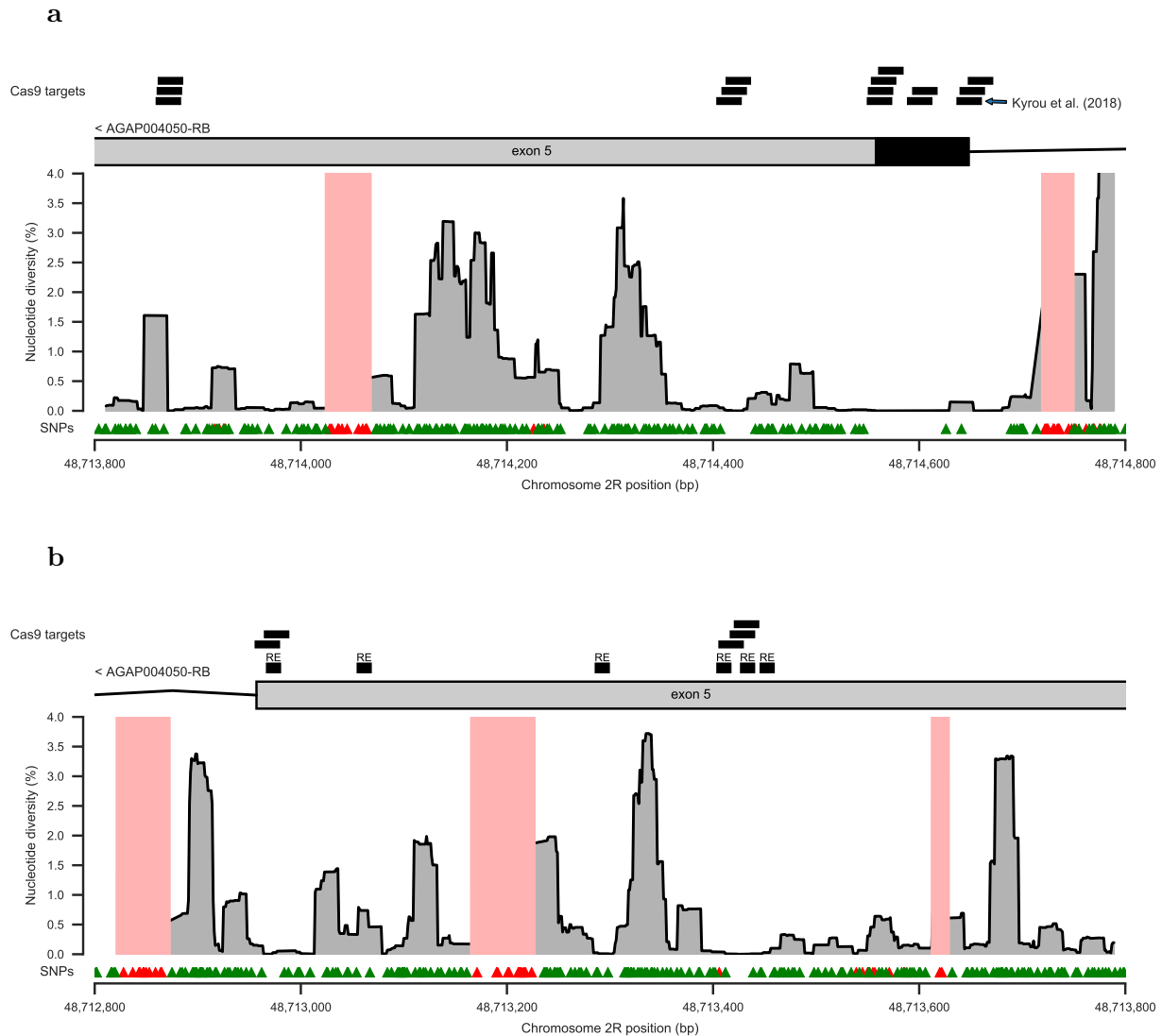


Figure 6. Nucleotide diversity within the female-specific exon 5 of the doublesex gene (*dsx*; AGAP004050), a key component of the sex determination pathway and a gene targeted for Cas9-based homing endonuclease gene drive [8]. In both plots, the location of exon 5 within the female-specific isoform (AGAP004050-RB; AgamP4.12 gene set) is shown above (black = coding sequence; grey = untranslated region), with additional annotations above to show the location of viable Cas9 target sequences containing at most 1 SNP, and the putative exon splice enhancing sequences (“RE”) reported in [40]. The main region of the plot shows nucleotide diversity averaged across all Ag1000G phase 2 populations, computed in 23 bp moving windows. Regions shaded pale red indicate regions not accessible to SNP calling. Triangle markers below show the locations of SNPs discovered in Ag1000G phase 2 (green = passed variant filters; red = failed variant filters). (a) exon5/intron4 boundary. (b) exon5/intron6 boundary.

247 The presence of highly conserved regions within *dsx* also provides an example of how
248 genetic variation data from natural populations can be relevant to the study of fundamental
249 molecular processes such as sex determination. The region of conservation containing the
250 Cas9 target site in fact extends over 200 bp, including 50 bp of untranslated sequence

251 within exon 5, the entire coding sequence of exon 5, and 50 bp of intron 4 (Figure 6a).
252 Such conservation of both coding and non-coding sites suggests that purifying selection
253 is acting here on the nucleotide sequence and not just on the protein sequence. This in
254 turn suggests that the nucleotide sequence serves as an important target for factors that
255 bind to DNA or pre-mRNA molecules. This is plausible because sex determination in
256 insects depends on sex-specific splicing of *dsx*, with exon 5 being included in the female
257 transcript and excluded in the male transcript [41]. The upstream regulatory factors that
258 control this differential splicing are not known in *An. gambiae* [40, 42], but in *Drosophila*
259 *melanogaster* it has been shown that female-specific factors bind to regulatory sequences
260 (*dsxREs*) within the exon 5 region of the *dsx* pre-mRNA and promote inclusion of exon
261 5 within the final transcript [43, 41]. Putative homologs of these (*dsxRE*) sequences are
262 present in *An. gambiae* [40], and five out of six *dsxREs* are located in tracts of near-
263 complete nucleotide conservation in our data, consistent with purifying selection due to
264 pre-mRNA binding (Figure 6b). However, the 200 bp region of conservation spanning
265 the intron 4/exon 5 boundary targeted for Cas9 gene drive remains mysterious, because
266 it is more than 1 kbp distant from any of these putative regulatory sequences. Overall
267 these data add further evidence for fundamental differences in the molecular biology of
268 sex determination between *Anopheles* and *Drosophila* and provide new clues for further
269 investigation of the molecular pathway upstream of *dsx* in *An. gambiae* [40, 42].

270 **Resistance to pyrethroid insecticides**

271 Malaria control in Africa depends heavily on mass distribution of long-lasting insecticidal
272 bed-nets (LLINs) impregnated with pyrethroid insecticides [44, 45, 46]. Entomological
273 surveillance programs regularly test malaria vector populations for pyrethroid resistance
274 using standardised bioassays, and these data have shown that pyrethroid resistance has
275 become widespread in *An. gambiae* [2, 3]. However, pyrethroid resistance can be con-
276 ferred by different molecular mechanisms, and it is not well understood which molecular
277 mechanisms are responsible for resistance in which mosquito populations. The nucleotide
278 variation data in this resource include 66 non-synonymous SNPs within the *Vgsc* gene that
279 encodes the binding target for pyrethroid insecticides, of which two SNPs (L995F, L995S)
280 are known to confer a pyrethroid resistance phenotype, and one SNP (N1570Y) has been

281 shown to substantially increase pyrethroid resistance when present in combination with
282 L995F [47]. These SNPs can serve as markers of target-site resistance to pyrethroids, but
283 knowledge of genetic markers of metabolic resistance in *An. gambiae* and *An. coluzzii* is
284 currently limited [48, 49]. Metabolic resistance to pyrethroids is mediated at least in part
285 by increased expression of cytochrome P450 (CYP) enzymes [50, 51, 52, 53], and we found
286 CNV hot-spots at two loci containing *Cyp* genes [19]. One of these loci occurs on chromo-
287 some arm 2R and overlaps a cluster of 10 *Cyp* genes, including *Cyp6p3* previously shown
288 to metabolise pyrethroids [54] and recently shown to confer pyrethroid resistance when
289 expression is increased in *An. gambiae* using the GAL4/UAS transgenic system [55]. The
290 second locus occurs on the X chromosome and spans a single *Cyp* gene, *Cyp9k1*, which has
291 also been shown to metabolise pyrethroids [53]. At each of these two loci we found a re-
292 markable allelic heterogeneity, with at least 15 distinct CNV alleles, several of which were
293 present in over 50% of individuals in some populations and were associated with signatures
294 of positive selection [19]. We also found CNVs at two other *Cyp* gene loci on chromosome
295 arm 3R containing genes previously associated with pyrethroid resistance, *Cyp6z1* [56] and
296 *Cyp6m2* [57], although there was only a single CNV allele at each locus. Overexpression
297 of *Cyp6m2* has been shown to confer resistance to pyrethroids but increased susceptibility
298 to the organophosphate malathion [55], and so the selection pressures at this locus may
299 be more complex. The precise phenotype of these CNVs remains to be characterised, but
300 given the multiple lines of evidence showing that increased expression of genes at these
301 loci confers pyrethroid resistance, it seems reasonable to assume that CNVs at these loci
302 can serve as a molecular marker of CYP-mediated metabolic resistance to pyrethroids.

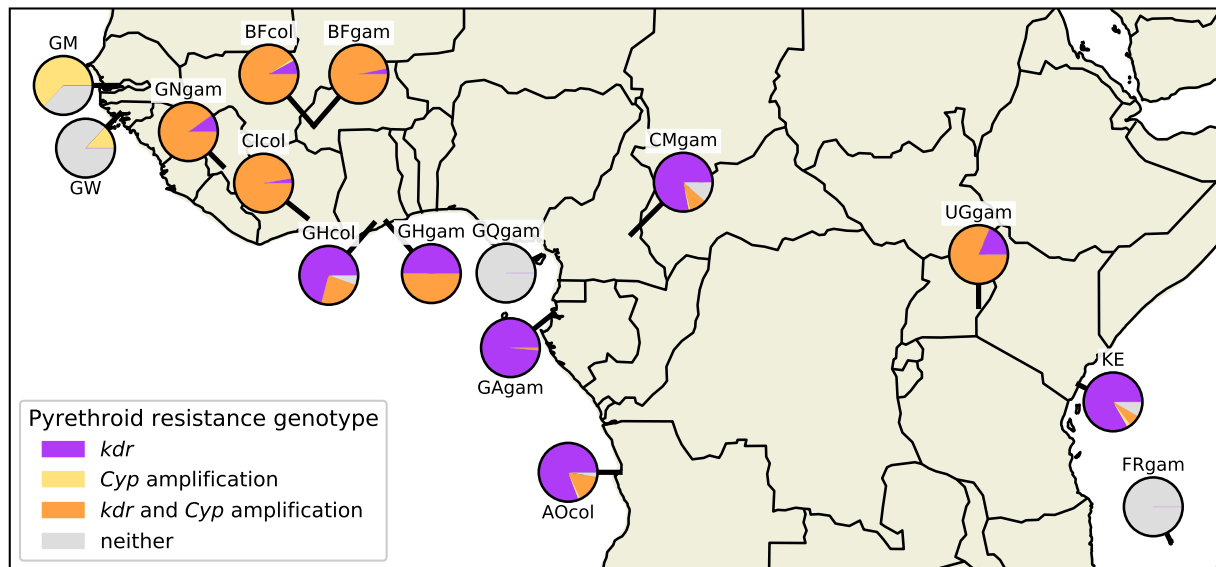


Figure 7. Pyrethroid resistance genotype frequencies. The geographical distribution of pyrethroid insecticide resistance genotypes are shown by population. Pie chart colours represent resistance genotype frequencies: purple - these individuals were either homozygous or heterozygous for one of the two *kdr* pyrethroid target site resistance alleles *Vgsc*-L995F/S; yellow - these individuals carried a copy number amplification within any of the *Cyp6p/aa*, *Cyp6m*, *Cyp6z* or *Cyp9k* gene clusters, but no *kdr* alleles; orange - these individuals carried at least one *kdr* allele and one *Cyp* gene amplification; grey - these individuals carried no known pyrethroid resistance alleles (no *kdr* alleles or *Cyp* amplifications). The Guinea *An. coluzzii* population is omitted due to small sample size.

303 We constructed an overview of the prevalence of these two pyrethroid resistance mecha-
304 nisms – target-site resistance and CYP-mediated metabolic resistance – within the Ag1000G
305 phase 2 cohort by combining the data on nucleotide and copy number variation (Figure
306 7). The sampling of these populations was conducted at different times in different loca-
307 tions, and the geographical sampling is relatively sparse, so we cannot draw any general
308 conclusions about the current distribution of resistance from our data. However, some pat-
309 terns were clear. For example, West African populations of both species (Burkina Faso,
310 Guinea, Côte d’Ivoire) all had more than 84% of individuals carrying both target-site
311 and metabolic resistance markers. In Ghana, Cameroon, Gabon and Angola, target-site
312 resistance was nearly fixed in all populations, but metabolic resistance markers were at
313 lower frequencies, and the samples from Bioko Island carried no resistance markers at all.
314 The Bioko samples were collected in 2002, and so the lack of resistance is likely due to
315 the fact that sampling predated any major scale-up of vector control interventions [53].
316 However, the Gabon samples were collected in 2000, and show that high levels of target-

317 site resistance were present in some populations at that time. In the "Far West" (Guinea
318 Bissau, The Gambia) [58], target-site resistance was absent, but *Cyp* gene amplifications
319 were present, and thus surveillance using only molecular assays that detect target site
320 resistance at those locations could be missing an important signal of metabolic resistance.
321 In East Africa, both Kenya and Uganda had high frequencies of target-site resistance (88%
322 and 100% respectively). However, 81% of Uganda individuals also had *Cyp* gene amplifi-
323 cations, whereas only 4% of Kenyans (two individuals) carried these metabolic resistance
324 markers. Denser spatio-temporal sampling and sequencing will enable us to build a more
325 complete picture of the prevalence and spread of these different resistance mechanisms,
326 and would be highly relevant to the design of insecticide resistance management plans.

327 **Discussion**

328 **Insecticide resistance surveillance**

329 The Ag1000G phase 2 data resource incorporates both nucleotide and copy number vari-
330 ation from the whole genomes of 1,142 mosquitoes collected from 13 countries spanning
331 the African continent. These data provide a battery of new genetic markers that can be
332 used to expand our capabilities for molecular surveillance of insecticide resistance. In-
333 secticide resistance management is a major challenge for malaria vector control, but the
334 availability of new vector control products is opening up new possibilities. However, new
335 products may be more expensive than products currently in use, so procurement decisions
336 have to be justified, and resources targeted to areas where they will have the greatest im-
337 pact. For example, next-generation LLINs are now available which combine a pyrethroid
338 insecticide with either a second insecticide or a synergist compound, piperonyl butoxide
339 (PBO), which partially ameliorates metabolic resistance by inhibiting CYP enzyme activ-
340 ity in the mosquito. However, CYP-mediated metabolic resistance is only one of several
341 possible mechanisms of pyrethroid resistance that may or may not be present in vector
342 populations being targeted. It would therefore be valuable to survey mosquito popula-
343 tions and determine the prevalence of different pyrethroid resistance mechanisms, both
344 before and after any change in vector control strategy. Our data resource includes CNVs
345 at four *Cyp* gene loci (*Cyp6p/aa*, *Cyp6m*, *Cyp6z* and *Cyp9k*) which could serve as molec-

346 ular markers of CYP-mediated metabolic resistance. Glutathione S-transferase enzymes
347 are also associated with metabolic resistance to pyrethroids [59, 55] as well as to other
348 insecticide classes [48, 60, 61, 55] and we found CNVs at the *Gste* locus which could serve
349 as molecular markers of this alternative resistance mechanism, which is not inhibited by
350 PBO. *Gste* CNVs were less prevalent in our dataset than *Cyp* CNVs, and the geographical
351 distribution also differed, suggesting they may be driven by different selection pressures
352 (Figure S7). Further work is needed to characterise the resistance phenotype associated
353 with these CNVs, but the allelic heterogeneity, the high population frequencies, and the
354 evidence for positive selection observed in our data, coupled with previous gene expres-
355 sion and functional studies [50, 51, 52, 53, 55], all support a metabolic role in insecticide
356 resistance.

357 To illustrate the potential for improved molecular surveillance of pyrethroid resistance,
358 we combined the data on known SNP markers of target-site resistance and the novel puta-
359 tive CNV markers of CYP-mediated metabolic resistance, and computed the frequencies
360 of these different resistance mechanisms in the populations we sampled (Figure 7). There
361 are clear heterogeneities, with some populations at high frequency for both resistance
362 mechanisms, particularly in West Africa. The presence of CYP-mediated pyrethroid resis-
363 tance in a population suggests that PBO LLINs might provide some benefit over standard
364 LLINs. However, if other resistance mechanisms are also at high frequency, the benefit of
365 the PBO synergist might be diminished. Current WHO guidance states that PBO LLINs
366 are recommended in regions with “intermediate levels” of pyrethroid resistance, but not
367 where resistance levels are high [62]. This guidance is based on modelling of bioassay data
368 and experimental hut trials, and it is not clear why PBO LLINs are predicted to provide
369 diminishing returns at higher resistance levels, although high levels of resistance presum-
370 ably correlate with the presence of multiple resistance mechanisms, including mechanisms
371 not inhibited by PBO [45]. Without molecular data, however, this guidance is hard to
372 evaluate or improve upon.

373 Ideally, molecular data on insecticide resistance mechanisms would be collected as part of
374 routine entomological surveillance, as well as in field trials of new vector control products,
375 alongside data from bioassays and other standard entomological monitoring procedures.
376 There are several options for scaling up surveillance of new genetic markers, including

377 both whole genome sequencing and targeted (amplicon) sequencing with several choices of
378 sequencing technology platform, as well as various PCR-based assays. Assays that target
379 specific genetic loci are attractive in the short term, because of the low cost and infras-
380 tructure requirements, and data from the Ag1000G project have been used successfully to
381 design multiplex assays for the Agena Biosciences iPLEX platform [63] and for Illumina
382 amplicon sequencing (manuscript in preparation). But targeted assays would need to be
383 updated regularly to ensure all current forms of insecticide resistance are covered, and to
384 capture new forms of resistance as they emerge. None of the samples sequenced in this
385 study were collected more recently than 2012, geographical sampling within each country
386 was limited, and many countries are not yet represented in the resource, therefore there
387 remain important gaps to be filled. The next phase of the Ag1000G project will expand
388 the resource to cover 18 countries, and will include another major malaria vector, *An.*
389 *arabiensis*, in addition to *An. gambiae* and *An. coluzzii*, and so will address some of these
390 gaps. Looking beyond the Ag1000G project, genomic surveillance of insecticide resistance
391 will require new sampling frameworks that incorporate spatial and ecological modelling
392 of vector distributions to improve future collections and guide sampling at appropriate
393 spatial scales [64]. To keep pace with vector populations, regular whole genome sequenc-
394 ing of contemporary populations from a well-chosen set of sentinel sites will be needed.
395 Fortunately mosquitoes are easy to transport, and the costs of whole genome sequencing
396 continue to fall, so it is reasonable to consider a mixed strategy that includes both whole
397 genome sequencing and targeted assays.

398 **Gene flow**

399 These data also cast some new, and in some cases contrasting, light on the question of gene
400 flow between malaria vector populations. The question is of practical interest because gene
401 flow is enabling the spread of insecticide resistance between species and across large geo-
402 graphical distances [12, 65]. This gene flow also needs to be quantified and modelled before
403 new vector control interventions based on the release of genetically modified mosquitoes
404 could be considered [66]. We found evidence that isolation by distance is greater for *An.*
405 *coluzzii* than for *An. gambiae*, at least within West Africa, suggesting that the effective
406 rate of migration could be lower in *An. coluzzii*. This result was supported by population

407 structure analyses, where all *An. coluzzii* individuals were clearly clustered by country in
408 the UMAP analysis, whereas *An. gambiae* individuals from Guinea, Burkina Faso, Ghana
409 and Cameroon could not be separated in any of the UMAP, PCA or admixture analy-
410 ses. A variety of anopheline species have recently been found to engage in long-distance
411 wind-assisted migration, including *An. coluzzii* but not *An. gambiae*, which would appear
412 to contradict our results, although the study was limited to a single location within the
413 Sahelian region [32]. If *An. coluzzii* does have a lower rate and/or range of dispersal than
414 *An. gambiae*, this is clearly not limiting the spread of insecticide resistance adaptations
415 between countries. For example, among the CNV alleles we discovered at the *Cyp6p/aa*,
416 *Cyp9k1* and *Gste* loci, 7/13 alleles found in *An. coluzzii* had spread to more than one
417 country, compared with 8/27 alleles in *An. gambiae* [19]. There is also an interesting con-
418 trast between the spread of pyrethroid target-site and metabolic resistance alleles. Our
419 previous analysis of haplotypes carrying target-site resistance alleles in the Ag1000G phase
420 1 cohort found that resistance haplotypes had spread to countries spanning the equatorial
421 rainforest and the Rift valley, and had moved between *An. gambiae* and *An. coluzzii*
422 [12, 65]. In the most extreme example, one haplotype (F1) had spread to countries as
423 distant as Guinea and Angola. In contrast, although CNV alleles were commonly found
424 in multiple countries, we did not observe any cases of CNV alleles crossing any of these
425 ecological or biological boundaries, apart from a single allele found in both Gabon and
426 Cameroon *An. gambiae* (*Gste* Dup5). There are multiple possible explanations for this
427 difference, including differences in the strength, timing or spatial distribution of selective
428 pressures, or intrinsic factors such as differences in fitness costs in the absence of posi-
429 tive selection. Further work is required to investigate the selective forces and biological
430 constraints affecting the spread of these different modes of adaptation to insecticide use.

431 The two island populations sampled in this project phase also provide an interesting
432 contrast. Samples from Mayotte are highly differentiated from mainland *An. gambiae*
433 and have patterns of reduced genetic diversity, consistent with a reduction in population
434 size and strong isolation. Bioko samples, on the other hand, are closely related to West
435 African *An. gambiae*, and have comparable levels of genetic diversity, suggesting ongoing
436 gene flow. Bioko is part of Equatorial Guinea administratively, and there are frequent fer-
437 ries to the mainland, which could provide opportunities for mosquito movement. However,

438 there are no pyrethroid resistance alleles in our Bioko samples and these were collected in
439 2002 at a time when target-site resistance alleles were present in mainland populations, so
440 the rate of contemporary migration between Bioko and mainland populations remains an
441 open question. A recent study of *An. gambiae* populations on the Lake Victoria islands,
442 separated from mainland Uganda by 4-50 km, found evidence for isolation between is-
443 land and mainland populations, as well as between individual islands [67]. However, some
444 selective sweeps at insecticide resistance loci had spread through both mainland and is-
445 land populations, thus isolation is not complete and some contemporary gene flow occurs.
446 Resolving these gene flow questions and apparent contradictions will require fitting quanti-
447 tative models of contemporary migration to genomic data. We previously fitted migration
448 models to pairs of populations using site frequency spectra, but the approach provides poor
449 resolution to differentiate recent from ancient migration rates [12]. In general, methods
450 that leverage information about haplotype sharing within and between populations should
451 provide the greatest resolution to disentangle ancient from recent demographic events, as
452 well as providing independent estimates for both migration rates and population densities.
453 There is promising recent work in this direction [68], but models have so far only been
454 applied to data from human populations. The haplotype data we have generated should
455 prove a useful resource for further work to evaluate whether these models can be applied
456 to malaria vector populations with sufficient accuracy to support real-world planning of
457 new vector control interventions.

458 **Conclusions**

459 Malaria has become a stubborn foe, frustrating global efforts towards elimination in both
460 low and high burden settings. However, new vector control tools offer hope, as does
461 the renewed focus on improving surveillance systems and using data to tailor interven-
462 tions. The genomic data resource we have generated provides a platform from which to
463 accelerate these efforts, demonstrating the potential for data integration on a continental
464 scale. Nevertheless, work remains to fill gaps in these data, by expanding geographical
465 coverage, including other malaria vector species and integrating genomic data collection
466 with routine surveillance of contemporary populations using quantitative sampling design.
467 We hope that the MalariaGEN data-sharing community and framework for international

468 collaboration can continue to serve as a model for coordinated action.

469 **Methods**

470 **Population sampling**

471 Ag1000G phase 2 mosquitoes were collected from natural populations at 33 sites in 13
472 sub-Saharan African countries (Figure 1 & Table S1). Throughout, we use species nomen-
473 clature following Coetzee *et al.* [13]; prior to Coetzee *et al.*, *An. gambiae* was known as
474 *An. gambiae sensu stricto* (S form) and *An. coluzzii* was known as *An. gambiae sensu*
475 *stricto* (M form). Details of the eighteen collection sites novel to Ag1000G phase 2 (dates,
476 collection and DNA extraction methods) can be found below. Information pertaining to
477 the collection of samples released as part of Ag1000G phase 1 can be found in the supple-
478 mentary information of [12]. Unless otherwise stated, the DNA extraction method used
479 for the collections described below was Qiagen DNeasy Blood and Tissue Kit (Qiagen
480 Science, MD, USA).

481 **Côte d’Ivoire:** Tiassalé (5.898, -4.823) is located in the evergreen forest zone of south-
482 ern Côte d’Ivoire. The primary agricultural activity is rice cultivation in irrigated fields.
483 High malaria transmission occurs during the rainy seasons, between May and November.
484 Samples were collected as larvae from irrigated rice fields by dipping between May and
485 September 2012. All larvae were reared to adults and females preserved over silica for
486 DNA extraction. Specimens from this site were all *An. coluzzii*, determined by PCR assay
487 [21].

488 **Bioko:** Collections were performed during the rainy season in September, 2002 by
489 overnight CDC light traps in Sacriba of Bioko island (3.7, 8.7). Specimens were stored
490 dry on silica gel before DNA extraction. Specimens contributed from this site were
491 *An. gambiae* females, genotype determined by two assays [69, 70]. All specimens had
492 the $2L^{+a}/2L^{+a}$ karyotype as determined by the molecular PCR diagnostics [71]. These
493 mosquitoes represent a population that inhabited Bioko Island before a comprehensive
494 malaria control intervention initiated in February 2004 [72]. After the intervention *An.*
495 *gambiae* was declining, and more recently almost only *An. coluzzii* can be found [73].

496 **Mayotte:** Samples were collected as larvae during March-April 2011 in temporary

497 pools by dipping, in Bouyouni (-12.738, 45.143), M'Tsamboro Forest Reserve (-12.703,
498 45.081), Combani (-12.779, 45.143), Mtsanga Charifou (-12.991, 45.156), Karihani Lake
499 forest reserve (-12.797, 45.122), Mont Benara (-12.857, 45.155) and Sada (-12.852, 45.104)
500 in Mayotte island. Larvae were stored in 80% ethanol prior to DNA extraction. All
501 specimens contributed to Ag1000G phase 2 were *An. gambiae* [70] with the standard
502 $2L^{+a}/2L^{+a}$ or inverted $2L^a/2L^a$ karyotype as determined by the molecular PCR diagnos-
503 tics [71]. The samples were identified as males or females by the sequencing read coverage
504 of the X chromosome using LookSeq [74].

505 **The Gambia:** Indoor resting female mosquitoes were collected by pyrethrum spray
506 catch from four hamlets around Njabakunda (-15.90, 13.55), North Bank Region, The
507 Gambia between August and October 2011. The four hamlets were Maria Samba Nyado,
508 Sare Illo Buya, Kerr Birom Kardo, and Kerr Sama Kuma; all are within 1 km of each
509 other. This is an area of unusually high rates of apparent hybridization between *An.*
510 *gambiae s.s.* and *An. coluzzii* [75, 76]. Njabakunda village is approximately 30 km to the
511 west of Farafenni town and 4 km away from the Gambia River. The vegetation is a mix
512 of open savannah woodland and farmland.

513 **Ghana:** Mosquitoes were collected from Twifo Praso (5.609, -1.549), a peri-urban com-
514 munity located in semi-deciduous forest in the Central Region of Ghana. It is an extensive
515 agricultural area characterised by small-scale vegetable growing and large-scale commer-
516 cial farms such as oil palm and cocoa plantations. Mosquito samples were collected as
517 larvae from puddles near farms between September and October, 2012. Madina (5.668,
518 -0.219) is a suburb of Accra within the coastal savanna zone of Ghana. It is an urban
519 community characterised by numerous vegetable-growing areas. The vegetation consists
520 of mainly grassland interspersed with dense short thickets often less than 5 m high with
521 a few trees. Specimens were sampled from puddles near roadsides and farms between
522 October and December 2012. Takoradi (4.912, -1.774) is the capital city of Western Re-
523 gion of Ghana. It is an urban community located in the coastal savanna zone. Mosquito
524 samples were collected from puddles near road construction and farms between August
525 and September 2012. Koforidua (6.094, -0.261) is the capital city of Eastern Region of
526 Ghana and is located in semi-deciduous forest. It is an urban community characterized
527 by numerous small-scale vegetable farms. Samples were collected from puddles near road

528 construction and farms between August and September 2012. Larvae from all collection
529 sites were reared to adults and females preserved over silica for DNA extraction. Both
530 *An. gambiae* and *An. coluzzii* were collected from these sites, determined by PCR assay
531 [21].

532 **Guinea-Bissau:** Mosquitoes were collected in October 2010 using indoor CDC light
533 traps, in the village of Safim (11.957, -15.649), ca. 11 km north of Bissau city, the capital
534 of the country. Malaria is hyperendemic in the region and transmitted by members of
535 the *Anopheles gambiae* complex [77]. *An. arabiensis*, *An. melas*, *An. coluzzii* and *An.*
536 *gambiae*, as well as apparent hybrids between the latter two species, are known to occur
537 in the region [78, 77]. Mosquitoes were preserved individually on 0.5ml micro-tubes filled
538 with silica gel and cotton. DNA extraction was performed by a phenol-chloroform protocol
539 [79].

540 **Lab crosses**

541 The Ag1000G phase 2 data release includes the genomes of seven additional lab colony
542 crosses, both parents and offspring (Table S2): cross 18-5 (Ghana mother x Kisumu/G3
543 father, 20 offspring); 37-3 (Kisumu x Pimperena, 20 offspring); 45-1 (Mali x Kisumu, 20
544 offspring); 47-6 (Mali x Kisumu, 20 offspring); 73-2 (Akron x Ghana, 19 offspring); 78-
545 2 (Mali x Kisumu/Ghana, 19 offspring); 80-2 (Kisumu x Akron, 20 offspring). Father
546 colonies with two names, e.g., "Kisumu/G3", signify that the father is from one of these
547 two colonies, but exactly which one is unknown. The colony labels, e.g., "18-5", are
548 identifiers used for each of the crosses within the project and have no particular meaning.
549 Information pertaining to the crosses released as part of Ag1000G phase 1 can be found in
550 the supplementary information of [12] as well as methods for cross creation and processing
551 that also apply to the crosses in phase 2.

552 **Whole genome sequencing**

553 Sequencing was performed on the Illumina HiSeq 2000 platform at the Wellcome Sanger
554 Institute. Paired-end multiplex libraries were prepared using the manufacturer's proto-
555 col, with the exception that genomic DNA was fragmented using Covaris Adaptive Fo-
556 cused Acoustics rather than nebulization. Multiplexes comprised 12 tagged individual

557 mosquitoes and three lanes of sequencing were generated for each multiplex to even out
558 variations in yield between sequencing runs. Cluster generation and sequencing were un-
559 dertaken per the manufacturer's protocol for paired-end 100 bp sequence reads with insert
560 size in the range 100-200 bp. Target coverage was 30X per individual.

561 **Genome accessibility**

562 For various population-genomic analyses, it is necessary to have a map of which positions
563 in the reference genome can be considered accessible, at which we can confidently call
564 nucleotide variation. For Ag1000G phase 2, we repeated the phase 1 genome accessibility
565 analyses [12] with 1,142 samples and the additional Mendelian error information provided
566 by the 11 crosses (in phase 1 there were four crosses). These analyses constructed a number
567 of annotations for each position in the reference genome, based on data from sequence read
568 alignments from all wild-caught samples, and additional data from repeat annotations.
569 These annotations were then analysed for their association with rates of variants with
570 one or more Mendelian errors in the crosses. Annotations and thresholds were chosen
571 to remove classes of variants that were enriched for Mendelian errors. Following these
572 analyses it was apparent that the accessibility classifications used in Ag1000G phase 1 were
573 also appropriate in application to phase 2. Reference genome positions were classified as
574 accessible if: Not repeat masked by DUST; No Coverage $\leq 0.1\%$ (at most 1 individual
575 had zero coverage); Ambiguous Alignment $\leq 0.1\%$ (at most 1 individual had ambiguous
576 alignments); High Coverage $\leq 2\%$ (at most 20 individuals had more than twice their
577 genome-wide average coverage); Low Coverage $\leq 10\%$ (at most 114 individuals had less
578 than half their genome-wide average coverage); Low Mapping Quality $\leq 10\%$ (at most
579 114 individuals had average mapping quality below 30).

580 We performed additional analyses to verify that there was no significant bias towards
581 one species or another given the use of a single reference genome AgamP3 [9] for alignment
582 of reads from all individuals. We found that the genomes of *An. coluzzii* and *An. gambiae*
583 individuals were similarly diverged from the reference genome (Fig. S8). The similarity in
584 levels of divergence is likely to reflect the mixed ancestry of the PEST strain from which
585 the reference genome was derived [9, 10]. An exception to this was the pericentromeric
586 region of the X chromosome, a known region of divergence between the two species [12]

587 where the reference genome is closer to *An. coluzzii* than to *An. gambiae*. The similarity
588 of this region to *An. coluzzii* may be due to artificial selection for the X-linked pink eye
589 mutation in the reference strain [9], as this originated in the *An. coluzzii* parent it may
590 have led to the removal of any *An. gambiae* ancestry in this region.

591 **Sequence analysis and variant calling**

592 SNP calling methods were unchanged from phase 1 of the Anopheles 1000 genomes project
593 [12]. Briefly, sequence reads were aligned to the AgamP3 reference genome [9, 10] using
594 **bwa** version 0.6.2, duplicate reads marked [80], reads realigned around putative indels,
595 and SNPs discovered using GATK version 2.7.4 Unified Genotyper following best practice
596 recommendations [81].

597 **Sample quality control**

598 A total of 1,285 individual mosquitoes were sequenced as part of Ag1000G phase 2 and
599 included in the cohort for variant discovery. After variant discovery, quality-control (QC)
600 steps using coverage and contamination filters alongside principal component analysis and
601 metadata concordance were performed to exclude individuals with poor quality sequence
602 and/or genotype data as detailed in [12]. A total of 143 individuals were excluded at this
603 stage, retaining 1,142 individuals for downstream analyses. Any SNPs with variant alleles
604 found only in excluded samples were then also excluded.

605 **Variant Filtering**

606 Following Ag1000G phase 1 [12], we applied the following SNP filters to reduce the number
607 of false SNP discoveries. We filtered any SNP that occurred at a genome position classified
608 as inaccessible as described in the section on genome accessibility above, thus removing
609 SNPs with evidence for excessively high or low coverage or ambiguous alignment. We
610 then applied additional filters using variant annotations produced by GATK based on an
611 analysis of Mendelian error in all 11 crosses present in phase 2 and Ti/Tv ratio, similar to
612 that described above for the genome accessibility analysis. We filtered any SNP that failed
613 any of the following criteria: QD <5; FS >100; ReadPosRankSum <-8; BaseQRankSum
614 <-50.

615 Of 105,486,698 SNPs reported in the raw callset, 57,837,885 passed all quality filters,
616 13,760,984 (23.8%) of which were multi-allelic (three or more non-reference alleles). To
617 produce an analysis-ready VCF file for each chromosome arm, we first removed all non-
618 SNP variants. We then removed genotype calls for individuals excluded by the sample
619 QC analysis described above, then removed any variants that were no longer variant after
620 excluding individuals. We then added INFO annotations with genome accessibility metrics
621 and added FILTER annotations per the criteria defined above. Finally, we added INFO
622 annotations with information about functional consequences of mutations using SNPEFF
623 version 4.1b [82].

624 **Haplotype estimation**

625 Haplotype estimation, also known as phasing, was performed on all phase 2 wild-caught
626 individuals using unchanged methodology from phase 1 of the Anopheles 1000 genomes
627 project [12]. In short, SHAPEIT2 was used to perform statistical phasing with information
628 from sequence reads [83].

629 **Population structure**

630 Ancestry informative marker (AIM), F_{ST} , doubleton sharing and SNP PCA were con-
631 ducted following methods defined in [12]. The PCA and UMAP analyses were performed
632 on 131,679 SNPs from euchromatic regions of chromosome arms 3L and 3R obtained
633 from the full dataset via random downsampling to 100,000 non-singleton SNPs from
634 each chromosome arm then performing LD-pruning. To generate the UMAP projec-
635 tion shown in Figure 2, each country and species was downsampled to a maximum of
636 50 individuals, to provide a projection that was less warped by differences in sample
637 size. The UMAP analysis was also performed on the full set of individuals, which gave
638 qualitatively identical results in terms of the clustering of individuals. UMAP was per-
639 formed using the umap-learn Python package [26] with the following parameter settings:
640 $n_neighbors = 15$; $min_dist = 2$; $spread = 5$; $metric = euclidean$. Other parameter
641 values for $n_neighbours$ and min_dist were also performed, all producing qualitatively
642 identical results. One population (Guinea *An. coluzzii*, n=4) was excluded from F_{ST}
643 analysis and three populations (Guinea *An. coluzzii*, n=4; Bioko *An. gambiae*, n=9;

644 Ghana *An. gambiae*, n=12) were excluded from doubleton sharing analysis due to small
645 sample size. All analyses of geographical population structure using SNP data were con-
646 ducted on euchromatic regions of Chromosome 3 (3R:1-37 Mbp, 3L:15-41 Mbp), which
647 avoids regions of polymorphic inversions, reduced recombination and unequal divergence
648 from the reference genome [12]. Unscaled CNV variation PCAs were built from the CNV
649 presence/absence calls [19], using the *prcomp* function in R [84].

650 Admixture models were fitted using the program LEA version 2.0 [85] in R version 3.6.1
651 [84]. Ten independent sets of SNPs were generated by selecting SNPs from euchromatic
652 regions of Chromosome 3 with minor allele frequency greater than 1%, then randomly
653 selecting 100,000 SNPs from each chromosome arm, then applying the same LD pruning
654 methodology as used for PCA, then combining back together remaining SNPs from both
655 chromosome arms. The resulting files were exported in .geno format, which were then
656 analyzed using the *snmf* method (sparse non-negative matrix factorization [28]) to obtain
657 ancestry estimates to each cluster (K) tested. We tested all K values from 2 to 15. Ten
658 replicates of the analysis with *snmf* were run for each dataset, which meant that 100 runs
659 were performed for each K. We assessed the convergence and replicability of the results
660 across the 100 runs (ten different datasets, each one replicated ten times dataset) using
661 CLUMPAK [86]. CLUMPAK was used to summarize the results, identify the major and
662 minor clustering solutions identified at each K (if they occurred), and estimate the average
663 ancestry proportions for the major solution which was used to interpret the results. We
664 assessed how the clustering solution fitted with the data using the cross-entropy criterion.
665 The lower this criterion is, the better is the model fit to the data.

666 Genetic diversity

667 Analyses of genetic diversity, including nucleotide diversity, Tajima's D, ROH and IBD
668 (identity by descent), were conducted following methods defined in [12] but using the phase
669 2 data release of 1,142 samples. In short, scikit-allel version 1.2.0 was used to calculate
670 windowed averages of nucleotide diversity and Tajima's D [87], IBDseq version r1206 [88]
671 was used to calculate IBD, and an HMM implemented in scikit-allel was used to calculate
672 ROH.

673 **The *Anopheles gambiae* 1000 Genomes Consortium**

674 Please address correspondence to Alistair Miles <alistair.miles@bdi.ox.ac.uk> and Do-
675 minic Kwiatkowski <dominic@sanger.ac.uk>.

676 Chris S. Clarkson and Alistair Miles jointly led curation of the phase 2 data resource
677 and wrote the paper.

678 **Data analysis group**

679 Chris S. Clarkson¹, Alistair Miles^{2,1}, Nicholas J. Harding², Eric R. Lucas³, C. J. Battey⁴,
680 Jorge Edouardo Amaya-Romero^{5,6}, Andrew D. Kern⁴, Michael C. Fontaine^{5,6}, Martin J.
681 Donnelly^{3,1}, Mara K. N. Lawniczak¹ and Dominic P. Kwiatkowski^{1,2} (chair).

682 **Partner working group**

683 Martin J. Donnelly^{3,1} (chair), Diego Ayala^{7,5}, Nora J. Besansky⁸, Austin Burt⁹, Beni-
684 amino Caputo¹⁰, Alessandra della Torre¹⁰, Michael C. Fontaine^{5,6}, H. Charles J. God-
685 fray¹¹, Matthew W. Hahn¹², Andrew D. Kern⁴, Dominic P. Kwiatkowski^{2,1}, Mara K. N.
686 Lawniczak¹, Janet Midega¹³, Samantha O'Loughlin⁹, João Pinto¹⁴, Michelle M. Riehle¹⁵,

¹Parasites and Microbes Programme, Wellcome Sanger Institute, Hinxton, Cambridge CB10 1SA, UK.

²MRC Centre for Genomics and Global Health, University of Oxford, Oxford OX3 7BN, UK.

³Department of Vector Biology, Liverpool School of Tropical Medicine, Pembroke Place, Liverpool L3 5QA, UK.

⁴Institute for Ecology and Evolution, University of Oregon, 301 Pacific Hall, Eugene, OR 97403, USA.

⁵Laboratoire MIVEGEC (Université de Montpellier, CNRS 5290, IRD 229), Centre IRD de Montpellier, 911, Avenue Agropolis BP 64501, 34395 Montpellier Cedex 5, France.

⁶Groningen Institute for Evolutionary Life Sciences (GELIFES), University of Groningen, PO Box 11103 CC, Groningen, The Netherlands.

⁷Unit d'Ecologie des Systèmes Vectoriels, Centre International de Recherches Médicales de Franceville, Franceville, Gabon.

⁸Eck Institute for Global Health, Department of Biological Sciences & University of Notre Dame, IN 46556, USA.

⁹Department of Life Sciences, Imperial College, Silwood Park, Ascot, Berkshire SL5 7PY, UK.

¹⁰Istituto Pasteur Italia - Fondazione Cenci Bolognetti, Dipartimento di Sanita Pubblica e Malattie Infettive, Università di Roma SAPIENZA, Rome, Italy.

¹¹Department of Zoology, University of Oxford, 11a Mansfield Road, Oxford OX1 3SZ, UK.

¹²Department of Biology and School of Informatics and Computing, Indiana University, Bloomington, IN 47405, USA.

¹³KEMRI-Wellcome Trust Research Programme, PO Box 230, Bofa Road, Kilifi, Kenya.

¹⁴Global Health and Tropical Medicine, GHTM, Instituto de Higiene e Medicina Tropical, IHMT, Universidade Nova de Lisboa, UNL, Rua da Junqueira 100, 1349-008 Lisbon, Portugal.

¹⁵Department of Microbiology and Immunology, Medical College of Wisconsin, Milwaukee, WI 53226, USA.

687 Igor Sharakhov^{16,17}, Daniel R. Schrider¹⁸, Kenneth D. Vernick¹⁹, David Weetman³, Craig
688 S. Wilding²⁰ and Bradley J. White²¹.

689 Population sampling

690 **Angola:** Arlete D. Troco²², João Pinto¹⁴; **Bioko:** Jorge Cano²³; **Burkina Faso:** Ab-
691 doulaye Diabaté²⁴, Samantha O'Loughlin⁹, Austin Burt⁹; **Cameroon:** Carlo Costan-
692 tini^{5,25}, Kyanne R. Rohatgi⁸, Nora J. Besansky⁸; **Côte d'Ivoire:** Edi Constant²⁶, David
693 Weetman³; **Gabon:** Nohal Elissa²⁷, João Pinto¹⁴; **Gambia:** Davis C. Nwakanma²⁸, Musa
694 Jawara²⁸; **Ghana:** John Essandoh²⁹, David Weetman³; **Guinea:** Boubacar Coulibaly³⁰,
695 Michelle M. Riehle¹⁵, Kenneth D. Vernick¹⁹; **Guinea-Bissau:** João Pinto¹⁴, João Di-
696 nis³¹; **Kenya:** Janet Midega¹³, Charles Mbogo¹³, Philip Bejon¹³; **Mayotte:** Gilbert Le
697 Goff⁵, Vincent Robert⁵; **Uganda:** Craig S. Wilding²⁰, David Weetman³, Henry D. Mawe-
698 jje³², Martin J. Donnelly³; **Lab crosses:** David Weetman³, Craig S. Wilding²⁰, Martin
699 J. Donnelly³.

700 Sequencing and data production

701 Jim Stalker³³, Kirk A. Rockett², Eleanor Drury¹, Daniel Mead¹, Anna E. Jeffreys²,
702 Christina Hubbard², Kate Rowlands², Alison T. Isaacs³, Dushyanth Jyothi³⁴, Cinzia

¹⁶Department of Entomology, Virginia Tech, Blacksburg, VA 24061, USA.

¹⁷Department of Cytology and Genetics, Tomsk State University, Tomsk 634050, Russia.

¹⁸Department of Genetics, University of North Carolina, 5111 Genetic Medicine Building, 7264, Chapel Hill, NC 27599-7264, USA.

¹⁹Unit for Genetics and Genomics of Insect Vectors, Institut Pasteur, Paris, France.

²⁰School of Biological and Environmental Sciences, Liverpool John Moores University, Liverpool L3 3AF, UK.

²¹Verily Life Sciences, 269 E Grand Ave, South San Francisco, CA 94080, USA.

²²Programa Nacional de Controle da Malária, Direção Nacional de Saúde Pública, Ministério da Saúde, Luanda, Angola.

²³London School of Hygiene & Tropical Medicine. Keppel St, Bloomsbury, London WC1E 7HT, UK.

²⁴Institut de Recherche en Sciences de la Santé (IRSS), Bobo Dioulasso, Burkina Faso.

²⁵Laboratoire de Recherche sur le Paludisme, Organisation de Coordination pour la lutte contre les Endémies en Afrique Centrale (OCEAC), Yaoundé, Cameroon.

²⁶Centre Suisse de Recherches Scientifiques. Yopougon, Abidjan - 01 BP 1303 Abidjan, Côte d'Ivoire.

²⁷Institut Pasteur de Madagascar, Avaradoha, BP 1274, 101, Antananarivo, Madagascar.

²⁸Medical Research Council Unit The Gambia at the London School of Hygiene & Tropical Medicine (MRCG at LSHTM), Atlantic Boulevard, Fajara, P.O. Box 273, Banjul, The Gambia.

²⁹Department of Wildlife and Entomology, University of Cape Coast, Cape Coast, Ghana.

³⁰Malaria Research and Training Centre, Faculty of Medicine and Dentistry, University of Mali.

³¹Instituto Nacional de Saúde Pública, Ministério da Saúde Pública, Bissau, Guiné-Bissau.

³²Infectious Diseases Research Collaboration, 2C Nakasero Hill Road, PO Box 7475, Kampala, Uganda.

³³Microbiotica Limited, Biodata, Innovation Centre, Wellcome Genome Campus, Cambridge, CB10 1DR, UK.

³⁴European Bioinformatics Institute, Hinxton, Cambridge CB10 1SA, UK.

703 Malangone³⁴ and Maryam Kamali^{35,16}.

704 **Project coordination**

705 Victoria Simpson², Christa Henrichs² and Dominic P. Kwiatkowski^{1,2}.

706 **Acknowledgments**

707 The authors would like to thank the staff of the Wellcome Sanger Institute Sample Logis-
708 tics, Sequencing and Informatics facilities for their contributions. The sequencing, anal-
709 ysis, informatics and management of the *Anopheles gambiae* 1000 Genomes Project are
710 supported by Wellcome through Sanger Institute core funding (098051), core funding
711 to the Wellcome Centre for Human Genetics (203141/Z/16/Z), and a strategic award
712 (090770/Z/09/Z); and by the MRC Centre for Genomics and Global Health which is
713 jointly funded by the Medical Research Council and the Department for International
714 Development (DFID) (G0600718; M006212). M.K.N.L. was supported by MRC grant
715 G1100339. S.O.L. and A.B. were supported by a grant from the Foundation for the Na-
716 tional Institutes of Health through the Vector-Based Control of Transmission: Discovery
717 Research (VCTR) program of the Grand Challenges in Global Health initiative of the Bill
718 and Melinda Gates Foundation. D.W., C.S.W., H.D.M. and M.J.D. were supported by
719 Award Numbers U19AI089674 and R01AI082734 from the National Institute of Allergy
720 and Infectious Diseases (NIAID). The content is solely the responsibility of the authors
721 and does not necessarily represent the official views of the NIAID or NIH.

722 **Data availability**

723 Sequence read alignments and variant calls from Ag1000G phase 2 are available from the
724 European Nucleotide Archive under study accession PRJEB36277 (ENA - <http://www.ebi.ac.uk/ena>).
725 Sequence read alignments for samples in Ag1000G phase 1 are available under study ac-
726 cession PRJEB18691.

³⁵Department of Medical Entomology and Parasitology, Faculty of Medical Sciences, Tarbiat Modares University, Tehran, Iran.

727 All variation data from Ag1000G phase 2 can also be downloaded from the Ag1000G
728 public FTP site via the MalariaGEN website (<https://www.malariagen.net/resource/27>).

729 **Supplementary figures and tables**

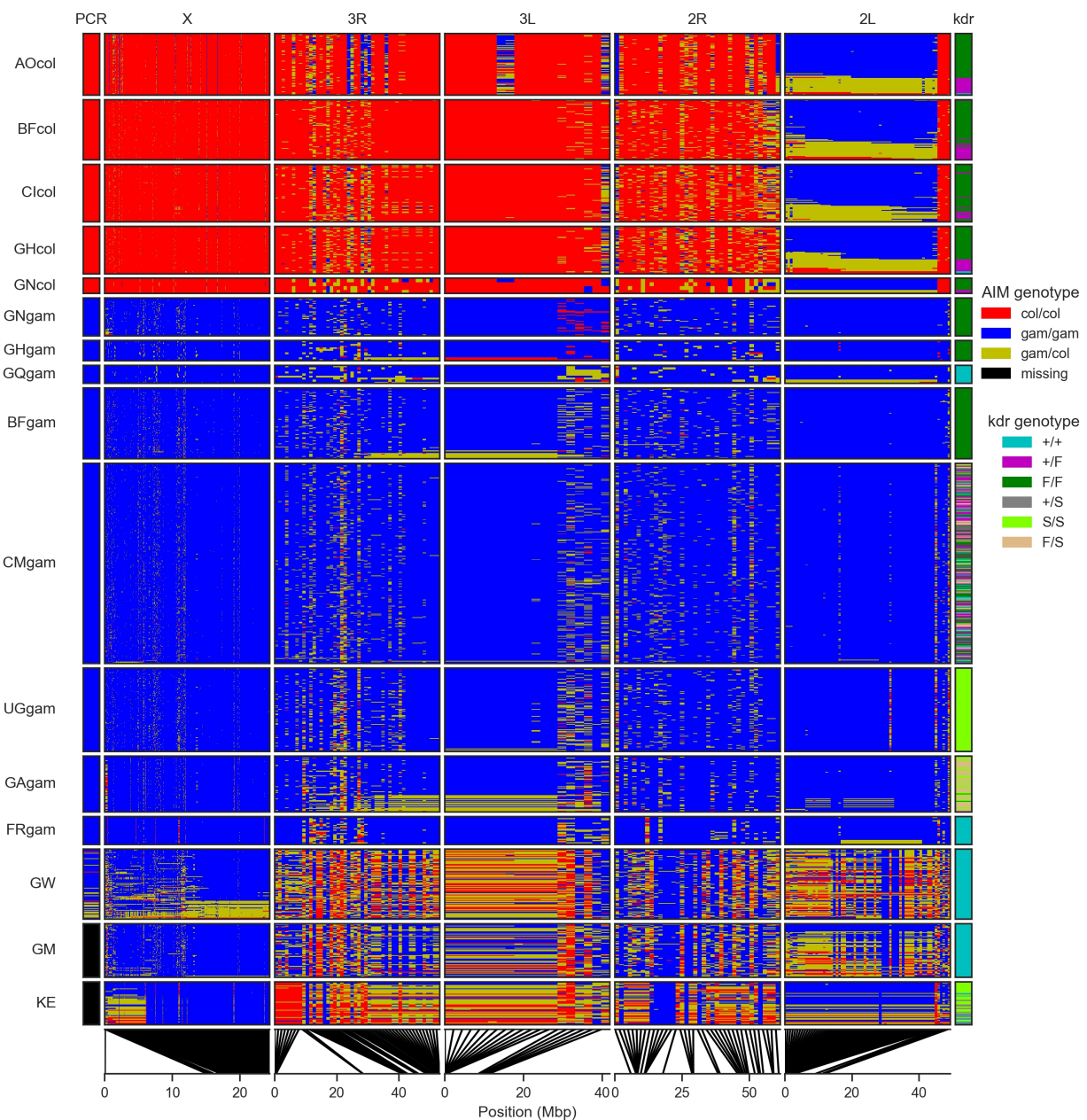


Figure S1. Ancestry informative markers (AIM). Rows represent individual mosquitoes (grouped by population) and columns represent SNPs (grouped by chromosome arm). Colours represent species genotype. The column at the far left (“PCR”) shows the species assignment according to the conventional molecular test based on a single marker on the X chromosome, which was performed for all populations except The Gambia (GM) and Kenya (KE). The column at the far right shows the genotype for *kdr* variants in *Vgsc* codon 995. Lines at the lower edge show the physical locations of the AIM SNPs.

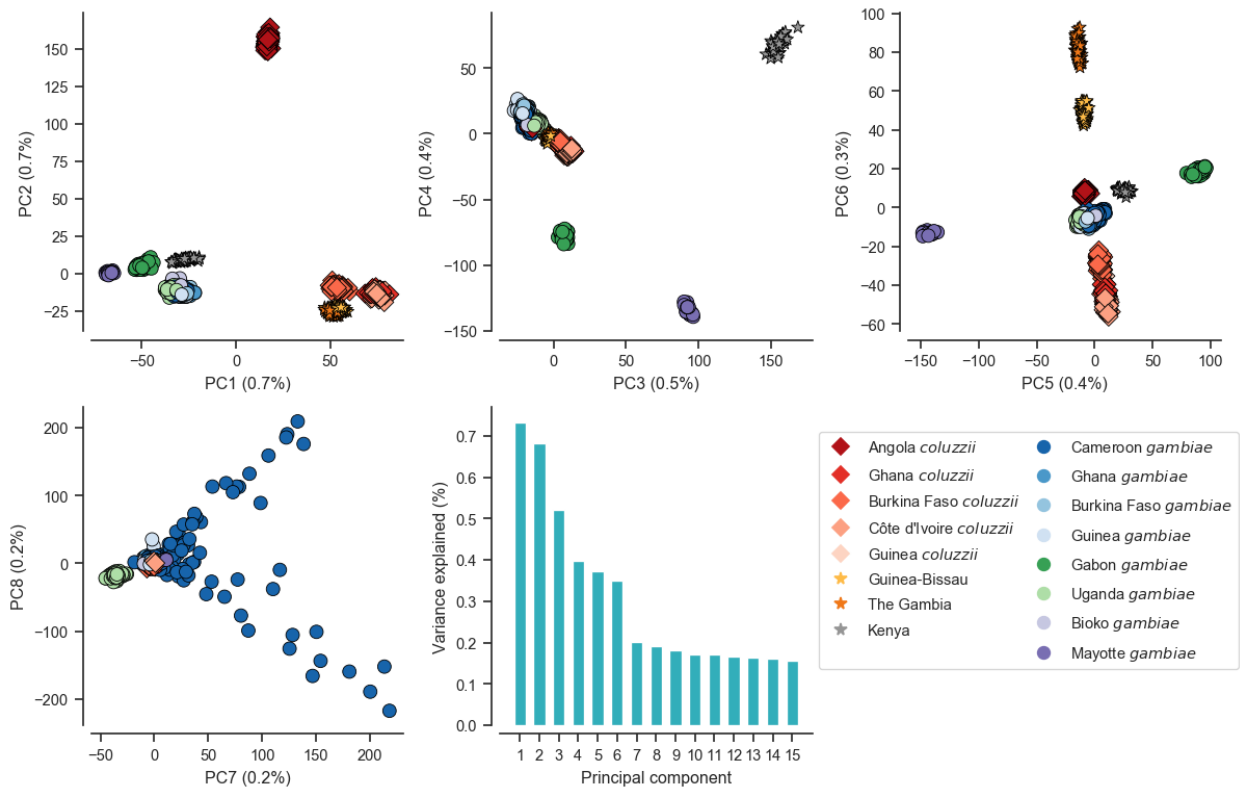


Figure S2. Principal component analysis of the 1,142 wild-caught mosquitoes using biallelic SNPs from euchromatic regions of Chromosome 3. Scatter plots show relationships of principle components 1-8 where each marker represents an individual mosquito. Marker shape and colour denotes population. The bar chart shows the percentage of variance explained by each principal component.

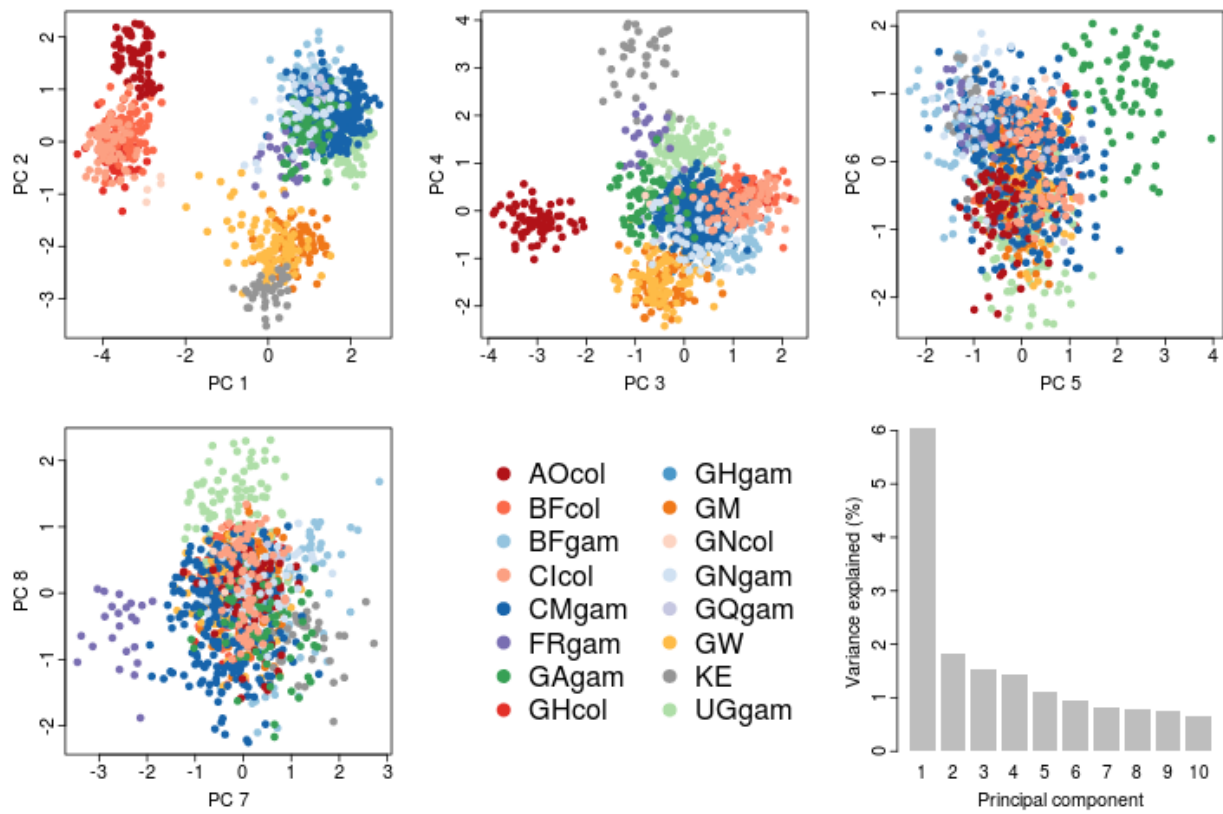


Figure S3. Principal component analysis of the 1,142 wild-caught mosquitoes using copy number variant calls. Bar chart shows the percentage of variance explained by each component.

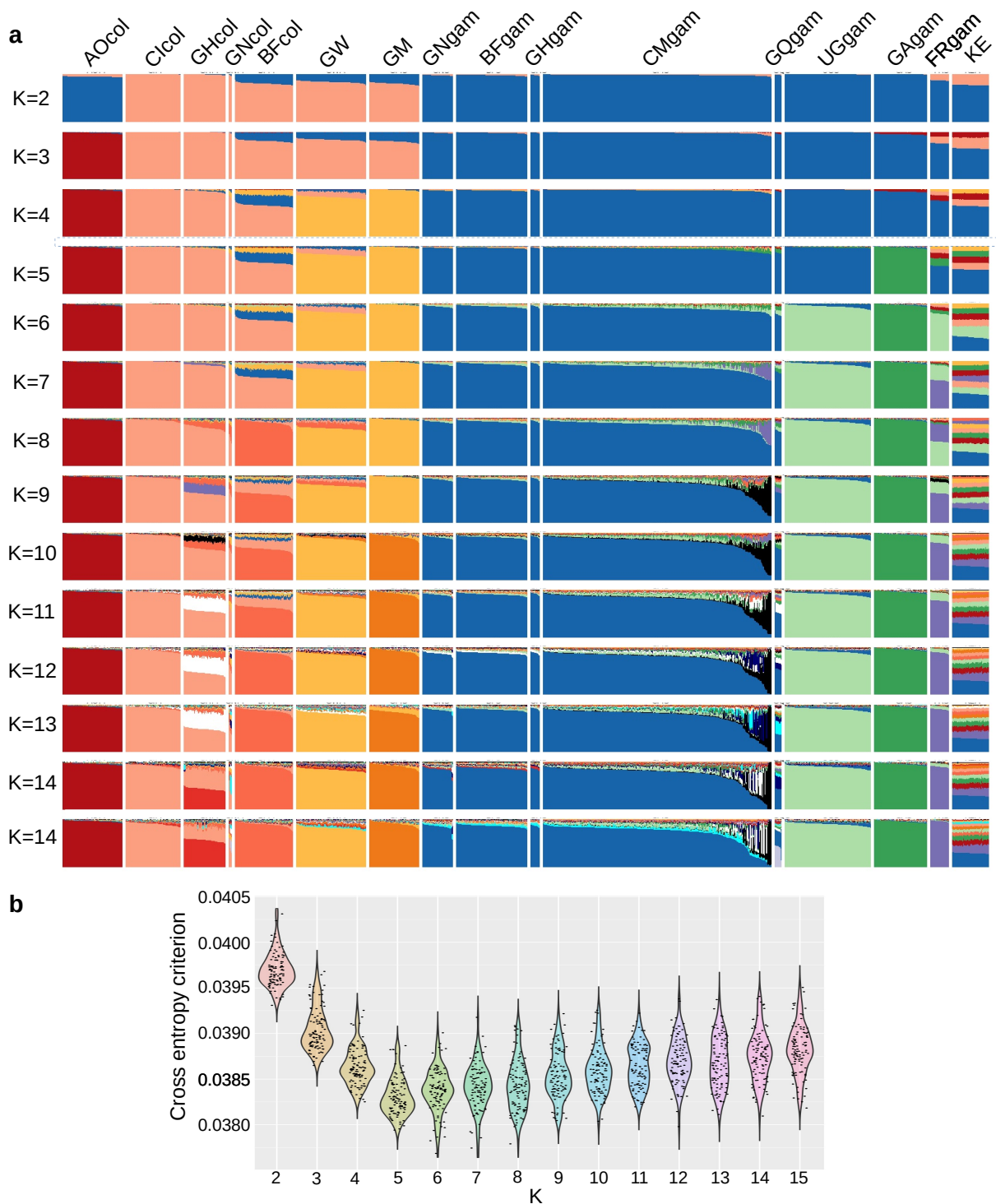


Figure S4. Analysis of population structure and admixture. **(a)** Each row shows results of modelling ancestry in sampled individuals assuming a given number K of ancestral populations [85]. Within each row, individual mosquitoes are represented as vertical bars, grouped according to sampling location and species, and coloured according to the proportion of the genome inherited from each ancestral population. **(b)** Cross-entropy criterion values obtained for each value of K ancestral populations, where lower values imply a better fit of the model to the data.

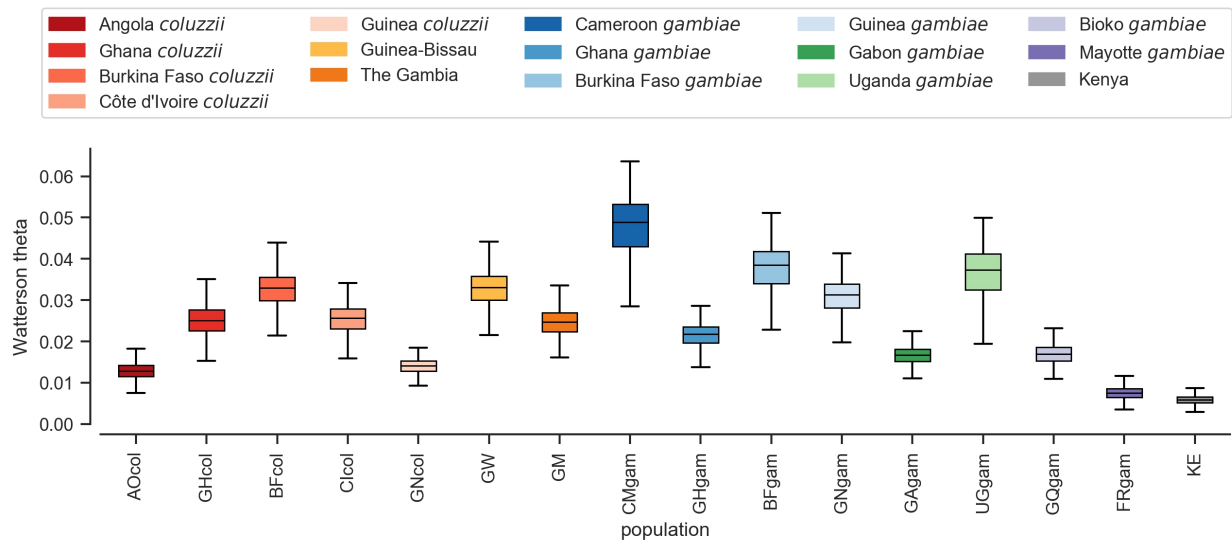


Figure S5. Watterson's theta (θ_W), the density of segregating sites, calculated in non-overlapping 20 kbp genomic windows using SNPs from euchromatic regions of Chromosome 3.

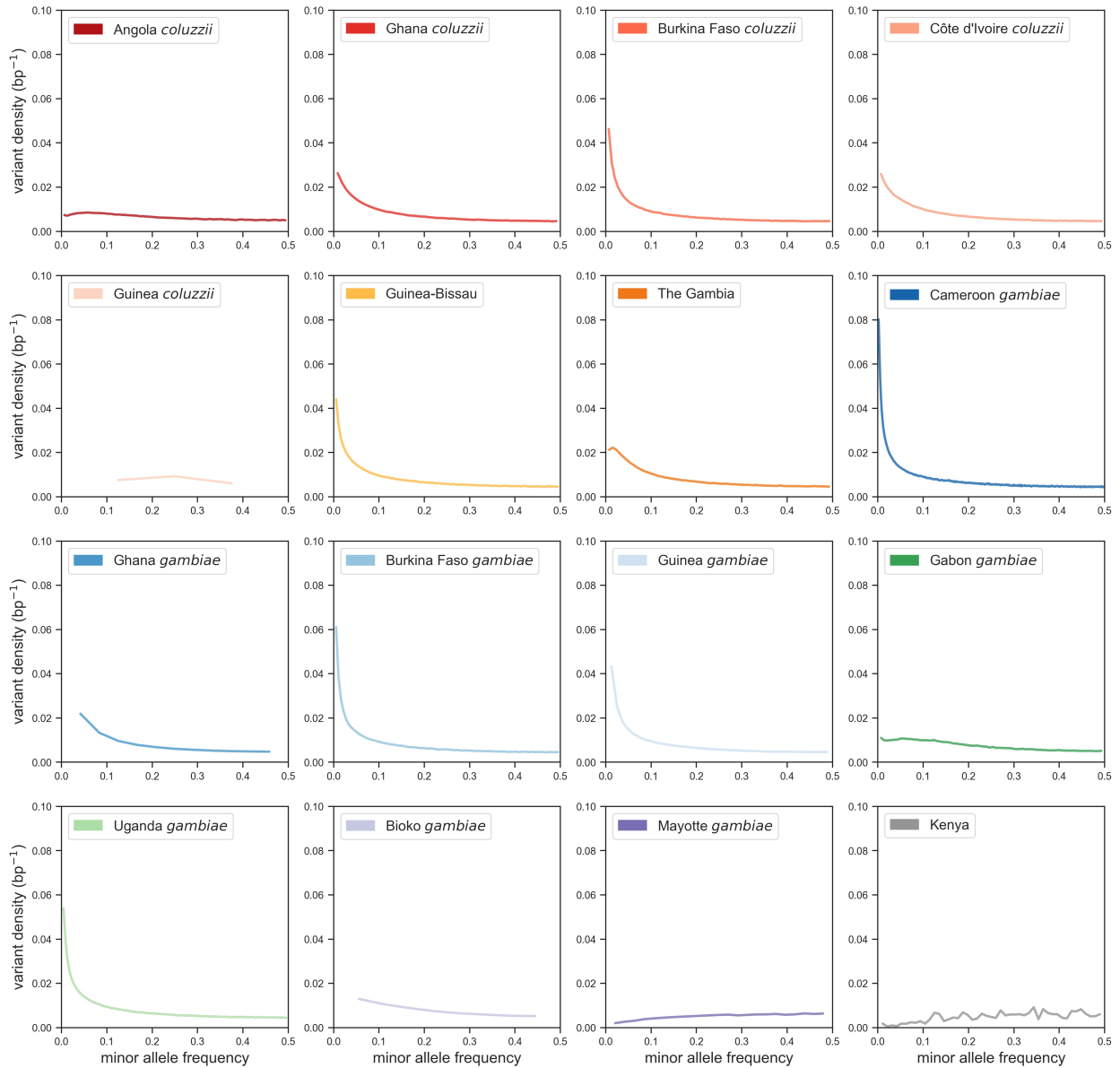


Figure S6. SNP density. Plots depict the distribution of allele frequencies (site frequency spectrum) for each population, scaled such that a population with constant size over time is expected to have a constant SNP density over all allele frequencies.

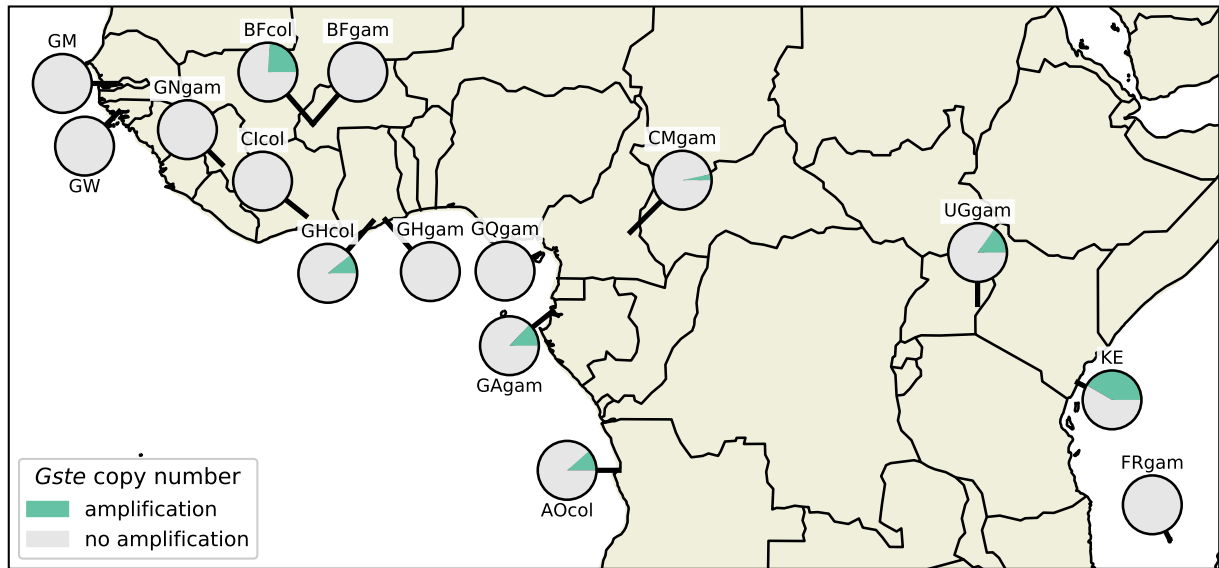


Figure S7. Prevalence of copy number amplifications at the *Gste* locus. Each pie shows the frequency of individuals from a given population carrying an amplification spanning at least one gene in the *Gste* gene cluster. The Guinea *An. coluzzii* population is omitted due to small sample size.

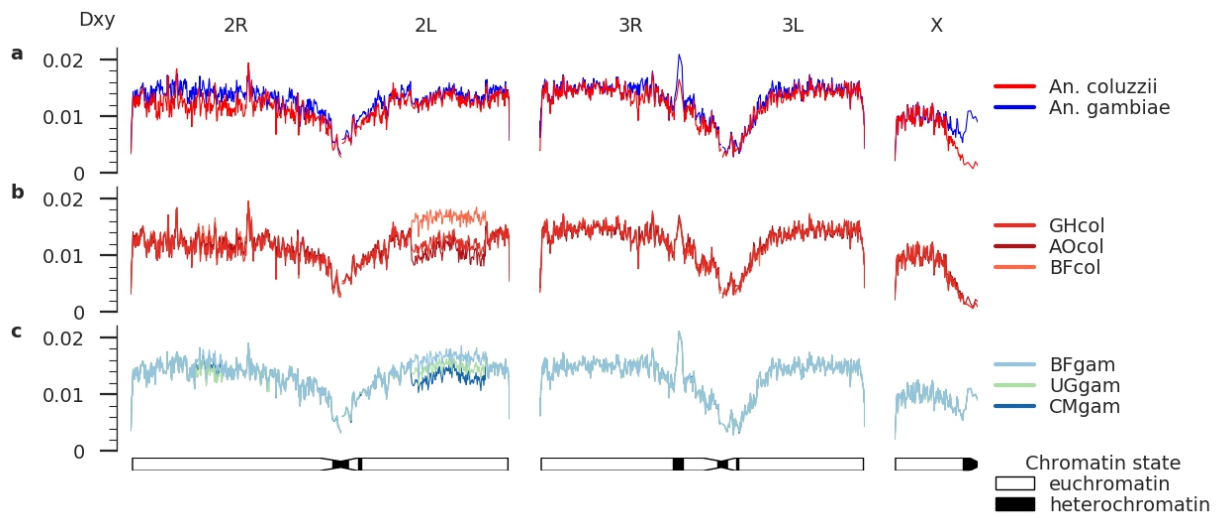


Figure S8. Divergence from the AgamP3 reference genome, calculated as D_{xy} , is largely similar for *An. coluzzii* and *An. gambiae*, with the exception of the centromere of the X chromosome (a). Comparing three populations of *An. coluzzii* (b) or *An. gambiae* (c) highlights the strong effect of the 2La chromosomal inversion on the accumulation of genetic variation.

Table S1. Ag1000G phase 2 sampling locations.

Country	Collection			Year	Latitude	Longitude	Sample size		
	Location	Site					Total	Female	Male
Angola	Luanda			2009	-8.821	13.291	78	78	0
Burkina Faso	Bana			2012	11.233	-4.472	60	40	20
	Pala			2012	11.150	-4.235	56	48	8
	Souroukoudinga			2012	11.235	-4.535	51	51	0
Cameroon	Daignuene			2009	4.777	13.844	96	81	15
	Gado Badzere			2009	5.747	14.442	73	58	15
	Mayos			2009	4.341	13.558	105	91	14
	Zembe Borongo			2009	5.747	14.442	23	23	0
Cote d'Ivoire	Tiassale			2012	5.898	-4.823	71	71	0
Equatorial Guinea	Bioko			2002	3.700	8.700	9	9	0
France	Mayotte	Bouyouni		2011	-12.738	45.142	1	1	0
		Combani		2011	-12.779	45.143	5	2	3
		Karihani Lake		2011	-12.797	45.122	3	3	0
		Mont Benara		2011	-12.857	45.155	2	1	1
		Mtsamboro Forest Reserve		2011	-12.703	45.081	1	1	0
		Mtsanga Charifou		2011	-12.991	45.156	8	3	5
		Sada		2011	-12.852	45.104	4	1	3
Gabon	Libreville			2000	0.384	9.455	69	69	0
Gambia, The	Njabakunda	Kerr Birom Kardo		2011	13.550	-15.900	19	19	0
		Kerr Sama Kuma		2011	13.550	-15.900	8	8	0
		Maria Samba Nyado		2011	13.550	-15.900	18	18	0
		Sare Illo Buya		2011	13.550	-15.900	20	20	0
Ghana	Koforidua			2012	6.094	-0.261	1	1	0
				2012	5.668	-0.219	24	24	0
				2012	4.912	-1.774	20	20	0
				2012	5.609	-1.549	22	22	0
Guinea	Koraboh			2012	9.250	-9.917	22	22	0
				2012	8.500	-9.417	22	22	0
Guinea-Bissau	Antula			2010	11.891	-15.582	58	58	0
				2010	11.957	-15.649	33	33	0
Kenya	Kilifi	Junju		2012	-3.862	39.745	16	16	0
		Mbogolo		2012	-3.635	39.858	32	32	0
Uganda	Tororo	Nagongera		2012	0.770	34.026	112	112	0

Table S2. Colony crosses.

Cross ID	Mother Colony	Father Colony	N progeny
18-5	Ghana	Kisumu/G3	20
29-2	Ghana	Kisumu	20
36-9	Ghana	Mali	20
37-3	Kisumu	Pimperena	20
42-4	Mali	Kisumu/Ghana	14
45-1	Mali	Kisumu	20
46-9	Pimperena	Mali	20
47-6	Mali	Kisumu	20
73-2	Akron	Ghana	19
78-2	Mali	Kisumu/Ghana	19
80-2	Kisumu	Akron	20

730 References

- 731 [1] *World malaria report 2019*. Tech. rep. World Health Organization, 2019.
- 732 [2] Janet Hemingway et al. ‘Averting a malaria disaster: Will insecticide resistance derail
733 malaria control?’ In: *The Lancet* 387.10029 (2016), pp. 1785–1788. ISSN: 1474547X.
- 734 [3] *Global report on insecticide resistance in malaria vectors: 2010–2016*. Tech. rep.
735 World Health Organization, 2018.
- 736 [4] *Global Technical Strategy for Malaria 2016–2030*. Tech. rep. World Health Organi-
737 zation, 2015.
- 738 [5] Deus S. Ishengoma et al. ‘Deployment and utilization of next-generation sequenc-
739 ing of *Plasmodium falciparum* to guide anti-malarial drug policy decisions in sub-
740 Saharan Africa: opportunities and challenges’. In: *Malaria Journal* 18 (2019).
- 741 [6] Richard M. Oxborough et al. ‘Susceptibility testing of *Anopheles* malaria vectors
742 with the neonicotinoid insecticide clothianidin; results from 16 African countries,
743 in preparation for indoor residual spraying with new insecticide formulations’. In:
744 *Malaria Journal* (2019).
- 745 [7] Rosemary Lees et al. ‘A testing cascade to identify repurposed insecticides for next-
746 generation vector control tools: screening a panel of chemistries with novel modes of
747 action against a malaria vector’. In: *Gates Open Research* (2019).
- 748 [8] Kyros Kyrou et al. ‘A CRISPR–Cas9 gene drive targeting doublesex causes com-
749 plete population suppression in caged *Anopheles gambiae* mosquitoes’. In: *Nature*
750 *Biotechnology* 36.11 (2018), p. 1062.
- 751 [9] R A Holt et al. ‘The genome sequence of the malaria mosquito *Anopheles gambiae*’.
752 In: *Science* 298.5591 (2002), pp. 129–149. ISSN: 0036-8075.
- 753 [10] Maria V Sharakhova et al. ‘Update of the *Anopheles gambiae* PEST genome assem-
754 bly’. In: *Genome Biology* 8.1 (2007), R5.
- 755 [11] Gloria I Giraldo-Calderón et al. ‘VectorBase: an updated bioinformatics resource for
756 invertebrate vectors and other organisms related with human diseases’. In: *Nucleic*
757 *Acids Research* 43.D1 (2014), pp. D707–D713.
- 758 [12] *Anopheles gambiae* 1000 Genomes Consortium et al. ‘Genetic diversity of the African
759 malaria vector *Anopheles gambiae*’. In: *Nature* 552.7683 (2017), p. 96.

- 760 [13] Maureen Coetzee et al. ‘*Anopheles coluzzii* and *Anopheles amharicus*, new members
761 of the *Anopheles gambiae* complex’. In: *Zootaxa* 3619.3 (2013), pp. 246–274.
- 762 [14] Antoinette Wiebe et al. ‘Geographical distributions of African malaria vector sibling
763 species and evidence for insecticide resistance’. In: *Malaria Journal* 16.1 (2017),
764 p. 85.
- 765 [15] Robert T Schimke et al. ‘Gene amplification and drug resistance in cultured murine
766 cells’. In: *Science* 202.4372 (1978), pp. 1051–1055.
- 767 [16] Alan L Devonshire and Linda M Field. ‘Gene amplification and insecticide resis-
768 tance’. In: *Annual Review of Entomology* 36.1 (1991), pp. 1–21.
- 769 [17] David Weetman et al. ‘Contemporary evolution of resistance at the major insecti-
770 cide target site gene *Ace-1* by mutation and copy number variation in the malaria
771 mosquito *Anopheles gambiae*’. In: *Molecular Ecology* 24.11 (2015), pp. 2656–2672.
- 772 [18] R. G. et al. Sayre. ‘A New Map of Standardized Terrestrial Ecosystems of Africa’.
773 In: *American Association of Geographers* (2013).
- 774 [19] Eric R Lucas et al. ‘Whole-genome sequencing reveals high complexity of copy num-
775 ber variation at insecticide resistance loci in malaria mosquitoes’. In: *Genome Re-
776 search* 29.8 (2019), pp. 1250–1261.
- 777 [20] C Fanello, F Santolamazza and A Della Torre. ‘Simultaneous identification of species
778 and molecular forms of the *Anopheles gambiae* complex by PCR-RFLP’. In: *Medical
779 and Veterinary Entomology* 16.4 (2002), pp. 461–464.
- 780 [21] Federica Santolamazza et al. ‘Insertion polymorphisms of SINE200 retrotransposons
781 within speciation islands of *Anopheles gambiae* molecular forms’. In: *Malaria Journal*
782 7.1 (2008), p. 163.
- 783 [22] Mylène Weill et al. ‘The *kdr* mutation occurs in the Mopti form of *Anopheles gambiae*
784 s.s. through introgression’. In: *Insect Molecular Biology* 9.5 (2000), pp. 451–455.
- 785 [23] Abdoulaye Diabaté et al. ‘The spread of the Leu-Phe *kdr* mutation through *Anophe-*
786 *les gambiae* complex in Burkina Faso: genetic introgression and de novo phenomena’.
787 In: *Tropical Medicine & International Health* 9.12 (2004), pp. 1267–1273.

- 788 [24] Chris S Clarkson et al. ‘Adaptive introgression between *Anopheles* sibling species
789 eliminates a major genomic island but not reproductive isolation’. In: *Nature Com-*
790 *munications* 5 (2014), p. 4248.
- 791 [25] Laura C Norris et al. ‘Adaptive introgression in an African malaria mosquito coinci-
792 dent with the increased usage of insecticide-treated bed nets’. In: *Proceedings of the*
793 *National Academy of Sciences* 112.3 (2015), pp. 815–820.
- 794 [26] Leland McInnes, John Healy and James Melville. *UMAP: Uniform Manifold Ap-*
795 *proximation and Projection for Dimension Reduction*. 2018. arXiv: 1802.03426
796 [stat.ML].
- 797 [27] ‘Population structure and eigenanalysis’. In: *PLoS Genetics* 2.12 (2006), pp. 2074–
798 2093.
- 799 [28] Eric Frichot et al. ‘Fast and efficient estimation of individual ancestry coefficients’.
800 In: *Genetics* (2014). ISSN: 19432631.
- 801 [29] Daniel J. Lawson, Lucy van Dorp and Daniel Falush. ‘A tutorial on how not to over-
802 interpret STRUCTURE and ADMIXTURE bar plots’. In: *Nature Communications*
803 (2018). ISSN: 20411723.
- 804 [30] Ace R North, Austin Burt and H Charles J Godfray. ‘Modelling the potential of
805 genetic control of malaria mosquitoes at national scale’. In: *BMC Biology* 17.1 (2019),
806 p. 26.
- 807 [31] A Dao et al. ‘Signatures of aestivation and migration in Sahelian malaria mosquito
808 populations’. In: *Nature* 516.7531 (2014), p. 387.
- 809 [32] Diana L Huestis et al. ‘Windborne long-distance migration of malaria mosquitoes in
810 the Sahel’. In: *Nature* 574.7778 (2019), pp. 404–408.
- 811 [33] Sewall Wright. ‘Isolation by distance under diverse systems of mating’. In: *Genetics*
812 31.1 (1946), p. 39.
- 813 [34] François Rousset. ‘Genetic differentiation and estimation of gene flow from F-statistics
814 under isolation by distance’. In: *Genetics* 145.4 (1997), pp. 1219–1228.

- 815 [35] Austin Burt. ‘Site-specific selfish genes as tools for the control and genetic engineer-
816 ing of natural populations’. In: *Proceedings of the Royal Society of London. Series*
817 *B: Biological Sciences* 270.1518 (2003), pp. 921–928.
- 818 [36] Valentino M Gantz et al. ‘Highly efficient Cas9-mediated gene drive for population
819 modification of the malaria vector mosquito *Anopheles stephensi*’. In: *Proceedings of*
820 *the National Academy of Sciences* 112.49 (2015), E6736–E6743.
- 821 [37] Andrew Hammond et al. ‘A CRISPR-Cas9 gene drive system targeting female repro-
822 duction in the malaria mosquito vector *Anopheles gambiae*’. In: *Nature Biotechnology*
823 34.1 (2016), p. 78.
- 824 [38] Philip A Eckhoff et al. ‘Impact of mosquito gene drive on malaria elimination in a
825 computational model with explicit spatial and temporal dynamics’. In: *Proceedings*
826 *of the National Academy of Sciences* 114.2 (2017), E255–E264.
- 827 [39] Robert L Unckless, Andrew G Clark and Philipp W Messer. ‘Evolution of resistance
828 against CRISPR/Cas9 gene drive’. In: *Genetics* 205.2 (2017), pp. 827–841.
- 829 [40] Christina Scali et al. ‘Identification of sex-specific transcripts of the *Anopheles gam-*
830 *biae* doublesex gene’. In: *Journal of Experimental Biology* (2005). ISSN: 00220949.
- 831 [41] Tanja Gempe and Martin Beye. *Function and evolution of sex determination mech-*
832 *anisms, genes and pathways in insects*. 2011.
- 833 [42] Elzbieta Krzywinska et al. ‘A maleness gene in the malaria mosquito *Anopheles*
834 *gambiae*’. In: *Science* (2016). ISSN: 10959203.
- 835 [43] Thomas W. Cline and and Barbara J. Meyer. ‘VIVE LA DIFFÉRENCE: Males vs
836 Females in Flies vs Worms’. In: *Annual Review of Genetics* (1996). ISSN: 0066-4197.
- 837 [44] Hilary Ranson and Natalie Lissenden. ‘Insecticide resistance in African *Anopheles*
838 mosquitoes: a worsening situation that needs urgent action to maintain malaria
839 control’. In: *Trends in Parasitology* 32.3 (2016), pp. 187–196.
- 840 [45] Thomas S Churcher et al. ‘The impact of pyrethroid resistance on the efficacy and
841 effectiveness of bednets for malaria control in Africa’. In: *Elife* 5 (2016), e16090.

- 842 [46] S. Bhatt et al. ‘The effect of malaria control on *Plasmodium falciparum* in Africa
843 between 2000 and 2015’. In: *Nature* 526.7572 (2015), pp. 207–211. ISSN: 0028-0836.
844 arXiv: [arXiv:1011.1669v3](https://arxiv.org/abs/1011.1669v3).
- 845 [47] Christopher M Jones et al. ‘Footprints of positive selection associated with a mu-
846 tation (N1575Y) in the voltage-gated sodium channel of *Anopheles gambiae*’. In:
847 *Proceedings of the National Academy of Sciences* 109.17 (2012), pp. 6614–6619.
- 848 [48] Sara N Mitchell et al. ‘Metabolic and target-site mechanisms combine to confer
849 strong DDT resistance in *Anopheles gambiae*’. In: *PLoS One* 9.3 (2014), e92662.
- 850 [49] David Weetman et al. ‘Candidate-gene based GWAS identifies reproducible DNA
851 markers for metabolic pyrethroid resistance from standing genetic variation in East
852 African *Anopheles gambiae*’. In: *Scientific Reports* 8.1 (2018), p. 2920.
- 853 [50] R. M. Kwiatkowska et al. ‘Dissecting the mechanisms responsible for the multiple
854 insecticide resistance phenotype in *Anopheles gambiae* s.s., M form, from Vallée du
855 Kou, Burkina Faso’. In: *Gene* 519.1 (2013), pp. 98–106.
- 856 [51] Constant V Edi et al. ‘CYP6 P450 enzymes and ACE-1 duplication produce extreme
857 and multiple insecticide resistance in the malaria mosquito *Anopheles gambiae*’. In:
858 *PLoS Genetics* 10.3 (2014), e1004236.
- 859 [52] C. Ngufor et al. ‘Insecticide resistance profile of *Anopheles gambiae* from a phase II
860 field station in Cové, southern Benin: implications for the evaluation of novel vector
861 control products’. In: *Malaria Journal* 14.1 (2015), p. 464.
- 862 [53] John Vontas et al. ‘Rapid selection of a pyrethroid metabolic enzyme CYP9K1 by
863 operational malaria control activities’. In: *Proceedings of the National Academy of
864 Sciences* 115.18 (2018), pp. 4619–4624.
- 865 [54] Pie Müller et al. ‘Field-caught permethrin-resistant *Anopheles gambiae* overexpress
866 CYP6P3, a P450 that metabolises pyrethroids’. In: *PLoS Genetics* 4.11 (2008),
867 e1000286.
- 868 [55] Adriana Adolphi et al. ‘Functional genetic validation of key genes conferring insecticide
869 resistance in the major African malaria vector, *Anopheles gambiae*’. In: *Proceedings*

- 870 *of the National Academy of Sciences* (2019). ISSN: 0027-8424. eprint: <https://www.pnas.org/content/early/2019/12/03/1914633116.full.pdf>.
- 871
- 872 [56] Dimitra Nikou, Hilary Ranson and Janet Hemingway. ‘An adult-specific CYP6 P450
873 gene is overexpressed in a pyrethroid-resistant strain of the malaria vector, *Anopheles*
874 *gambiae*’. In: *Gene* 318 (2003), pp. 91–102.
- 875 [57] Bradley J Stevenson et al. ‘Cytochrome P450 6M2 from the malaria vector *Anopheles*
876 *gambiae* metabolizes pyrethroids: sequential metabolism of deltamethrin revealed’.
877 In: *Insect Biochemistry and Molecular Biology* 41.7 (2011), pp. 492–502.
- 878 [58] Beniamino Caputo et al. ‘The “far-west” of *Anopheles gambiae* molecular forms’. In:
879 *PloS One* 6.2 (2011), e16415.
- 880 [59] Kevin Ochieng’Opondo et al. ‘Does insecticide resistance contribute to heterogeneities
881 in malaria transmission in The Gambia?’ In: *Malaria Journal* 15.1 (2016), p. 166.
- 882 [60] Jacob M Riveron et al. ‘A single mutation in the GSTe2 gene allows tracking of
883 metabolically based insecticide resistance in a major malaria vector’. In: *Genome*
884 *Biology* 15.2 (2014), R27.
- 885 [61] Nena Pavlidi, John Vontas and Thomas Van Leeuwen. *The role of glutathione S-*
886 *transferases (GSTs) in insecticide resistance in crop pests and disease vectors*. 2018.
- 887 [62] *Conditions for deployment of mosquito nets treated with a pyrethroid and piperonyl*
888 *butoxide*. Tech. rep. World Health Organization, 2017.
- 889 [63] Eric R Lucas et al. ‘A high throughput multi-locus insecticide resistance marker
890 panel for tracking resistance emergence and spread in *Anopheles gambiae*’. In: *Sci-*
891 *entific Reports* 9.1 (2019), pp. 1–10.
- 892 [64] Luigi Sedda et al. ‘Improved spatial ecological sampling using open data and stan-
893 dardization: an example from malaria mosquito surveillance’. In: *Journal of the Royal*
894 *Society Interface* 16.153 (2019), p. 20180941.
- 895 [65] Chris S. Clarkson et al. ‘The genetic architecture of target-site resistance to pyrethroid
896 insecticides in the African malaria vectors *Anopheles gambiae* and *Anopheles coluzzii*’.
897 In: *BioRxiv* (2018). eprint: [https://www.biorxiv.org/content/early/2018/08/](https://www.biorxiv.org/content/early/2018/08/06/323980.full.pdf)
898 [06/323980.full.pdf](https://www.biorxiv.org/content/early/2018/08/06/323980.full.pdf).

- 899 [66] Ace R. North and H. Charles J. Godfray. ‘Modelling the persistence of mosquito
900 vectors of malaria in Burkina Faso’. In: *Malaria Journal* (2018). ISSN: 14752875.
- 901 [67] Christina M. Bergey et al. ‘Assessing connectivity despite high diversity in island
902 populations of a malaria mosquito’. In: *BioRxiv* (2019). eprint: <https://www.biorxiv.org/content/early/2019/02/28/430702.full.pdf>.
- 903
- 904 [68] Hussein Al-Asadi et al. ‘Estimating recent migration and population-size surfaces’.
905 In: *PLoS Genetics* (2019). ISSN: 15537404.
- 906 [69] Julie A Scott, William G Brogdon and Frank H Collins. ‘Identification of single
907 specimens of the *Anopheles gambiae* complex by the polymerase chain reaction’. In:
908 *The American Journal of Tropical Medicine and Hygiene* 49.4 (1993), pp. 520–529.
- 909 [70] Federica Santolamazza, Alessandra della Torre and Adalgisa Caccone. ‘A new poly-
910 merase chain reaction-restriction fragment length polymorphism method to identify
911 *Anopheles arabiensis* from *An. gambiae* and its two molecular forms from degraded
912 DNA templates or museum samples’. In: *The American Journal of Tropical Medicine
913 and Hygiene* 70.6 (2004), pp. 604–606.
- 914 [71] Bradley J White et al. ‘Molecular karyotyping of the 2La inversion in *Anopheles
915 gambiae*’. In: *The American Journal of Tropical Medicine and Hygiene* 76.2 (2007),
916 pp. 334–339.
- 917 [72] Brian L Sharp et al. ‘Malaria vector control by indoor residual insecticide spraying
918 on the tropical island of Bioko, Equatorial Guinea’. In: *Malaria Journal* 6.1 (2007),
919 p. 52.
- 920 [73] Hans J Overgaard et al. ‘Malaria transmission after five years of vector control on
921 Bioko Island, Equatorial Guinea’. In: *Parasites & Vectors* 5.1 (2012), p. 253.
- 922 [74] Heinrich Magnus Manske and Dominic P Kwiatkowski. ‘LookSeq: a browser-based
923 viewer for deep sequencing data’. In: *Genome Research* 19.11 (2009), pp. 2125–2132.
- 924 [75] Beniamino Caputo et al. ‘*Anopheles gambiae* complex along The Gambia river, with
925 particular reference to the molecular forms of *An. gambiae* ss’. In: *Malaria Journal*
926 7.1 (2008), p. 182.

- 927 [76] Davis C Nwakanma et al. ‘Breakdown in the process of incipient speciation in
928 *Anopheles gambiae*’. In: *Genetics* (2013), pp. 1221–1231.
- 929 [77] José L Vicente et al. ‘Massive introgression drives species radiation at the range
930 limit of *Anopheles gambiae*’. In: *Scientific Reports* 7 (2017), p. 46451.
- 931 [78] Vasco Gordicho et al. ‘First report of an exophilic *Anopheles arabiensis* population
932 in Bissau City, Guinea-Bissau: recent introduction or sampling bias?’ In: *Malaria*
933 *Journal* 13.1 (2014), p. 423.
- 934 [79] MJ Donnelly et al. ‘Population structure in the malaria vector, *Anopheles arabiensis*
935 Patton, in East Africa’. In: *Heredity* 83.4 (1999), p. 408.
- 936 [80] Heng Li and Richard Durbin. ‘Fast and accurate short read alignment with Burrows–
937 Wheeler transform’. In: *Bioinformatics* 25.14 (2009), pp. 1754–1760.
- 938 [81] Geraldine A Van der Auwera et al. ‘From FastQ data to high-confidence variant
939 calls: the genome analysis toolkit best practices pipeline’. In: *Current Protocols in*
940 *Bioinformatics* 43.1 (2013), pp. 11–10.
- 941 [82] Pablo Cingolani et al. ‘A program for annotating and predicting the effects of single
942 nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster*
943 strain w1118; iso-2; iso-3’. In: *Fly* 6.2 (2012), pp. 80–92. ISSN: 19336942.
- 944 [83] Olivier Delaneau et al. ‘Haplotype estimation using sequencing reads’. In: *American*
945 *Journal of Human Genetics* 93.4 (2013), pp. 687–696. ISSN: 00029297.
- 946 [84] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foun-
947 dation for Statistical Computing. Vienna, Austria, R.3.4.4 2019.
- 948 [85] Eric Frichot and Olivier François. ‘LEA: An R package for landscape and ecological
949 association studies’. In: *Methods in Ecology and Evolution* (2015). ISSN: 2041210X.
- 950 [86] Naama M. Kopelman et al. ‘Clumpak: A program for identifying clustering modes
951 and packaging population structure inferences across K’. In: *Molecular Ecology Re-*
952 *sources* (2015). ISSN: 17550998.
- 953 [87] A Miles and N Harding. *scikit-allel-Explore and analyse genetic variation. In., 1.*
954 2018.

- 955 [88] Sharon R Browning and Brian L Browning. ‘Accurate non-parametric estimation
956 of recent effective population size from segments of identity by descent’. In: *The*
957 *American Journal of Human Genetics* 97.3 (2015), pp. 404–418.