

# Genetic diversity of the African malaria vector *Anopheles gambiae*

The *Anopheles gambiae* 1000 Genomes Consortium\*

**The sustainability of malaria control in Africa is threatened by the rise of insecticide resistance in *Anopheles* mosquitoes, which transmit the disease<sup>1</sup>. To gain a deeper understanding of how mosquito populations are evolving, here we sequenced the genomes of 765 specimens of *Anopheles gambiae* and *Anopheles coluzzii* sampled from 15 locations across Africa, and identified over 50 million single nucleotide polymorphisms within the accessible genome. These data revealed complex population structure and patterns of gene flow, with evidence of ancient expansions, recent bottlenecks, and local variation in effective population size. Strong signals of recent selection were observed in insecticide-resistance genes, with several sweeps spreading over large geographical distances and between species. The design of new tools for mosquito control using gene-drive systems will need to take account of high levels of genetic diversity in natural mosquito populations.**

Blood-sucking mosquitoes of the *An. gambiae* species complex are the principal vectors of *Plasmodium falciparum* malaria in Africa. Substantial reductions in malaria morbidity and mortality have been achieved by the use of insecticide-based interventions<sup>2</sup>, but increasing levels of insecticide resistance and other adaptive changes in mosquito populations threaten to reverse these gains<sup>1</sup>. A better understanding of the molecular, ecological and evolutionary processes driving these changes is essential to maximize the active lifespan of existing insecticides, and to accelerate the development of new strategies and tools for vector control. The *Anopheles gambiae* 1000 Genomes Project (Ag1000G; <http://www.malariagen.net/ag1000g>) was established to provide a foundation for detailed investigation of mosquito genome variation and evolution. Here we report the first phase of the project, which analysed 765 wild-caught specimens of *An. gambiae* sensu stricto and *An. coluzzii*. These two species account for the majority of malaria transmission in Africa, and are morphologically indistinguishable and often sympatric, but are genetically distinct<sup>3,4</sup> and differ in geographical range<sup>5</sup>, larval ecology<sup>6</sup>, behaviour<sup>7</sup> and strategies for surviving the dry season<sup>8</sup>. The specimens were collected at 15 locations across 8 African countries, spanning a range of ecologies including rainforest, inland savannah and coastal biomes, and thus provide a broad sample in which to explore factors shaping mosquito population variation (Extended Data Fig. 1; Supplementary Information 1).

Specimens were sequenced using the Illumina HiSeq platform, and single nucleotide polymorphisms (SNPs) were identified by alignment against the AgamP3 reference genome (Methods; Supplementary Information 2). A rigorous evaluation of data quality, including the use of experimental genetic crosses to quantify error rates, identified genomic regions totalling 141 megabases (Mb; 61% of the reference genome) that were accessible for the analysis of population variation (Supplementary Information 3; Extended Data Fig. 2). We identified 52,525,957 high-quality SNPs, of which 21% had three or more alleles, an average of one variant allele every 2.2 bases of the accessible genome (Fig. 1a). Individual mosquitoes carried between 1.7 and 2.7 million variant alleles, with no systematic difference observed between the two species (Extended Data Fig. 3a). In most populations, nucleotide

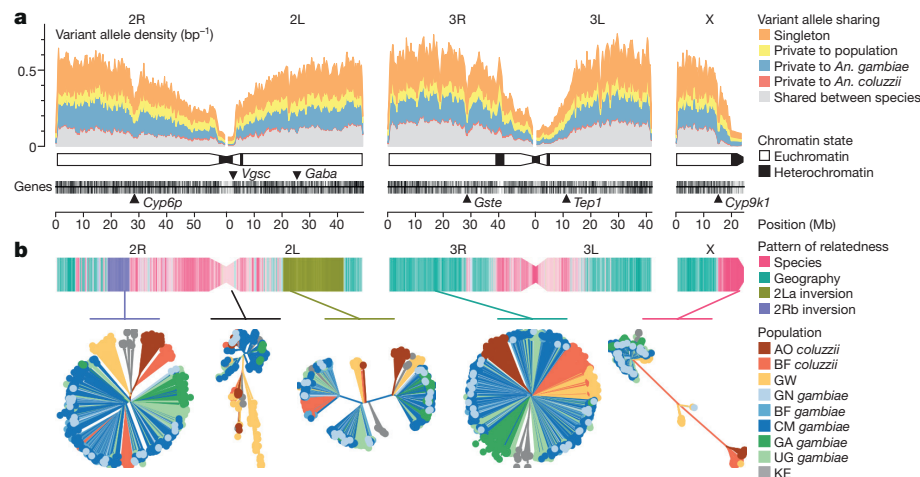
diversity was 1.5% on average (Extended Data Fig. 3b) and more than 3% at synonymous coding sites (Extended Data Fig. 3c), confirming that these are among the most genetically diverse eukaryotic species<sup>9</sup>.

High levels of natural diversity have practical implications for the development of gene-drive technologies for mosquito control<sup>10</sup>. CRISPR–Cas9 gene drives can be designed to edit a specific gene and confer a phenotype such as female sterility, which could suppress mosquito populations and thereby reduce disease transmission. However, naturally occurring polymorphisms within the approximately 21-base-pair (bp) Cas9 target site could prevent target recognition, and thus undermine gene-drive efficacy in the field. We found viable Cas9 targets in 11,625 protein-coding genes, but only 5,474 genes remained after excluding target sites with nucleotide variation in any of the 765 genomes sequenced here (Extended Data Fig. 3d; Supplementary Information 5). Resistance to gene drive could be countered by designing constructs that target multiple sites within the same gene, and we identified 863 genes that each contain at least 10 non-overlapping conserved target sites, including 13 putative sterility genes<sup>10</sup> (Supplementary Information 5.2). However, clearly more variants remain to be discovered (Extended Data Fig. 3d), and extensive sampling of multiple populations will be needed to inform the design of gene drives that are robust to natural genetic variation.

*An. gambiae* and *An. coluzzii* have a geographical range that spans sub-Saharan Africa and encompasses a variety of ecological settings<sup>5</sup>. Previous studies have found evidence that populations are locally adapted, and that migration between populations is limited by both geographical distance and major ecological discontinuities, notably the Congo Basin tropical rainforest and the East African rift system<sup>11–14</sup>. As a starting point for the analysis of population structure, we constructed neighbour-joining trees to explore patterns of genetic similarity between individuals (Fig. 1b; Supplementary Information 6.1). We observed four contrasting patterns of relatedness, associated with different regions of the genome. Within pericentromeric regions of chromosomes X, 3 and arm 2R, mosquitoes segregated into two highly distinct clades, largely corresponding to the two species as determined by conventional molecular diagnostics, consistent with previous studies that showed that genome regions of reduced recombination are associated with stronger differentiation between closely related species<sup>15</sup>. The large chromosomal inversions 2La and 2Rb were each associated with a distinct pattern of relatedness, as expected if recombination is reduced between inversion karyotypes. In most of the remaining genome, there was evidence of clustering by geographical region but not by species. There were also some genome regions in which we found unusually short genetic distances between individuals from different countries and/or species, indicating the influence of recent selective sweeps and adaptive gene flow.

To investigate geographical sub-divisions in more detail, we focused on euchromatic regions of chromosome 3, which are free from polymorphic inversions and regions of reduced recombination (Supplementary Information 6). ADMIXTURE models and principal component analysis (PCA) supported five major ancestral populations,

\*Lists of participants and their affiliations appear at the end of the paper.



**Figure 1 | Patterns of genomic variation.** **a**, Density of nucleotide variation in 200-kb windows over the genome. **b**, Variation in the pattern of relatedness between individual mosquitoes over the genome. The three chromosomes are painted using colours to represent the major pattern of relatedness found within each 100-kb window. Bottom,

neighbour-joining trees are shown from a selection of genomic windows that are representative of the four major patterns of relatedness found, as well as for the window spanning the *Vgsc* gene. AO, Angola; BF, Burkina Faso; CM, Cameroon; GA, Gabon; GN, Guinea; GW, Guinea-Bissau; KE, Kenya; UG, Uganda.

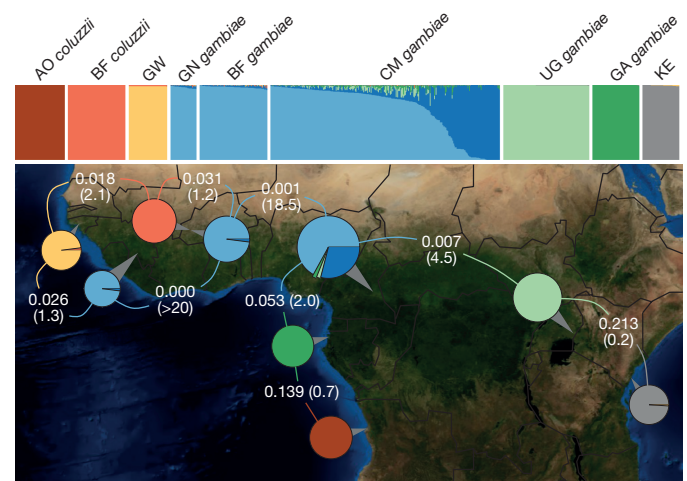
represented by: (i) *An. gambiae* from Guinea, Burkina Faso, Cameroon and Uganda; (ii) *An. gambiae* from Gabon; (iii) all individuals from Kenya; (iv) *An. coluzzii* from Angola; and (v) *An. coluzzii* from Burkina Faso and all individuals from Guinea-Bissau (Fig. 2; Extended Data Figs 4, 5). Within each species, we found relatively high allele frequency differentiation across the Congo Basin rainforest, exceeding differentiation between the two species at a single location (Extended Data Fig. 5b). There were also more subtle distinctions within and between populations. For example, in Cameroon, mosquitoes were sampled along a cline from savannah into forest, and there was some population structure associated with these different ecologies. However, among *An. gambiae* populations north of the Congo Basin, differentiation was extremely weak overall, despite considerable distances between populations, suggesting substantial gene flow.

Earlier studies concluded that purposeful movement of *Anopheles* mosquitoes is limited to short-range dispersal of up to 5 km<sup>16</sup>; however, recent evidence has emerged for long-distance seasonal migration in *An. gambiae*<sup>8</sup>. To explore evidence for migration, we computed joint site frequency spectra for selected population pairs and fitted models of population history (Methods; Supplementary Information 8). For all pairs examined, models with migration provided a better fit than models without migration (Supplementary Table 2). The inferred rate of migration was high between *An. gambiae* savannah populations, but some migration was also inferred between species and across both the Congo Basin rainforest and the East African rift. Although these analyses do not allow us to infer the timing or direction of gene flow events, they suggest that mosquito migration between different parts of the continent could affect the spread of insecticide resistance and dynamics of disease transmission.

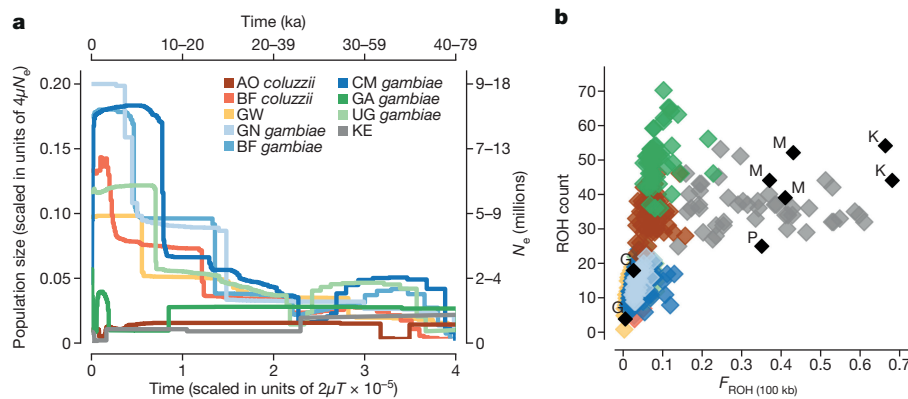
A key question in mosquito evolution concerns the extent and effect of gene flow between species, and *An. gambiae* and *An. coluzzii* are known to undergo hybridization at a rate that varies over space and time<sup>17</sup>. To study this phenomenon, we analysed 506 SNPs previously found to be highly differentiated between the two species<sup>18</sup> (Extended Data Fig. 6; Supplementary Information 6.6). These ancestry-informative markers (AIMs) showed that a genomic region on chromosome arm 2L has introgressed from *An. gambiae* into *An. coluzzii* in Burkina Faso and Angola. This region spans the *Vgsc* gene, in which introgression of insecticide-resistance alleles has been reported in Ghana<sup>19</sup> and Mali<sup>20</sup>, although this is the first evidence, to our knowledge, that introgressed alleles have spread to *An. coluzzii* south of the Congo Basin. AIMs also highlighted two populations with uncertain species status. In Guinea-Bissau, mosquitoes carried a

mixture of alleles from both species on all chromosomes. These individuals were sampled from the coast, within a region of West Africa that is thought to be a zone of secondary contact because previous studies have found evidence for extensive introgression<sup>21,22</sup>. We also found that mosquitoes from coastal Kenya carried a mixture of both species' alleles on all chromosomes. This was unexpected, as the geographical range of *An. coluzzii* is not thought to extend beyond the East African rift. There are several possible explanations for the Kenyan data, including historical admixture between species and retention of ancestral variation, and further analysis and population sampling are required. However, our data demonstrate that a simple *gambiae/coluzzii* dichotomy is not adequate for describing malaria vector species composition in some parts of Africa, and caution against the use of any single marker to infer species ancestry or recent hybridization.

Historical fluctuations in effective population size ( $N_e$ ) can be inferred from the genomes of extant individuals. Analysis of our genome variation data indicated a major expansion in all populations



**Figure 2 | Geographical population structure and migration.** Top, each mosquito is depicted as a vertical bar painted by the proportion of the genome inherited from each of  $K = 8$  inferred ancestral populations. Pie charts on the map depict the same ancestry proportions summed over all individuals for each population. Text in white shows average fixation index ( $F_{ST}$ ) followed in parentheses by estimates of the population migration rate ( $2Nm$ ).



**Figure 3 | Population size history.** **a**, Stairway plot of changes in population size over time. Absolute values of time and  $N_e$  are shown on alternative axes as a range of values, assuming lower and upper limits for the mutation rate  $\mu$  as  $2.8 \times 10^{-9}$  and  $5.5 \times 10^{-9}$ , respectively, and  $t = 11$

north of the Congo Basin and west of the East African rift (Fig. 3a; Extended Data Fig. 7; Methods; Supplementary Information 8). Knowledge of the *Anopheles* mutation rate is required to date this expansion, and this has not yet been determined, but assuming it is similar to *Drosophila* then the onset of expansion would be within the range 7,000 to 25,000 years ago (Fig. 3a; Methods). Because *An. gambiae* and *An. coluzzii* are highly anthropophilic, mosquito population expansion could be linked to that of humans, and particularly to the expansion of agricultural Bantu-speaking groups that originate from north of the Congo Basin beginning approximately 5,000 years ago<sup>23</sup>. It is possible to reconcile this theory with our data if *Anopheles* has a higher mutation rate than *Drosophila*, causing us to overestimate the age of the expansion, but it is also possible that mosquito populations benefited from earlier human population growth, or that other factors such as climate change were involved.

We also observed genomic signatures of a major recent population decline of *An. gambiae* in coastal Kenya. All Kenyan specimens (but no specimens from other locations) had long runs of homozygosity that comprised 10–60% of the genome, indicating high levels of inbreeding consistent with a recent population bottleneck (Fig. 3b). In Kenya, free mass distribution of insecticide-treated nets (ITNs) starting in 2006 resulted in a major increase in ITN coverage<sup>24</sup>. The specimens in this study were collected in 2012, raising the question of whether the population decline of *An. gambiae* can be attributed to ITN usage. To address this question, we analysed sharing of genome regions that are identical by descent (IBD) (Methods; Extended Data Fig. 8a, b). We estimated that the *An. gambiae* population in Kenya has fallen in size by at least two orders of magnitude, to  $N_e < 1,000$  (Extended Data Fig. 8c; Supplementary Information 8.4). The beginning of this inferred decline occurred approximately 200 generations before the date of sampling, which would pre-date mass distributions of ITNs, assuming approximately 11 generations per year. This is consistent with other studies that have found evidence for low  $N_e$  values<sup>11</sup> and changes in mosquito species abundance<sup>25</sup> in the region before high levels of ITN coverage. Nevertheless, our data show that major demographic events leave genetic signatures that could be used to gain important information about the impact of vector control interventions.

Many genes have been associated with insecticide resistance in *Anopheles*, but different genetic variants may be responsible for resistance in different populations, and it is not yet clear where or how resistance is spreading. Genomic data can help to address these questions by identifying genes with evidence of recent evolutionary adaptation in one or more mosquito populations. We found strong signals of recent positive selection at several genes that are known to have a role in resistance, including: *Vgsc*, the target site for dichlorodiphenyl-trichloroethane (DDT) and pyrethroid insecticides<sup>26</sup>; *Gste*, a cluster of glutathione S-transferase genes including *Gste2*, previously implicated

generations per year. **b**, Runs of homozygosity (ROH) in individual mosquitoes, highlighting recent inbreeding in Kenyan (grey) and colony (black) mosquitoes. G, Ghana; K, Kisumu; M, Mali; P, Pimperena.

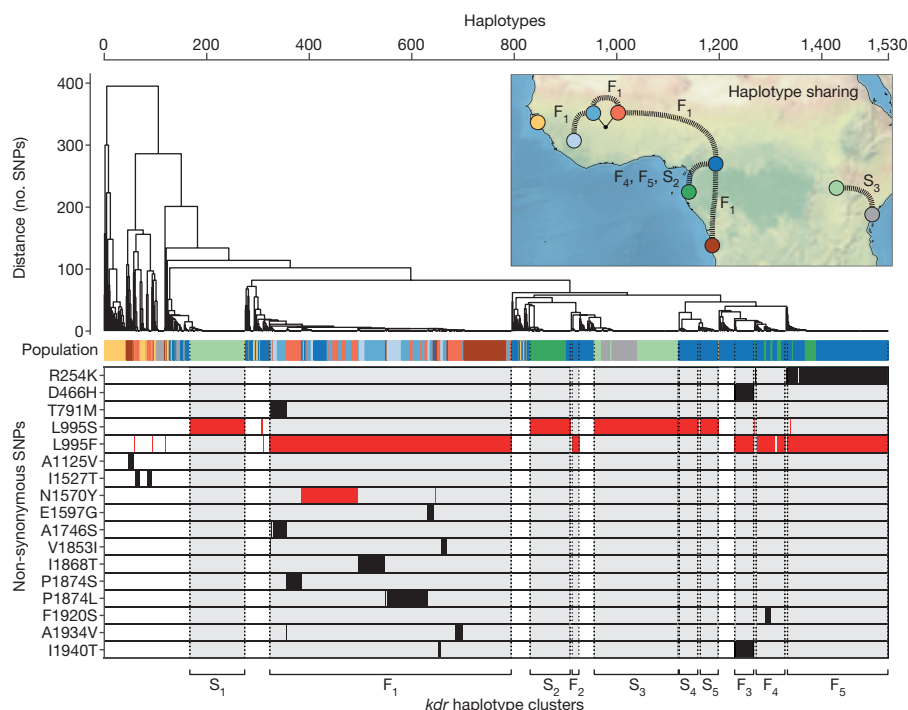
in the metabolism of DDT and pyrethroids<sup>27</sup>; and *Cyp6p*, a cluster of genes encoding cytochrome P450 enzymes, including *Cyp6p3*, which is upregulated in permethrin- and bendiocarb-resistant mosquitoes<sup>28</sup> (Extended Data Fig. 9; Supplementary Information 9). We also observed strong signals of selection at multiple loci with no known resistance genes, and these merit detailed investigation in future studies.

Mutations in *An. gambiae* *Vgsc* codon 995 (orthologous to *Musca domestica* *Vgsc* codon 1014), known as '*kdr*' owing to their knock-down-resistance phenotype, reduce susceptibility to DDT and pyrethroids<sup>26</sup>. We found the leucine-to-phenylalanine (Leu995Phe) *kdr* variant at high frequency in West and Central Africa (Guinea 100%; Burkina Faso 93%; Cameroon 53%; Gabon 36%; Angola 86%). A second *kdr* allele, a leucine-to-serine (Leu995Ser) variant, was present in Central and East Africa (Cameroon 15%; Gabon 65%; Uganda 100%; Kenya 76%). To investigate the evolution and spread of the two *kdr* alleles, we analysed the genetic backgrounds on which they were carried (Fig. 4; Supplementary Information 9.3). Leu995Phe occurred within five distinct haplotype clusters (labelled F<sub>1</sub>–F<sub>5</sub> in Fig. 4), whereas Leu995Ser was found in a further 5 haplotype clusters (labelled S<sub>1</sub>–S<sub>5</sub> in Fig. 4). Cluster F<sub>1</sub> contained individuals of both species and from four countries spanning the Congo Basin, proving that recent gene flow has carried resistance alleles between these populations. Three *kdr* haplotypes (F<sub>4</sub>, F<sub>5</sub> and S<sub>2</sub>) were found in both Cameroon and Gabon, providing multiple examples of recent gene flow between these two populations. The S<sub>3</sub> haplotype was present in both Uganda and coastal Kenya, thus resistance alleles can reach populations on both sides of the rift system.

Although the evolution of resistance in the *Vgsc* gene is clearly driven primarily by the two *kdr* alleles, we also found 15 other non-synonymous variants at a frequency of above 1% in our cohort (Fig. 4). Thirteen of these variants occurred almost exclusively on haplotypes carrying the Leu995Phe allele (Lewontin's  $D' > 0.96$ ). These included Asn1570Tyr, previously found on Leu995Phe haplotypes in West and Central Africa and shown to confer increased resistance<sup>29</sup>. Overall, there was a statistically significant enrichment for non-synonymous mutations on haplotypes carrying the Leu995Phe allele, indicating secondary selection on multiple variants that either enhance or compensate for the Leu995Phe phenotype (Supplementary Information 9.5).

Resistance due to genes that enhance insecticide metabolism is also a serious concern, as it has been implicated in extreme resistance phenotypes in some *Anopheles* populations<sup>27,28</sup>. Although several metabolic genes have been shown to be upregulated in resistant mosquitoes, only a single molecular marker of metabolic resistance (*Gste2*-Ile114Thr) has previously been identified in *An. gambiae* or *An. coluzzii*<sup>27</sup>. At both *Gste* and *Cyp6p*, we found evidence that resistance has emerged on several genetic backgrounds and is spreading between species and over considerable distances. At the *Gste* locus, we found at least four distinct





**Figure 4 | Evolution and spread of insecticide resistance in the *Vgsc* gene.** Top, a dendrogram obtained by hierarchical clustering of haplotypes from wild-caught individuals. The colour bar below shows the population of origin for each haplotype. Bottom, alleles carried by each haplotype at 17 non-synonymous SNPs with alternative allele frequency > 1%

(white denotes reference allele; black denotes alternative allele; red denotes previously known resistance allele). At the lower margin, we label 10 haplotype clusters carrying a *kdr* allele (either Leu995Phe or Leu995Ser). The inset map depicts haplotypes that are shared between populations, demonstrating the spread of insecticide resistance.

haplotypes under selection (Extended Data Fig. 10a). One of these haplotypes carried the known *Gste2*-Ile114Thr resistance allele, and this haplotype was found in all populations except Guinea-Bissau and Uganda, indicating a continent-wide spread. However, the other three haplotypes did not carry this allele, thus other genetic variants with a resistance phenotype must be present at this locus. At the *Cyp6p* locus, we found at least eight distinct haplotypes under selection, but limited spread between populations (Extended Data Fig. 10b). At both loci, we found multiple SNPs associated with haplotypes under selection that could be used as markers to track the spread of resistance and characterize resistance phenotypes (Extended Data Fig. 10).

In 1899, Ronald Ross proposed that malaria could be controlled by destroying breeding sites of the mosquitoes that transmit the disease<sup>30</sup>. *An. gambiae*, identified in the same year by Ross as a vector of malaria in Africa, has proved resilient to a century of attempts to repress it. The vector control armamentarium needs to be expanded, not only with new classes of insecticide and genetic control strategies, but also with tools to gather intelligence, to enable those responsible for planning and executing interventions to stay ahead of the mosquito's remarkable capacity for rapid evolutionary adaptation. There remain major knowledge gaps concerning the ecology and life history of *Anopheles* mosquitoes, such as the rate and range of migration, which are fundamental to understand both malaria transmission and the spread of insecticide resistance, and which will require spatiotemporal analysis of mosquito populations. Most importantly, it is essential to start collecting population genomic data prospectively as an integral part of vector control interventions, to identify which strategies are causing increased insecticide resistance, or what it takes to cause a population crash of the magnitude observed in our Kenyan data. By treating each intervention as an experiment, and by analysing its effect on both mosquito and parasite populations, we can aim to improve the efficacy and sustainability of future interventions, while at the same time learning about basic processes in ecology and evolution.

**Online Content** Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

**Received 22 December 2016; accepted 1 November 2017.**

**Published online 29 November 2017.**

- Hemingway, J. *et al.* Averting a malaria disaster: will insecticide resistance derail malaria control? *Lancet* **387**, 1785–1788 (2016).
- Bhatt, S. *et al.* The effect of malaria control on *Plasmodium falciparum* in Africa between 2000 and 2015. *Nature* **526**, 207–211 (2015).
- della Torre, A. *et al.* Molecular evidence of incipient speciation within *Anopheles gambiae* s.s. in West Africa. *Insect Mol. Biol.* **10**, 9–18 (2001).
- Lawniczak, M. K. N. *et al.* Widespread divergence between incipient *Anopheles gambiae* species revealed by whole genome sequences. *Science* **330**, 512–514 (2010).
- Tene Fossog, B. *et al.* Habitat segregation and ecological character displacement in cryptic African malaria mosquitoes. *Evol. Appl.* **8**, 326–346 (2015).
- Diabaté, A. *et al.* Larval development of the molecular forms of *Anopheles gambiae* (Diptera: Culicidae) in different habitats: a transplantation experiment. *J. Med. Entomol.* **42**, 548–553 (2005).
- Gimonneau, G. *et al.* A behavioral mechanism underlying ecological divergence in the malaria mosquito *Anopheles gambiae*. *Behav. Ecol.* **21**, 1087–1092 (2010).
- Dao, A. *et al.* Signatures of aestivation and migration in Sahelian malaria mosquito populations. *Nature* **516**, 387–390 (2014).
- Leffler, E. M. *et al.* Revisiting an old riddle: what determines genetic diversity levels within species? *PLoS Biol.* **10**, e1001388 (2012).
- Hammond, A. *et al.* A CRISPR-Cas9 gene drive system targeting female reproduction in the malaria mosquito vector *Anopheles gambiae*. *Nat. Biotechnol.* **34**, 78–83 (2016).
- Lehmann, T. *et al.* The Rift Valley complex as a barrier to gene flow for *Anopheles gambiae* in Kenya. *J. Hered.* **90**, 613–621 (1999).
- Lehmann, T. *et al.* Population structure of *Anopheles gambiae* in Africa. *J. Hered.* **94**, 133–147 (2003).
- Slotman, M. A. *et al.* Evidence for subdivision within the M molecular form of *Anopheles gambiae*. *Mol. Ecol.* **16**, 639–649 (2007).
- Pinto, J. *et al.* Geographic population structure of the African malaria vector *Anopheles gambiae* suggests a role for the forest-savannah biome transition as a barrier to gene flow. *Evol. Appl.* **6**, 910–924 (2013).
- Cruikshank, T. E. & Hahn, M. W. Reanalysis suggests that genomic islands of speciation are due to reduced diversity, not reduced gene flow. *Mol. Ecol.* **23**, 3133–3157 (2014).



16. Service, M. W. Mosquito (Diptera: Culicidae) dispersal—the long and short of it. *J. Med. Entomol.* **34**, 579–588 (1997).
17. Lee, Y. *et al.* Spatiotemporal dynamics of gene flow and hybrid fitness between the M and S forms of the malaria mosquito, *Anopheles gambiae*. *Proc. Natl Acad. Sci. USA* **110**, 19854–19859 (2013).
18. Neafsey, D. E. *et al.* SNP genotyping defines complex gene-flow boundaries among African malaria vector mosquitoes. *Science* **330**, 514–517 (2010).
19. Clarkson, C. S. *et al.* Adaptive introgression between *Anopheles* sibling species eliminates a major genomic island but not reproductive isolation. *Nat. Commun.* **5**, 4248 (2014).
20. Norris, L. C. *et al.* Adaptive introgression in an African malaria mosquito coincident with the increased usage of insecticide-treated bed nets. *Proc. Natl Acad. Sci. USA* **112**, 815–820 (2015).
21. Vicente, J. L. *et al.* Massive introgression drives species radiation at the range limit of *Anopheles gambiae*. *Sci. Rep.* **7**, 46451 (2017).
22. Nwakanma, D. C. *et al.* Breakdown in the process of incipient speciation in *Anopheles gambiae*. *Genetics* **193**, 1221–1231 (2013).
23. Li, S., Schlebusch, C. & Jakobsson, M. Genetic variation reveals large-scale population expansion and migration during the expansion of Bantu-speaking peoples. *Proc. R. Soc. Lond. B* **281**, 20141448 (2014).
24. Noor, A. M., Amin, A. A., Akhwale, W. S. & Snow, R. W. Increasing coverage and decreasing inequity in insecticide-treated bed net use among rural Kenyan children. *PLoS Med.* **4**, e255 (2007).
25. Mwangangi, J. M. *et al.* Shifts in malaria vector species composition and transmission dynamics along the Kenyan coast over the past 20 years. *Malar. J.* **12**, 13 (2013).
26. Davies, T. G. E., Field, L. M., Usherwood, P. N. R. & Williamson, M. S. A comparative study of voltage-gated sodium channels in the Insecta: implications for pyrethroid resistance in Anopheline and other Neopteran species. *Insect Mol. Biol.* **16**, 361–375 (2007).
27. Mitchell, S. N. *et al.* Metabolic and target-site mechanisms combine to confer strong DDT resistance in *Anopheles gambiae*. *PLoS ONE* **9**, e29662 (2014).
28. Edi, C. V. *et al.* CYP6 P450 enzymes and ACE-1 duplication produce extreme and multiple insecticide resistance in the malaria mosquito *Anopheles gambiae*. *PLoS Genet.* **10**, e1004236 (2014).
29. Jones, C. M. *et al.* Footprints of positive selection associated with a mutation (N1575Y) in the voltage-gated sodium channel of *Anopheles gambiae*. *Proc. Natl Acad. Sci. USA* **109**, 6614–6619 (2012).
30. Ross, R. Inaugural lecture on the possibility of eradicating malaria from certain localities by a new method. *BMJ* **2**, 1–4 (1899).

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

**Acknowledgements** The authors would like to thank the staff of the Wellcome Trust Sanger Institute Sample Logistics, Sequencing and Informatics facilities for their contributions. This work was supported by the Wellcome Trust (090770/Z/09/Z; 090532/Z/09/Z; 098051) and Medical Research Council UK and the Department for International Development (DFID) (MR/M006212/1). M.K.N.L. was supported by MRC grant G1100339. S.O.'L. and A.B. were supported by a grant from the Foundation for the National Institutes of Health through the Vector-Based Control of Transmission: Discovery Research (VCTR) program of the Grand Challenges in Global Health initiative of the Bill & Melinda Gates Foundation. D.W., C.S.W., H.D.M. and M.J.D. were supported by Award Numbers U19AI089674 and R01AI082734 from the National Institute of Allergy and Infectious Diseases (NIAID). The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIAID or NIH. T.A. was supported by a Sir Henry Wellcome Postdoctoral Fellowship.

**Author Contributions** Details of author contributions are given in the consortium author list.

**Author Information** Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations. Correspondence and requests for materials should be addressed to A.M. (alistair.miles@bdi.ox.ac.uk), M.K.N.L. (mara@sanger.ac.uk), M.J.D. (Martin.Donnely@lstmed.ac.uk) or D.P.K. (dominic@sanger.ac.uk).

**Reviewer Information** *Nature* thanks J. Pool and the other anonymous reviewer(s) for their contribution to the peer review of this work.

## The *Anopheles gambiae* 1000 Genomes Consortium

**Data analysis group** Alistair Miles<sup>1,2</sup> (project lead), Nicholas J. Harding<sup>2</sup>, Giordano Bottà<sup>2,3</sup>, Chris S. Clarkson<sup>1,4</sup>, Tiago Antão<sup>2,4,5</sup>, Krzysztof Kozak<sup>1</sup>, Daniel R. Schrider<sup>6</sup>, Andrew D. Kern<sup>6</sup>, Seth Redmond<sup>7</sup>, Igor Sharakhov<sup>8,9</sup>, Richard D. Pearson<sup>1,2</sup>, Christina Bergey<sup>10</sup>, Michael C. Fontaine<sup>11</sup>, Martin J. Donnelly<sup>1,4</sup>, Mara K. N. Lawniczak<sup>1</sup>, Dominic P. Kwiatkowski<sup>1,2</sup> (chair)

**Partner working group** Martin J. Donnelly<sup>1,4</sup> (chair), Diego Ayala<sup>12,13</sup>, Nora J. Besansky<sup>10</sup>, Austin Burt<sup>14</sup>, Beniamino Caputo<sup>3</sup>, Alessandra della Torre<sup>3</sup>, Michael C. Fontaine<sup>11</sup>, H. Charles J. Godfray<sup>15</sup>, Matthew W. Hahn<sup>16</sup>, Andrew D. Kern<sup>6</sup>, Dominic P. Kwiatkowski<sup>1,2</sup>, Mara K. N. Lawniczak<sup>1</sup>, Janet Midega<sup>17</sup>, Daniel E. Neafsey<sup>7</sup>, Samantha O'Loughlin<sup>14</sup>, João Pinto<sup>18</sup>, Michelle M. Riehle<sup>19</sup>, Igor Sharakhov<sup>8,9</sup>, Kenneth D. Vernick<sup>20</sup>, David Weetman<sup>4</sup>, Craig S. Wilding<sup>4,21</sup>, Bradley J. White<sup>22</sup>

**Sample collections—Angola:** Arlete D. Troco<sup>23</sup>, João Pinto<sup>18</sup>, **Burkina Faso:** Abdoulaye Diabaté<sup>24</sup>, Samantha O'Loughlin<sup>14</sup>, Austin Burt<sup>14</sup>; **Cameroon:** Carlo Costantini<sup>13,25</sup>, Kyanne R. Rohatgi<sup>10</sup>, Nora J. Besansky<sup>10</sup>; **Gabon:** Nohal Elissa<sup>12</sup>, João Pinto<sup>18</sup>; **Guinea:** Boubacar Coulibaly<sup>26</sup>, Michelle M. Riehle<sup>19</sup>, Kenneth D. Vernick<sup>20</sup>; **Guinea-Bissau:** João Pinto<sup>18</sup>, João Dinis<sup>27</sup>; **Kenya:** Janet Midega<sup>17</sup>, Charles Mbogo<sup>17</sup>, Philip Bejon<sup>17</sup>; **Uganda:** Craig S. Wilding<sup>4,21</sup>, David Weetman<sup>4</sup>, Henry D. Mawejje<sup>28</sup>, Martin J. Donnelly<sup>1,4</sup>; **Crosses:** David Weetman<sup>4</sup>, Craig S. Wilding<sup>4,21</sup>, Martin J. Donnelly<sup>1,4</sup>

**Sequencing and data production** Jim Stalker<sup>1</sup>, Kirk Rockett<sup>2</sup>, Eleanor Drury<sup>1</sup>, Daniel Mead<sup>1</sup>, Anna Jeffreys<sup>2</sup>, Christina Hubbard<sup>2</sup>, Kate Rowlands<sup>2</sup>, Alison T. Isaacs<sup>4</sup>, Dushyanth Jyothi<sup>1</sup>, Cinzia Malangone<sup>1</sup>

**Web application development** Paul Vauterin<sup>2</sup>, Ben Jeffery<sup>2</sup>, Ian Wright<sup>2</sup>, Lee Hart<sup>2</sup>, Krzysztof Kluczyński<sup>2</sup>

**Project coordination** Victoria Cornelius<sup>2</sup>, Bronwyn MacInnis<sup>29</sup>, Christa Henrichs<sup>2</sup>, Rachel Giacomantonio<sup>1</sup> & Dominic P. Kwiatkowski<sup>1,2</sup>

<sup>1</sup>Malaria Programme, Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA, UK.

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

<sup>3</sup>Istituto Pasteur Italia – Fondazione Cenci Bolognietti, Dipartimento di Sanità Pubblica e Malattie Infettive, Università di Roma SAPIENZA, Rome, Italy. <sup>4</sup>Department of Vector Biology, Liverpool School of Tropical Medicine, Pembroke Place, Liverpool L3 5QA, UK. <sup>5</sup>University of Montana, Missoula, Montana 59812, USA. <sup>6</sup>Department of Genetics, Rutgers University, 604 Alison Road, Piscataway, New Jersey 08854, USA. <sup>7</sup>Genome Sequencing and Analysis Program, Broad Institute, 415 Main Street, Cambridge, Maryland 02142, USA. <sup>8</sup>Department of Entomology, Virginia Tech, Blacksburg, Virginia 24061, USA. <sup>9</sup>Laboratory of Ecology, Genetics and Environmental Protection, Tomsk State University, Tomsk 634050, Russia. <sup>10</sup>Eck Institute for Global Health, Department of Biological Sciences, University of Notre Dame, Indiana 46556, USA. <sup>11</sup>Groningen Institute for Evolutionary Life Sciences (GELIFES), Nijenborgh 7, 9747 AG Groningen, The Netherlands. <sup>12</sup>Unité d'Ecologie des Systèmes Vectoriels, Centre International de Recherches Médicales de Franceville, Franceville, Gabon. <sup>13</sup>Institut de Recherche pour le Développement (IRD), UMR MIVEGEC (UM1, UM2, CNRS 5290, IRD 224), Montpellier, France. <sup>14</sup>Department of Life Sciences, Imperial College, Silwood Park, Ascot, Berkshire SL5 7PY, UK.

<sup>15</sup>Department of Zoology, University of Oxford, The Tinbergen Building, South Parks Road, Oxford OX1 3PS, UK. <sup>16</sup>Department of Biology and School of Informatics and Computing, Indiana University, Bloomington, Indiana 47405, USA. <sup>17</sup>KEMRI-Wellcome Trust Research Programme, PO Box 230, Bofa Road, Kilifi, Kenya. <sup>18</sup>Global Health and Tropical Medicine, GHM, Instituto de Higiene e Medicina Tropical, IHMT, Universidade Nova de Lisboa, UNL, Rua da Junqueira 100, 1349-008 Lisbon, Portugal. <sup>19</sup>Department of Microbiology and Immunology, Microbial and Plant Genomics Institute, University of Minnesota, St Paul, Minnesota 55108, USA. <sup>20</sup>Unit for Genetics and Genomics of Insect Vectors, Institut Pasteur, Paris, France. <sup>21</sup>School of Natural Sciences and Psychology, Liverpool John Moores University, Liverpool L3 3AF, UK. <sup>22</sup>Department of Entomology, University of California, Riverside, California, USA. <sup>23</sup>Programa Nacional de Controle da Malária, Direção Nacional de Saúde Pública, Ministério da Saúde, Luanda, Angola. <sup>24</sup>Institut de Recherche en Sciences de la Santé (IRSS), Bobo Dioulasso, Burkina Faso. <sup>25</sup>Laboratoire de Recherche sur le Paludisme, Organisation de Coordination pour la lutte contre les Endémies en Afrique Centrale (OCEAC), Yaoundé, Cameroon. <sup>26</sup>Malaria Research and Training Centre (MRTC), University of Bamako, Bamako, Mali. <sup>27</sup>Instituto Nacional de Saúde Pública, Ministério da Saúde Pública, Bissau, Guiné-Bissau. <sup>28</sup>Infectious Diseases Research Collaboration, 2C Nakasero Hill Road, PO Box 7475, Kampala, Uganda. <sup>29</sup>The Broad Institute of Massachusetts Institute of Technology and Harvard, 415 Main Street, Cambridge, Massachusetts 02142, USA.

## METHODS

**Population sampling.** Mosquitoes were collected from natural populations at 15 sampling sites in 8 African countries (Extended Data Fig. 1). Sampling locations, dates, specimen collection methods and DNA extraction methods are given in Supplementary Information 1.1. We also performed genetic crosses between adult mosquitoes obtained from laboratory colonies (Supplementary Information 1.2). Parents and progeny of four crosses were contributed to Ag1000G phase 1 (Extended Data Fig. 1). No statistical methods were used to predetermine sample size.

**Whole-genome sequencing.** Sequencing was performed on the Illumina HiSeq 2000 platform at the Wellcome Trust Sanger Institute. Paired-end multiplex libraries were prepared using the manufacturer's protocol, with the exception that genomic DNA was fragmented using Covaris Adaptive Focused Acoustics rather than nebulization. Multiplexes consisted of 12 tagged individual mosquitoes, and 3 lanes of sequencing were generated for each multiplex to even out variation in yield between sequencing runs. Cluster generation and sequencing were undertaken per the manufacturer's protocol for paired-end 100-bp sequence reads with insert size in the range 100–200 bp.

**Sequence analysis and variant calling.** Sequence reads were aligned to the AgamP3 reference genome<sup>31</sup> using bwa<sup>32</sup> and SNPs were discovered using GATK following best practice recommendations<sup>33,34</sup> (Supplementary Information 3.1–3.2). After sample quality control, we analysed data on 765 wild-caught specimens and a further 80 specimens comprising parents and progeny from the four laboratory crosses (Supplementary Information 3.3). The alignments were also used to identify genome regions accessible to SNP calling, where short reads could be uniquely mapped and there was minimal evidence for structural variation (Supplementary Information 3.4). Mendelian errors in the crosses were used to guide the design of filters to remove poor-quality variant calls (Supplementary Information 3.5). We performed capillary sequencing of five genes in 58 individual mosquitoes to provide an estimate for the SNP false discovery rate, sensitivity and genotyping accuracy (Supplementary Information 3.6). We also performed genotyping by primer-extension mass spectrometry using the Sequenom MassARRAY platform at 158 SNPs in 229 individual mosquitoes to provide a second estimate for genotyping accuracy (Supplementary Information 3.7).

**Haplotype estimation.** We used SHAPEIT2 to perform statistical phasing with information from sequence reads<sup>35</sup> for all wild-caught individuals (Supplementary Information 4.1). We assessed phasing performance by comparison with haplotypes generated from the crosses and from male X chromosome haplotypes (Supplementary Information 4.2; Extended Data Fig. 2b, c).

**Population structure.** To investigate variation in patterns of relatedness along the genome, we performed a windowed analysis using genetic distance and neighbour-joining trees (NJT). We divided the genome into 1,418 contiguous non-overlapping windows, where each window contained 100 kb of accessible positions. Within each window, we computed the city-block distance between all pairs of individuals. We used these distance matrices to construct a NJT for each window. We then computed the Pearson correlation coefficient between all pairs of distance matrices, and performed a singular value decomposition (SVD) on the correlation matrix. The resulting SVD components were used to identify major patterns of relatedness (Supplementary Information 6.1). We analysed geographical population structure using ADMIXTURE<sup>36</sup> and PCA<sup>37</sup>. For these analyses, we used biallelic SNPs from within the regions 3R: 1–37 Mb and 3L: 15–41 Mb and with minor allele frequency  $\geq 1\%$ , then each chromosome arm was randomly down-sampled to 100,000 variants using 10 different random seeds to provide 10 replicate variant sets, and then each set was pruned to remove variants in linkage disequilibrium (Supplementary Information 6.2). For each of the 10 replicate variant sets, ADMIXTURE was run for  $K$  (number of ancestral populations) from 2 to 11 with fivefold cross-validation. Each ADMIXTURE analysis was repeated 10 times with different seeds, resulting in a total of 100 runs for each value of  $K$ . We then used CLUMPAK<sup>38</sup> to analyse the ADMIXTURE results and compute ancestry proportions (Supplementary Information 6.2). Average  $F_{ST}$  values were computed using Hudson's estimator and the ratio of averages, and standard errors were computed using a block-jackknife<sup>39</sup> (Supplementary Information 6.4). AIMS were ascertained by starting with SNPs previously discovered in Mali<sup>18</sup> with an allele frequency difference between *An. gambiae* and *An. coluzzii*  $> 0.9$ , then taking the intersection with biallelic SNPs discovered in this study, resulting in 506 AIMS (Supplementary Information 6.6).

**Population size history.** We inferred the scale and timing of historical changes in  $N_e$  using two methods, stairway plot<sup>40</sup> and  $\partial a \partial i$ <sup>41</sup>, both using site frequency spectra but taking different modelling approaches. To compute site frequency spectra, we used SNPs from within the regions 3R: 1–37 Mb and 3L: 15–41 Mb, taking only intergenic SNPs at least 5 kb from the nearest gene (Supplementary Information 8). We modified the stairway plot to include an additional parameter representing

the probability of ancestral misclassification for each SNP (Supplementary Information 8.1). We fitted a three-epoch (two  $N_e$  changes)  $\partial a \partial i$  model for each population singly, and fitted joint population models for selected pairs of populations (Supplementary Information 8.2). Scaling of parameters assumed that the *Anopheles* mutation rate is within the range of values estimated for *Drosophila*, in which estimates<sup>42,43</sup> range from  $2.8 \times 10^{-9}$  to  $5.5 \times 10^{-9}$ . For joint population models, we computed the joint site frequency spectrum for each pair of populations from the same set of SNPs used for single-population inferences. Joint population models allowed for a phase of exponential size change in the ancestral population up until the time of the population split, after which each of the daughter populations experienced their own exponential size change until the present. We fitted these models with and without the addition of a symmetric, bidirectional migration rate parameter following the split. To study recent population history in Kenya we used IBDseq<sup>44</sup> to infer genome tracts IBD, then ran the IBDNe program (<http://faculty.washington.edu/browning/ibdne.html>)<sup>45</sup> to infer population size history (Supplementary Information 8.4).

**Recent selection.** To scan the genome for signals of recent selection, we computed the H12 haplotype diversity statistic<sup>46</sup> for each population, and the cross-population extended haplotype homozygosity (XP-EHH) score<sup>47</sup> for selected pairs of populations. H12 was computed in non-overlapping windows over the genome, in which each window contained a fixed number of SNPs, and window sizes were calibrated separately for each population to account for differences in the extent of linkage disequilibrium (Supplementary Information 9.1). XP-EHH was computed for all SNPs with a minor allele frequency  $\geq 5\%$  in the union of both populations in each pair, and normalized within each chromosome (Supplementary Information 9.2). To study haplotype structure at the *Vgsc*, *Gste* and *Cyp6p* loci, we computed the Hamming distance between all pairs of haplotypes, then performed hierarchical clustering of haplotypes (Supplementary Information 9.3). To identify haplotype clusters resulting from recent selection, we cut the dendrograms at a small genetic distance (0.0004 SNP differences per accessible base pair) and studied the largest clusters obtained after cutting. To look for evidence that the haplotype clusters we identified were related via recombination events, we performed the same clustering analysis but in non-overlapping windows upstream and downstream of the target region and compared the resulting clusters.

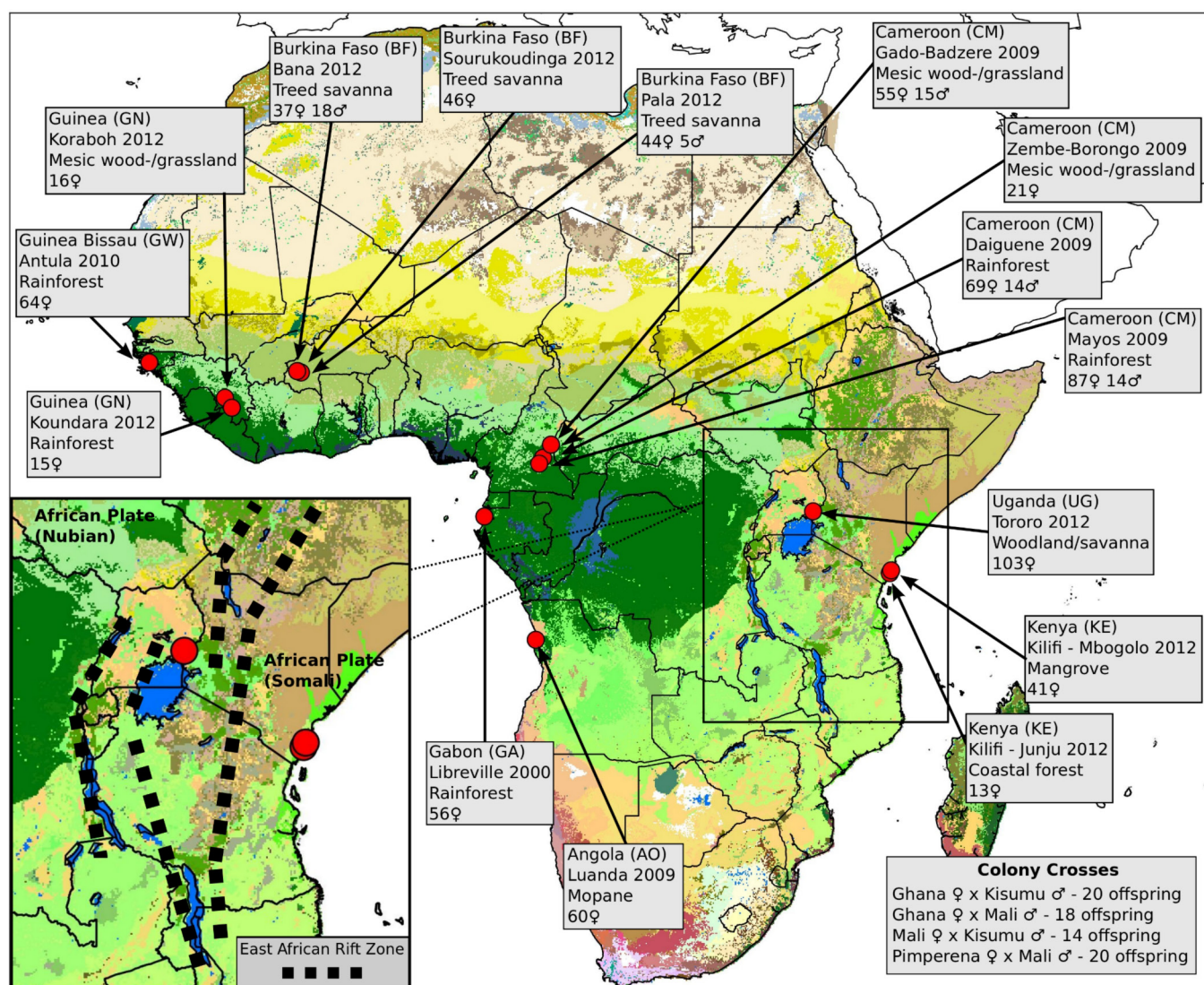
**Plotting and maps.** All figures were produced using the matplotlib package for Python<sup>48</sup>. The map component of Fig. 2 was produced via the matplotlib basemap package, using the NASA Blue Marble image as the map background. The map components of Fig. 4 and Extended Data Fig. 10 were plotted via the cartopy package, using the Natural Earth shaded relief raster as the map background. The map in Extended Data Fig. 1 was plotted via the cartopy package, using data from the map of standardized terrestrial ecosystems of Africa<sup>49</sup> as the map background.

**Data availability.** Sequence read alignments and variant calls from Ag1000G phase 1 are available from the European Nucleotide Archive (ENA; <http://www.ebi.ac.uk/ena>) under study PRJEB18691. Variant and haplotype calls and associated data from Ag1000G phase 1 can be explored via an interactive web application or downloaded via the MalariaGEN website (<https://www.malariagen.net/projects/ag1000g/data>). All other data are available from the corresponding authors upon reasonable request.

- Sharakhova, M. V. et al. Update of the *Anopheles gambiae* PEST genome assembly. *Genome Biol.* **8**, R5 (2007).
- Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
- DePristo, M. A. et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
- Van der Auwera, G. A. et al. From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinformatics* **11**, 11.10.1–11.10.33 (2013).
- Delaneau, O., Howie, B., Cox, A. J., Zagury, J.-F. & Marchini, J. Haplotype estimation using sequencing reads. *Am. J. Hum. Genet.* **93**, 687–696 (2013).
- Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* **19**, 1655–1664 (2009).
- Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS Genet.* **2**, e190 (2006).
- Kopelman, N. M., Mayzel, J., Jakobsson, M., Rosenberg, N. A. & Mayrose, I. Clumpak: a program for identifying clustering modes and packaging population structure inferences across *K*. *Mol. Ecol. Resour.* **15**, 1179–1191 (2015).
- Bhatia, G., Patterson, N., Sankararaman, S. & Price, A. L. Estimating and interpreting  $F_{ST}$ : the impact of rare variants. *Genome Res.* **23**, 1514–1521 (2013).
- Liu, X. & Fu, Y.-X. Exploring population size changes using SNP frequency spectra. *Nat. Genet.* **47**, 555–559 (2015).
- Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H. & Bustamante, C. D. Inferring the joint demographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genet.* **5**, e1000695 (2009).

42. Keightley, P. D., Ness, R. W., Halligan, D. L. & Haddrill, P. R. Estimation of the spontaneous mutation rate per nucleotide site in a *Drosophila melanogaster* full-sib family. *Genetics* **196**, 313–320 (2014).
43. Schrider, D. R., Houle, D., Lynch, M. & Hahn, M. W. Rates and genomic consequences of spontaneous mutational events in *Drosophila melanogaster*. *Genetics* **194**, 937–954 (2013).
44. Browning, B. L. & Browning, S. R. Detecting identity by descent and estimating genotype error rates in sequence data. *Am. J. Hum. Genet.* **93**, 840–851 (2013).
45. Browning, S. R. & Browning, B. L. Accurate non-parametric estimation of recent effective population size from segments of identity by descent. *Am. J. Hum. Genet.* **97**, 404–418 (2015).
46. Garud, N. R., Messer, P. W., Buzbas, E. O. & Petrov, D. A. Recent selective sweeps in North American *Drosophila melanogaster* show signatures of soft sweeps. *PLoS Genet.* **11**, e1005004 (2015).
47. Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive selection in human populations. *Nature* **449**, 913–918 (2007).
48. Hunter, J. D. Matplotlib: a 2D graphics environment. *Comput. Sci. Eng.* **9**, 90–95 (2007).
49. Sayre, R. G. *et al.* *A New Map of Standardized Terrestrial Ecosystems of Africa* (American Association of Geographers, 2013).
50. Sharakhova, M. V. *et al.* Genome mapping and characterization of the *Anopheles gambiae* heterochromatin. *BMC Genomics* **11**, 459 (2010).

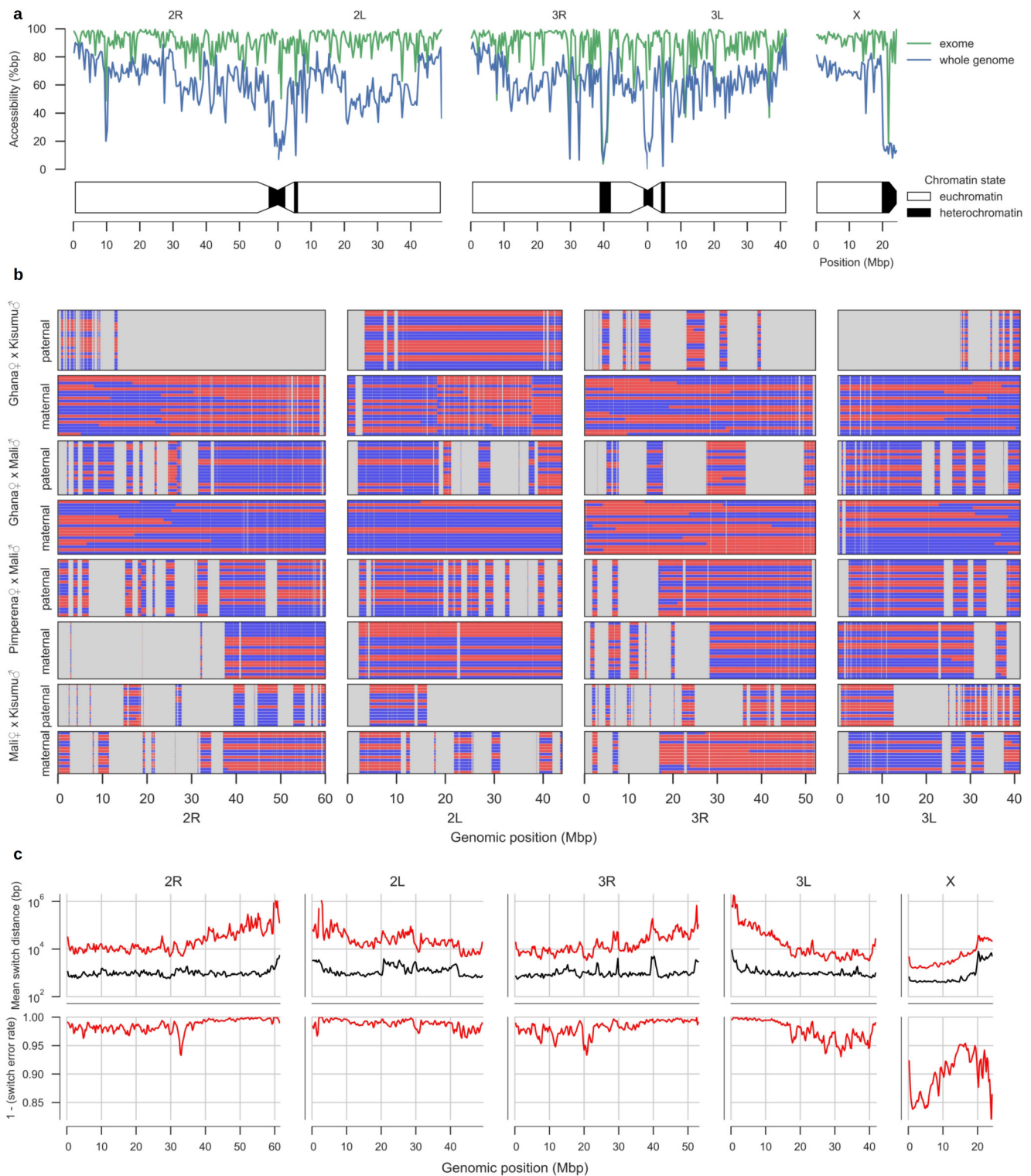




Country	Species	No. individuals	Sequencing depth (X)
Angola	<i>An. coluzzii</i>	60♀	30 (22 – 38)
Burkina Faso	<i>An. coluzzii</i>	66♀, 3♂	33 (26 – 44)
Guinea-Bissau		46♀	30 (24 – 41)
Guinea	<i>An. gambiae</i>	31♀	29 (19 – 39)
Burkina Faso	<i>An. gambiae</i>	61♀, 20♂	31 (22 – 68)
Cameroon	<i>An. gambiae</i>	232♀, 43♂	29 (19 – 53)
Gabon	<i>An. gambiae</i>	56♀	31 (21 – 37)
Uganda	<i>An. gambiae</i>	103♀	31 (27 – 39)
Kenya		44♀	31 (18 – 57)

**Extended Data Figure 1 | Overview of population sampling.** Red circles show sampling locations for wild-caught mosquitoes. Colours in the map represent ecosystem classes; dark green represents forest ecosystems; see figure 9 in ref. 49 for a complete colour legend. The Congo Basin tropical rainforest is the large region of dark green in Central Africa. Sampling details for each site are shown in light grey boxes, including country (two-letter country code), location and year of collection, predominant ecosystem classification for the local region, and number and sex of

individuals sequenced. For colony crosses, the direction of cross (colony of origin of mother and father) and number of offspring is shown. The inset map depicts geological fault lines in the East African rift system ([http://pubs.usgs.gov/publications/text/East\\_Africa.html](http://pubs.usgs.gov/publications/text/East_Africa.html)). Species assignment for Guinea-Bissau and Kenya specimens is uncertain, see main text. Sequencing depth per individual is shown as median (5th–95th percentile) for each population.

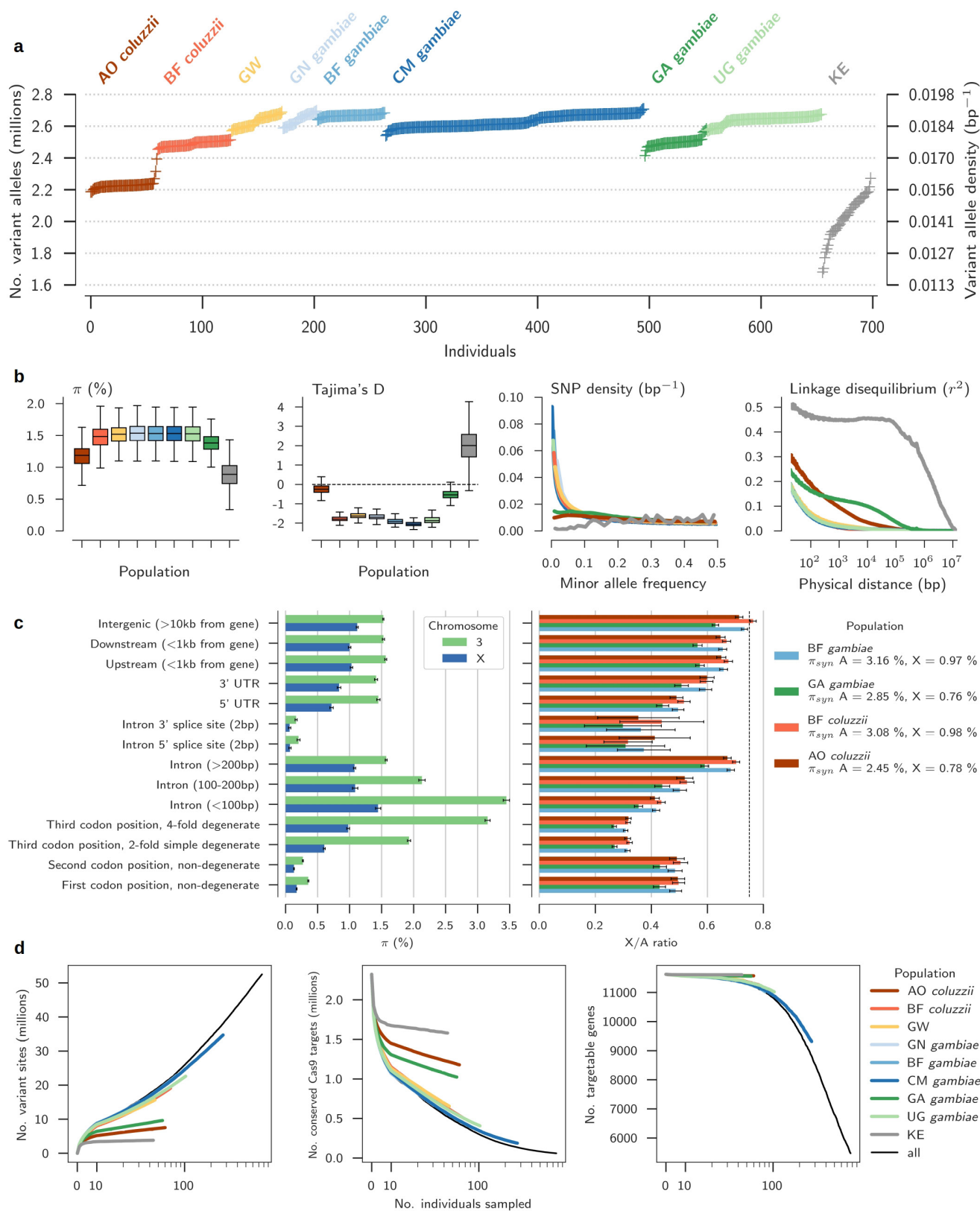


### Extended Data Figure 2 | Genome accessibility and haplotype

**validation.** **a**, Percentage of accessible bases in non-overlapping 400-kb windows. The schematic of chromosomes below shows chromatin state predictions from ref. 50. **b**, Haplotypes inferred in the crosses. Each panel shows either maternal or paternal haplotypes from a single cross. Each row within a panel represents a single progeny haplotype. Haplotypes are coloured by parental inheritance (blue denotes allele from parent's first chromosome; red denotes allele from parent's second chromosome).

Switches between colours along a haplotype indicate recombination events. Regions that were within a run of homozygosity in the parent and thus not informative for haplotype validation are masked in grey. **c**, Error rate estimates for haplotypes inferred in wild-caught individuals. Top plots show estimates for the mean switch distance (red line), compared to the mean switch distance if heterozygotes were phased randomly (black line). Bottom plots show the switch error rate (probability of a switch error occurring between two adjacent heterozygous genotype calls).



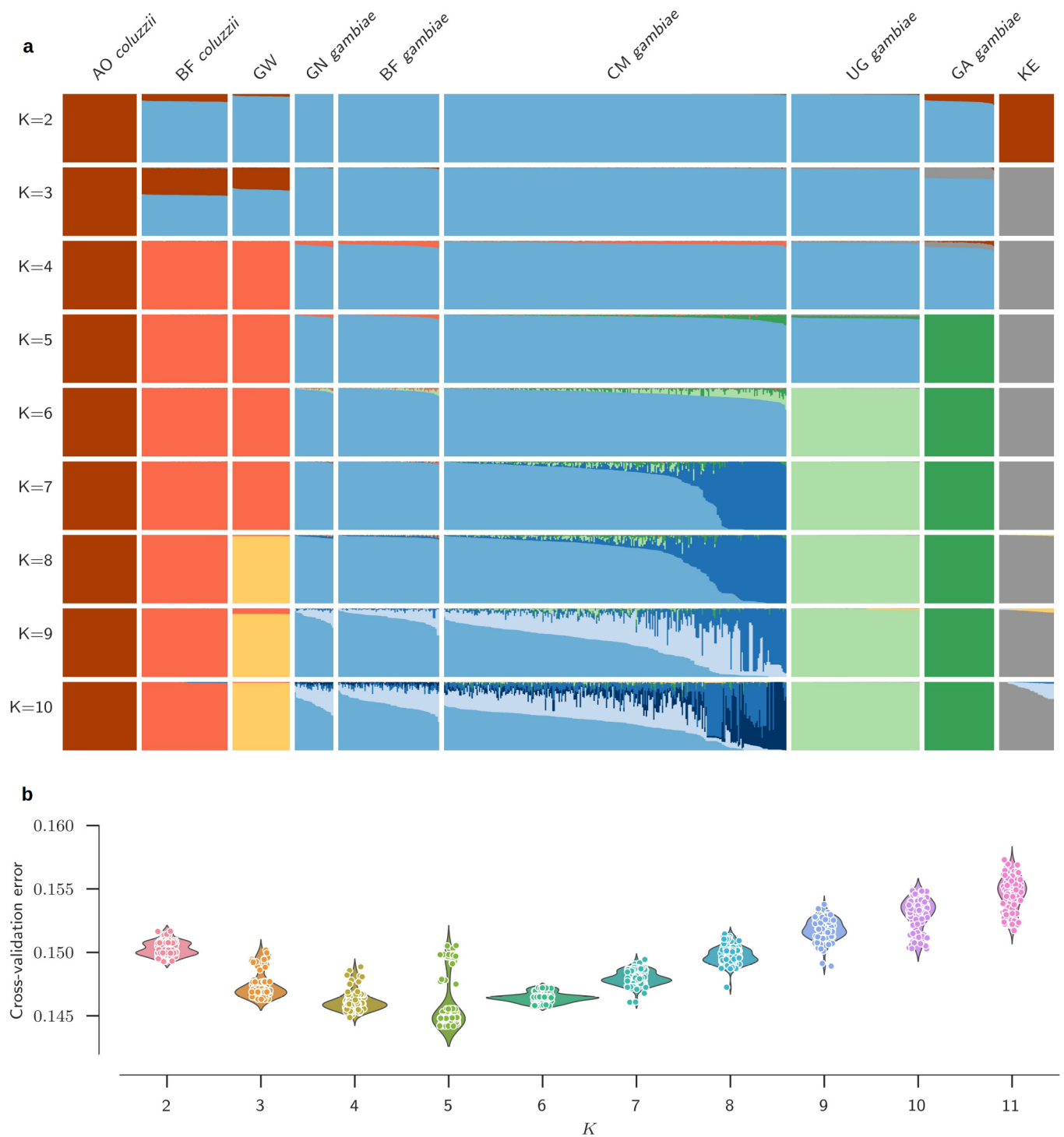


### Extended Data Figure 3 | Variant discovery and nucleotide diversity.

**a**, Number of variant alleles discovered per individual mosquito. Only females are plotted. **b**, Genetic diversity within populations. Nucleotide diversity ( $\pi$ ) and Tajima's  $D$  were calculated in non-overlapping 20-kb genomic windows. SNP density depicts 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

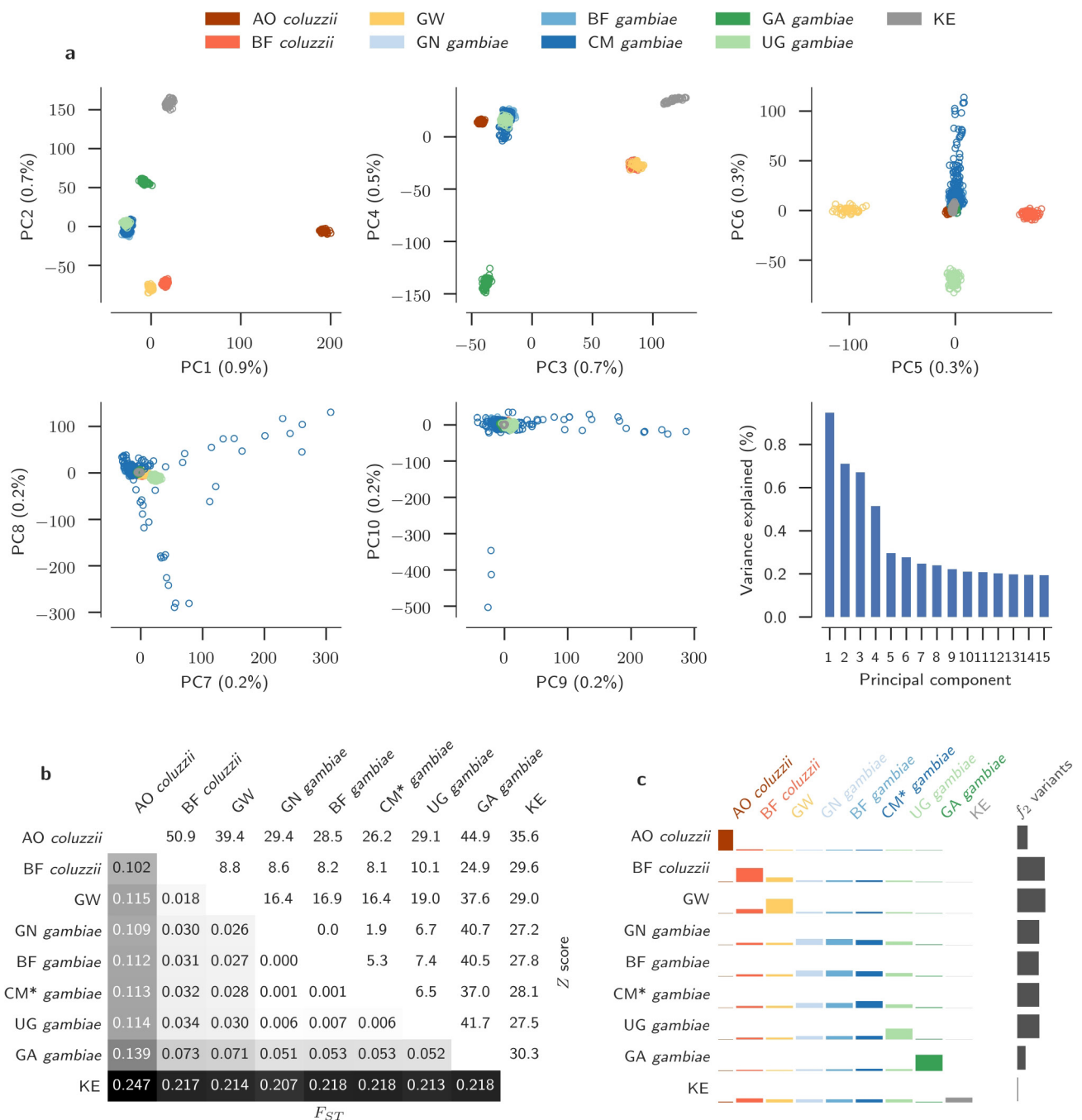
constant SNP density over all allele frequencies. **c**, Average nucleotide diversity ( $\pi$ ) and ratio of diversity between sex-linked (X) and autosomal (A) chromosomes in relation to gene architecture. **d**, Relationship between number of individuals sampled and the cumulative number of variant sites discovered (left), availability of conserved Cas9 target sites within genes (centre), and number of genes containing at least 1 conserved Cas9 target site which could thus be 'targetable' for gene drive (right).





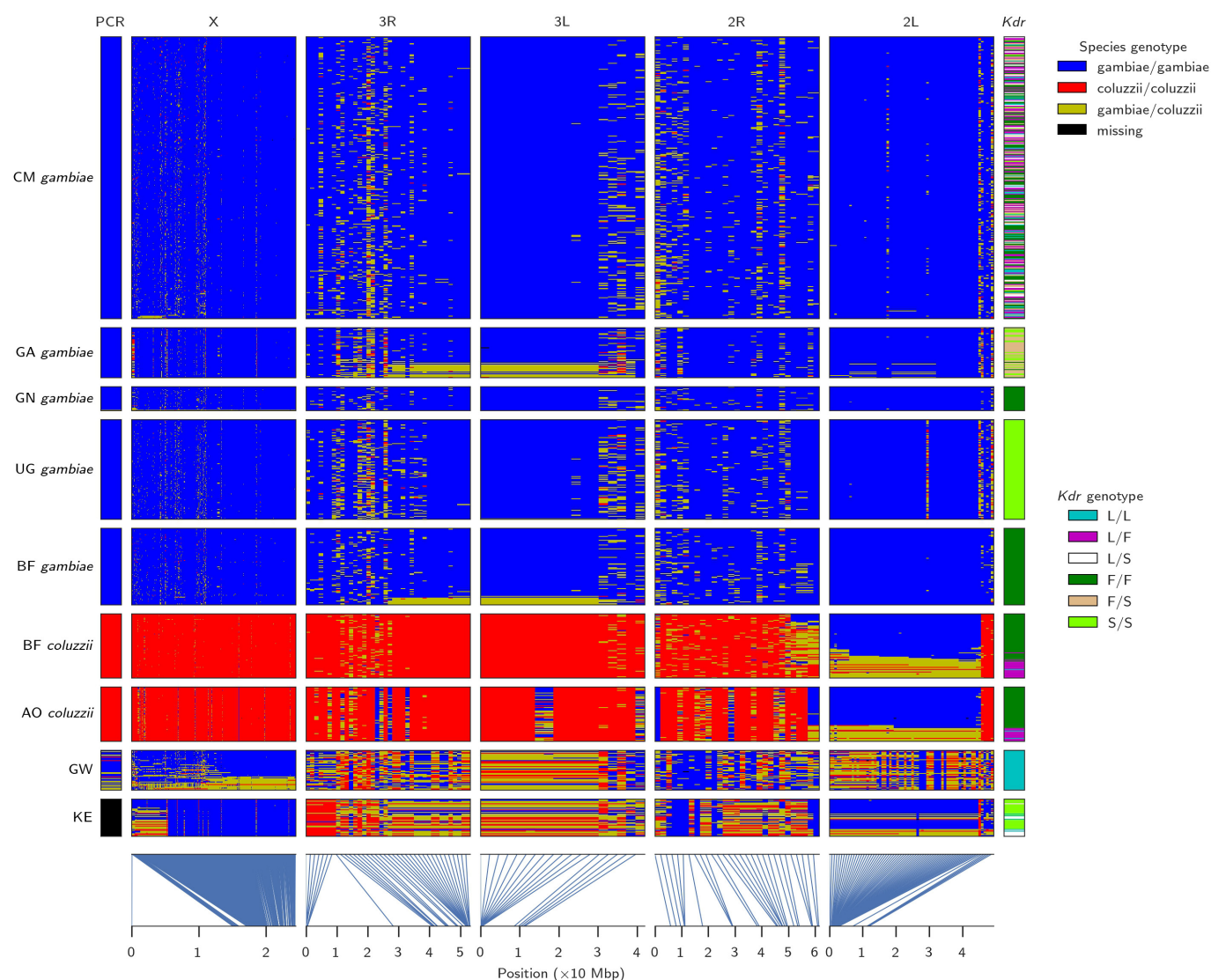
**Extended Data Figure 4 | ADMIXTURE analysis. a,** Ancestry proportions within individual mosquitoes for ADMIXTURE models from  $K=2$  to  $K=10$  ancestral populations. Each vertical bar represents the proportion of ancestry within a single individual, with colours

corresponding to ancestral populations. These data are the average of the major  $q$ -matrix clusters derived by CLUMPAK analysis. **b,** Violin plot of cross-validation error for each of 100 replicates for each  $K$  value.



**Extended Data Figure 5 | Population structure and differentiation.**  
**a**, Principal components analysis of the 765 wild-caught mosquitoes.  
**b**, 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

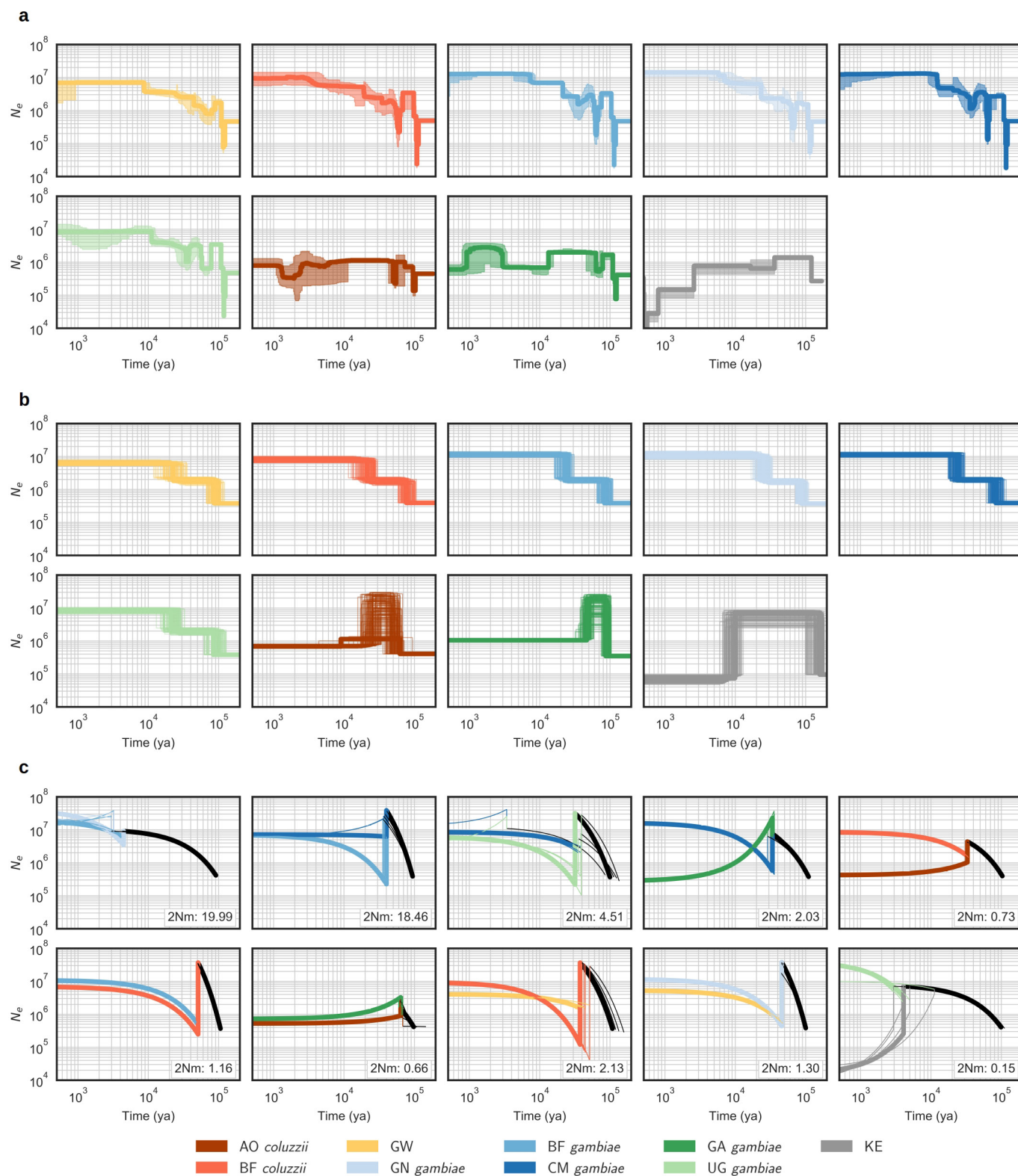
Z score for each  $F_{ST}$  value estimated via a block-jackknife procedure<sup>39</sup>. CM\* denotes Cameroon savannah sampling site only. **c**, Allele sharing in doubleton ( $f_2$ ) variants. The height of the coloured bars represent the probability of sharing a doubleton allele between two populations. Heights are normalized row-wise for each population.



**Extended Data Figure 6 | Ancestry informative markers.** 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 shows the species assignment according to the conventional molecular test based on a single marker on

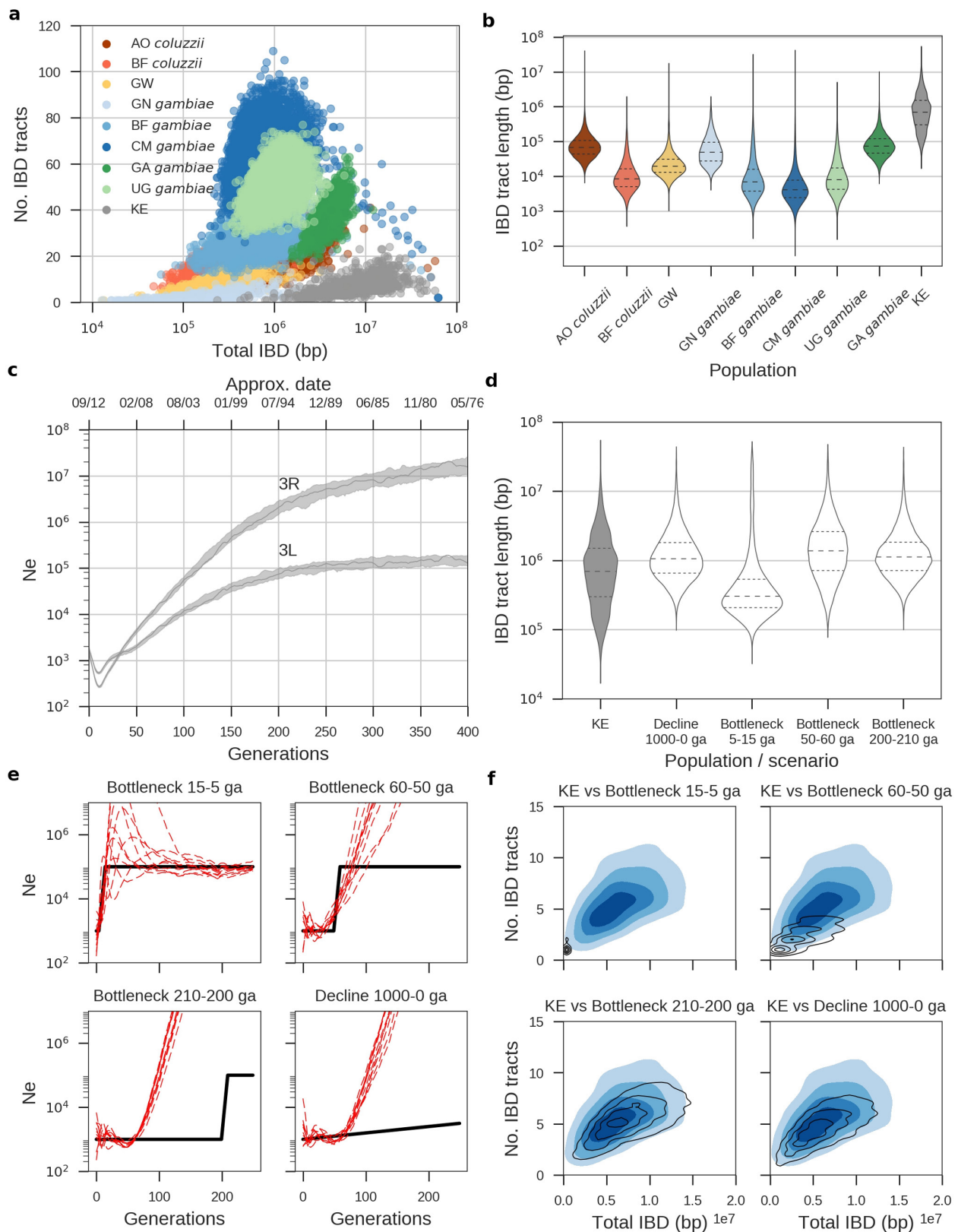
the X chromosome, which was performed for all individuals except Kenya (KE). The column at the far right shows the genotype for *kdr* variants in *Vgsc* codon 995. Lines at the lower edge show the physical locations of the AIM SNPs.





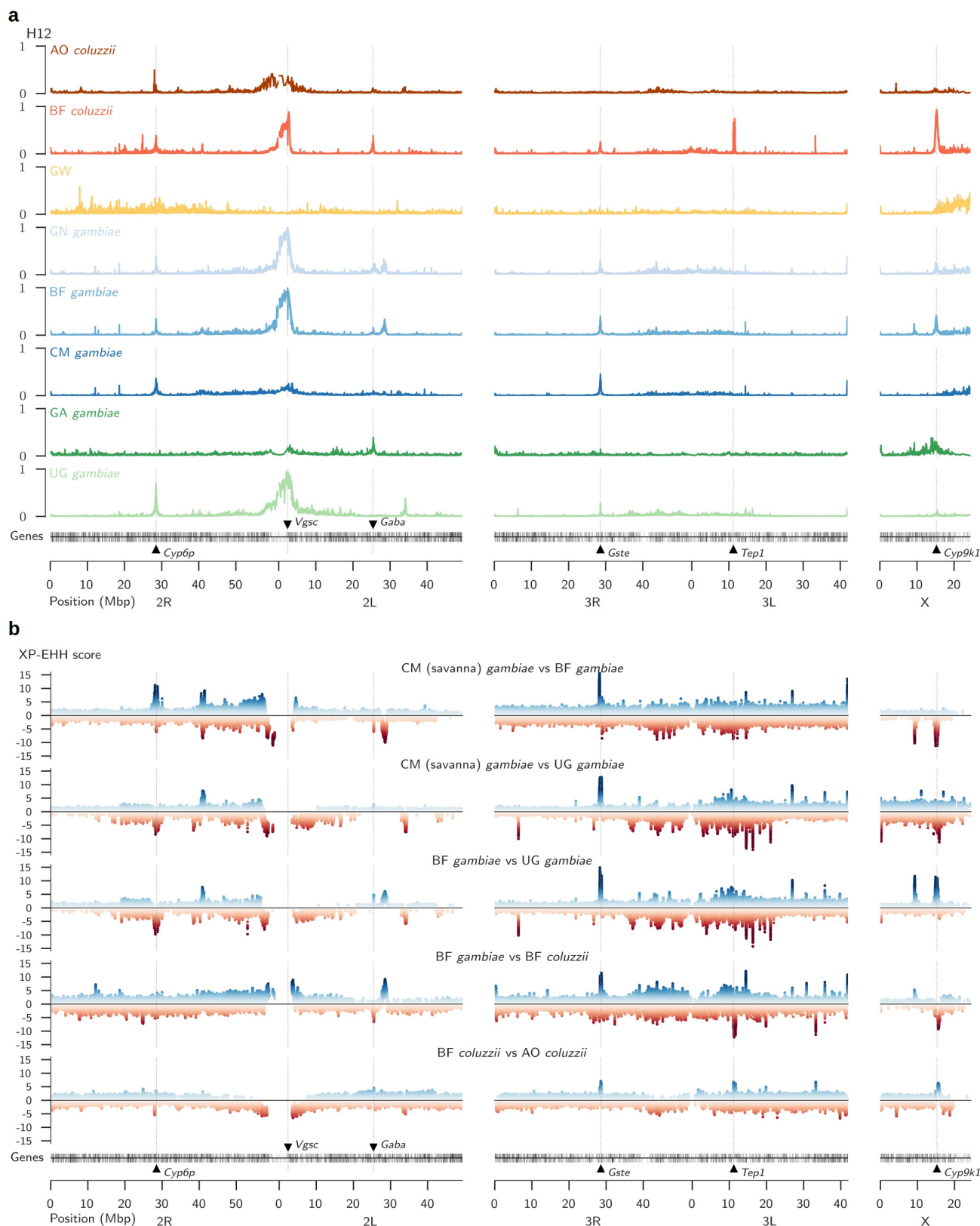
**Extended Data Figure 7 | Population size history.** **a**, Stairway plot of inferred histories for each population. The shaded area shows the 95% confidence interval from 199 bootstrap replicates. **b**, Inferred histories from three-epoch  $\partial a\partial i$  models<sup>41</sup>. The thick line shows the history with the highest likelihood found by optimization; thin lines show 100 histories with the highest likelihoods from even sampling of the model parameter space. **c**, Inferred histories from  $\partial a\partial i$  two-population models allowing for

migration. For each population pair, solutions from 5 optimization runs with the highest likelihoods are shown, with the thick line showing the history with the highest likelihood. In all panels, time and  $N_e$  are scaled assuming 11 generations per year and a mutation rate of  $\mu = 3.5 \times 10^{-9}$ . Scaling of time and  $N_e$  is proportional to  $1/\mu$ , for example, if the true mutation rate is twice as high then estimates of time and  $N_e$  would be halved. ya, years ago.



**Extended Data Figure 8 | Identity by descent and recent effective population size history.** **a**, Patterns of IBD sharing within populations. Each marker represents a pair of individuals. **b**, The distribution of IBD tract lengths within populations. **c**, Recent population size history for the Kenyan population inferred by the IBDNe program<sup>45</sup>. **d**, Comparison of the IBD tract length distribution between Kenya and four simulated

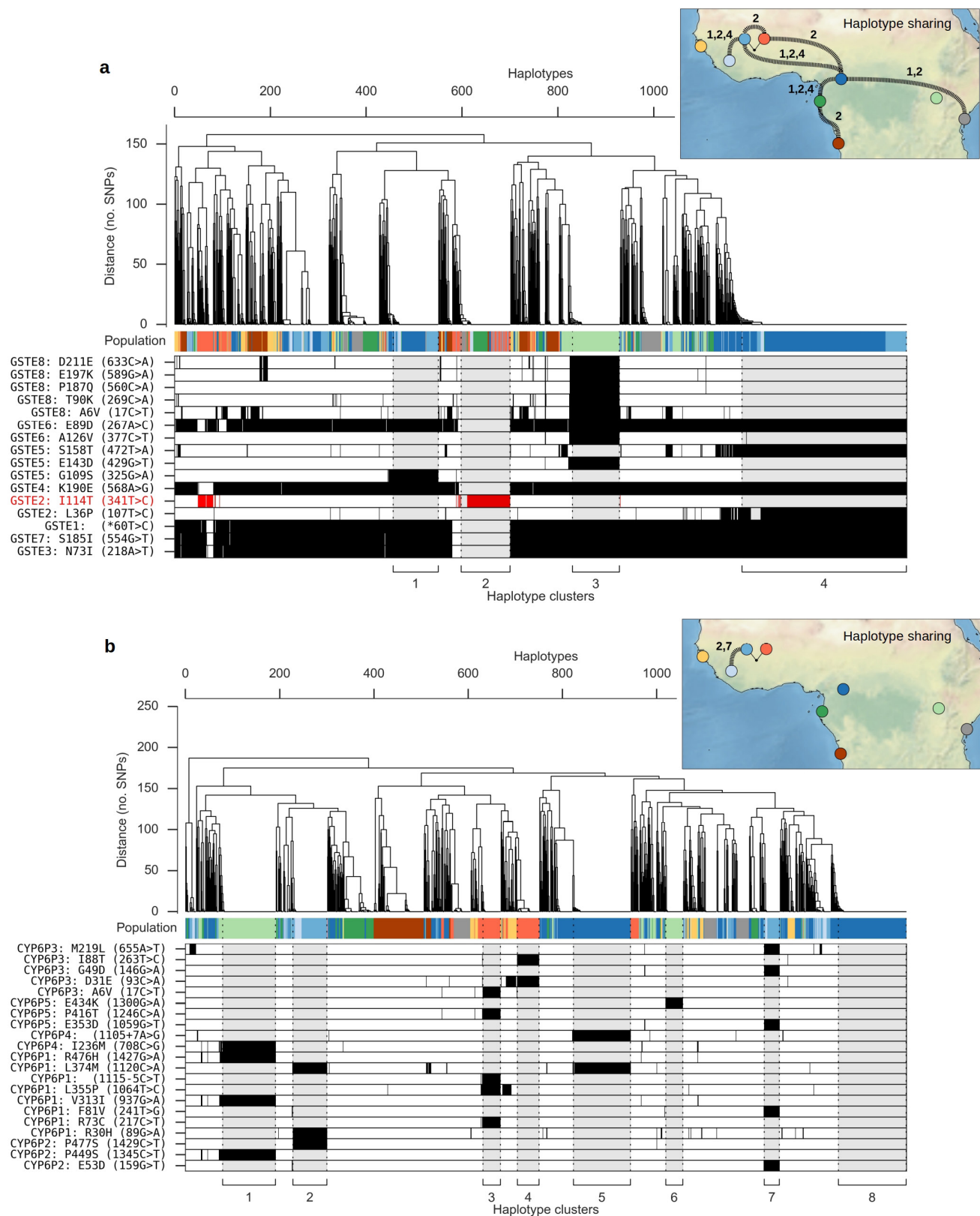
demographic scenarios. **e**, Population size histories inferred by IBDNe (red dashed lines) from data generated by simulations (black line shows the simulated population size history). **f**, Comparison of patterns of IBD sharing generated by simulations (black contour lines) with Kenyan data (filled blue contours). See Supplementary Information 8.4 for details of simulations. ga, generations ago.



**Extended Data Figure 9 | Genome scans for signatures of recent selection. a,** Haplotype diversity. Each track plots the H12 statistic in non-overlapping windows over the genome. A value of 1 indicates low haplotype diversity within a window, expected if one or two haplotypes have risen to high frequency owing to recent selection. A value of 0 indicates high haplotype diversity, expected in neutral regions. **b,** XP-EHH

scans. For each population comparison (for example, BF *gambiae* versus BF *coluzzii*), positive scores indicate longer haplotypes and therefore recent selection in the first population (for example, BF *gambiae*), and negative scores indicate selection in the second population (for example, BF *coluzzii*).





**Extended Data Figure 10 | Haplotype structure at metabolic insecticide-resistance loci.** Plot components are as described for Fig. 4. For both loci, SNPs shown in the bottom panel are all either non-synonymous or splice site variants, and are associated with one or more haplotypes under

selection. **a**, Haplotype clustering using 1,375 SNPs within the region 3R: 28,591,663–28,602,280 spanning 8 genes (*Gste1–Gste8*). **b**, Haplotype clustering using 1,844 SNPs within the region 2R: 28,491,415–28,502,910 spanning 5 genes (*Cyp6p1–Cyp6p5*).

## Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form is intended for publication with all accepted life science papers and provides structure for consistency and transparency in reporting. Every life science submission will use this form; some list items might not apply to an individual manuscript, but all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### ► Experimental design

#### 1. Sample size

Describe how sample size was determined.

Supplementary Text 1.1.

#### 2. Data exclusions

Describe any data exclusions.

Supplementary Text 3.3.

#### 3. Replication

Describe whether the experimental findings were reliably reproduced.

N/A

#### 4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

N/A

#### 5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

N/A

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

#### 6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or in the Methods section if additional space is needed).

- |                                     |  |
|-------------------------------------|--|
| n/a                                 | Confirmed  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The <u>exact sample size</u> ( <i>n</i> ) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)                               |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A statement indicating how many times each experiment was replicated  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section) |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as an adjustment for multiple comparisons  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The test results (e.g. <i>P</i> values) given as exact values whenever possible and with confidence intervals noted  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A clear description of statistics including <u>central tendency</u> (e.g. median, mean) and <u>variation</u> (e.g. standard deviation, interquartile range)  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> Clearly defined error bars   |

See the web collection on [statistics for biologists](#) for further resources and guidance.

## ► Software

Policy information about [availability of computer code](#)

### 7. Software

Describe the software used to analyze the data in this study.

See Supplementary Text.

For manuscripts utilizing custom algorithms or software that are central to the paper but not yet described in the published literature, software must be made available to editors and reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). *Nature Methods* [guidance for providing algorithms and software for publication](#) provides further information on this topic.

## ► Materials and reagents

Policy information about [availability of materials](#)

### 8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

N/A

### 9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

N/A

### 10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

N/A

b. Describe the method of cell line authentication used.

N/A

c. Report whether the cell lines were tested for mycoplasma contamination.

N/A

d. If any of the cell lines used are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

N/A

## ► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

### 11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

Extended Data Figure 1; Supplementary Text 1.

Policy information about [studies involving human research participants](#)

### 12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

N/A