

1 Assessing connectivity despite high diversity in island  
2 populations of the malaria mosquito *Anopheles*  
3 *gambiae*

4 Christina M. Bergey<sup>1,2,3,\*</sup>, Martin Lukindu<sup>1,2</sup>, Rachel M. Wiltshire<sup>1,2</sup>,  
5 Michael C. Fontaine<sup>4,5</sup>, Jonathan K. Kayondo<sup>6</sup>, and Nora J. Besansky<sup>1,2,\*</sup>

6 <sup>1</sup>Department of Biological Sciences, University of Notre Dame, Notre Dame, IN 46556, USA

7 <sup>2</sup>Eck Institute for Global Health, University of Notre Dame, Notre Dame, IN 46556, USA

8 <sup>3</sup>Departments of Anthropology and Biology, Pennsylvania State University, University Park, PA  
9 16802, USA.

10 <sup>4</sup>Groningen Institute for Evolutionary Life Sciences (GELIFES), University of Groningen, PO  
11 Box 11103 CC, Groningen, The Netherlands.

12 <sup>5</sup>MIVEGEC, IRD, CNRS, University of Montpellier, Montpellier, France.

13 <sup>6</sup>Department of Entomology, Uganda Virus Research Institute (UVRI), Entebbe, Uganda.

14 \**To whom correspondence should be addressed.*

15 September 28, 2018

## 16 Abstract:

17 Modeling and cage experiments suggest that mosquito gene drive systems will en-  
18 able malaria eradication, but establishing safety and efficacy requires field-testing  
19 in isolated populations. Documenting genetic isolation is notoriously difficult for  
20 species with vast polymorphic populations like the principal African malaria vec-  
21 tor *Anopheles gambiae*. Using genome-wide variation, we assess Lake Victoria  
22 islands as candidate field-testing sites. One island, 30 kilometers offshore, is as  
23 differentiated from mainland samples as populations from across the continent,  
24 and we confirm isolation using adaptive variation as a powerful assay of connec-  
25 tivity. Collectively, our results suggest sufficient contemporary isolation of these  
26 islands to warrant consideration as field-testing locations.

## 27 Introduction

28 In efforts to reduce the approximately 445,000 annual deaths attributable to malaria [1],  
29 conventional vector control techniques may soon be augmented with releases of genetically  
30 modified mosquitoes. The most promising involve introducing transgenes into the mosquito  
31 genome or its endosymbionts that interrupt pathogen transmission coupled with a gene drive  
32 system to propagate the effector genes through a population [2–4], and such systems have  
33 recently been successfully engineered in the laboratory for the major African malaria vector  
34 *Anopheles gambiae sensu stricto* (henceforth *An. gambiae*) [5]. Effective implementation of  
35 any genetic control method, not least a gene drive system designed to spread in a super-  
36 Mendelian fashion, will benefit from a detailed understanding of population structure and  
37 connectivity.

38 Several life history characteristics of *An. gambiae* complicate efforts to estimate con-  
39 nectivity from genetic data, however. High fecundity and dispersal potential result in large

40 interconnected populations exhibiting shallow population structure [6, 7], manifested in the  
41 genome as high levels of polymorphism [6] that impede accurate estimation of connectivity [8]  
42 and discernment of demographic independence from panmixia [9]. Population genetic meth-  
43 ods for estimating migration between *An. gambiae* populations using neutral markers may  
44 have limited utility when such a high proportion of diversity is shared between populations,  
45 a failing that is only partially redressed with the high quantity of markers available from  
46 massively parallel sequencing. The most powerful window into migration may instead be the  
47 distribution of selected variants [10]: Adaptive introgression of beneficial haplotypes indi-  
48 cates migration occurred, while the absence of a selective sweep signature that is otherwise  
49 widespread would suggest barriers to gene flow.

50 Islands present natural laboratories for disentangling the determinants of population  
51 structure, as gene flow—likely important in post-dry season recolonization [11]—is reduced.  
52 In addition to evolutionary insight, investigations of island population structure have practi-  
53 cal rationales. Geographically-isolated islands have been proposed as initial field sites to test  
54 the dynamics of transgene spread while limiting their movement beyond the study popula-  
55 tion [12–15]. Antecedent studies of population structure and connectivity of potential release  
56 sites are necessary to evaluate the success of such field trials, as well as to quantify the chance  
57 of migration of transgenic insects carrying constructs designed to propagate across mosquito  
58 populations and country borders.

59 We analyzed genome-wide variation in *An. gambiae* mosquitoes living near and on the  
60 Ssese archipelago of Lake Victoria in Uganda (Fig. 1) to understand the determinants of  
61 their genetic variation, recent and long-term connectivity and demographic history, and the  
62 spread of adaptive variants across the region. In addition to the high malaria prevalence  
63 of the islands (44% in children; 30% in children country-wide; [16]), we were motivated by  
64 the potential of such an island to be a field site for future tests of gene-drive vector control  
65 strategies.

## 66 Results

67 The Ssesse Islands are approximately 4-50 km from the mainland, and vary in size, infras-  
68 tructure, and accessibility. Sampled islands range from Banda—a small, largely forested  
69 island of approximately 1 square kilometer with a single settlement—to Bugala—296 square  
70 kilometers, site of a 10,000 ha oil palm plantation [17], and linked to the mainland via ferry  
71 service [18]. To explore the partitioning of *An. gambiae* genetic variation in the Lake Victo-  
72 ria Basin (LVB), we sequenced the genomes of 116 mosquitoes from 5 island and 4 mainland  
73 localities (Fig. 1, Supplementary Table S1). We sequenced 10-23 individuals per site to an  
74 average depth of  $17.6 \pm 4.6$  (Supplementary Table S2). After filtering, we identified 28.6  
75 million high quality Single Nucleotide Polymorphisms (SNPs). We merged our dataset with  
76 that of the *An. gambiae* 1000 Genomes project (Ag1000G; [6]) for a combined dataset of  
77 12.54 million SNPs (9.86 million after linkage disequilibrium pruning) in 881 individuals.



Figure 1: Map of Lake Victoria Basin study area.

Map of study area showing sampling localities on Ssesse Islands (blue) and mainland localities (red) in Lake Victoria Basin. The Ag1000G reference population, Nagongera, Tororo District, is not shown, but lies 111 km NE of Kiyindi, 57 km from the shore of Lake Victoria. Map data copyright 2018 Google.

## 78 Genetic structure

79 We analyzed LVB population structure with context from continent-wide populations [6]  
80 of *An. gambiae* and sister species *Anopheles coluzzii* mosquitoes (formerly known as *An.*  
81 *gambiae* M molecular form [19]). Both Bayesian clustering ([20]; Fig. 2a) and principal  
82 component analysis (PCA; Fig. 2d) showed LVB individuals closely related to the Ugan-  
83 dan reference population (Nagongera, Tororo; 0°46'12.0"N, 34°01'34.0"E; ~57 km from Lake  
84 Victoria; Fig. 1). With  $\geq 6$  clusters (which optimized predictive accuracy in the clustering  
85 analysis; Supplementary Fig. S2), island samples had distinct ancestry proportions (Fig.  
86 2a), and with  $k = 9$  clusters, we observed additional subdivision in LVB samples and the as-  
87 signment of the majority of Ssesse individuals' ancestry to a largely island-specific component  
88 (Figs. 2a, 2b, 2c, and Supplementary Fig. S1).

89 PCA of only LVB individuals indicated little differentiation among mainland samples  
90 in the first two components and varying degrees of differentiation on islands, with Banda,  
91 Sserinya, and Bukasa the most extreme (Fig. 2e). Twelve of 23 individuals from Bugala, the  
92 largest, most developed, and most connected island, exhibited affinity to mainland individ-  
93 uals instead of ancestry typical of the islands (Supplementary Fig. S3). As both PCA and  
94 clustering analyses revealed this differentiation, we split the Bugala sample into mainland-  
95 and island-like subsets for subsequent analyses (hereafter referenced as "Bugala (M)" and  
96 "Bugala (I)," respectively). Individuals with partial ancestry attributable to the component  
97 prevalent on the mainland and the rest to the island-specific component were present on all  
98 islands except Banda.

99 Differentiation concurred with observed population structure. Mean  $F_{ST}$  between sam-  
100 pling localities (range: 0.001-0.034) was approximately 0 ( $\leq 0.003$ ) for mainland-mainland  
101 comparisons and was highest in comparisons involving small island Banda (Fig. 2f). Ge-  
102 ographic distances and  $F_{ST}$  were uncorrelated (Mantel  $p = 0.88$ ;  $R^2 = 0.08$ ,  $p = 0.048$ ;  
103 Supplementary Fig. S4). Island samples showed greater within- and between-locality shar-

104 ing of genomic regions identical by descent (IBD), with sharing between nearby islands  
105 Sserinya, Banda and Bugala (Fig. 2g). Importantly, Banda Island shared no IBD regions  
106 with mainland sites, underscoring its contemporary isolation from the mainland.

## 107 Genetic diversity

108 Consistent with the predicted decrease in genetic variation for semi-isolated island popula-  
109 tions due to inbreeding and smaller effective population sizes ( $N_e$ ), islands displayed lower  
110 nucleotide diversity ( $\pi$ ), slightly higher proportion of shared to rare variants (Tajima's  $D$ ),  
111 more variance in inbreeding coefficient ( $F$ ), more linkage among SNPs (LD;  $r^2$ ), longer runs  
112 of homozygosity ( $F_{ROH}$ ), and longer IBD tracks (Fig. 3). Small island Banda was the most  
113 extreme in these measures.

## 114 Demographic history

115 To test islands for isolation and demographic independence from the mainland, we inferred  
116 the population history of LVB samples by estimating long-term and recent trends in  $N_e$  using  
117 stairway plots [21] based on the site frequency spectrum (SFS; Fig. 4a) and patterns of IBD  
118 sharing ([22]; Fig. 4b), respectively. Short-term final mainland sizes were unrealistically  
119 high, likely due to per-locality sample sizes, but island-mainland differences were nonetheless  
120 informative. In both, islands had consistently lower  $N_e$  compared to mainland populations  
121 extending back 500 generations ( $\sim 50$  years) and often severely fluctuated, particularly in  
122 the last 250 generations ( $\sim 22$  years). Mainland sites Wamala and Kaazi had island-like  
123 recent histories, with Wamala abruptly switching to an island-like pattern.

124 To all pairs of LVB localities we fit an isolation-with-migration (IM) demographic model  
125 using  $\delta a \delta i$ , in which an ancestral population splits into two populations, allowing exponential  
126 growth and continuous asymmetrical migration between the daughter populations (Supple-



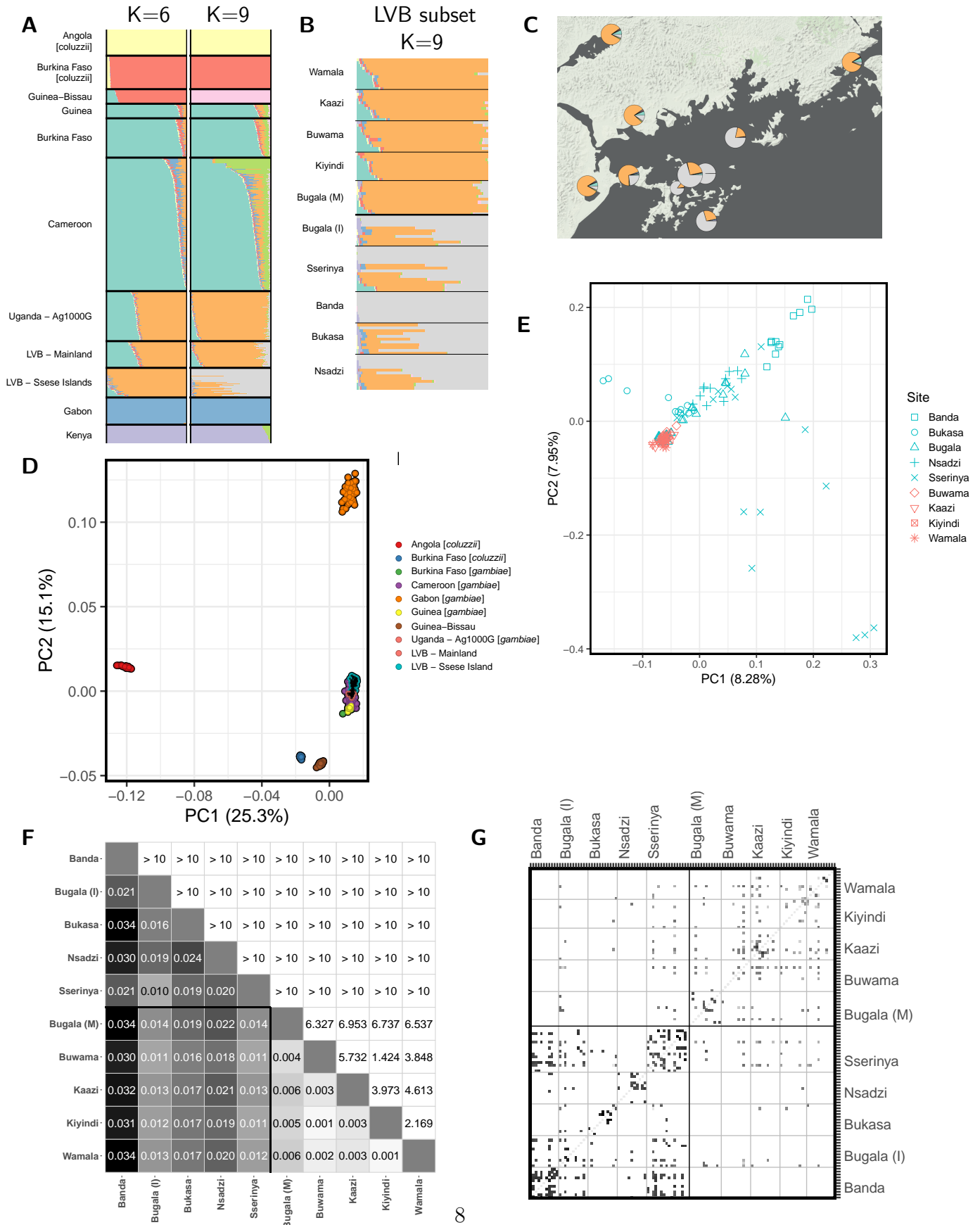


Figure 2: (Caption on next page.)



Figure 2: Population structure in the Lake Victoria Basin.

Analyses are based on chromosome 3 to avoid segregating inversions on other chromosome, unless otherwise noted. (A) ADMIXTURE-inferred ancestry of individuals in Lake Victoria Basin. Results based on analysis of LVB and Ag1000G merged dataset. Analysis is restricted to *A. gambiae s. s.*. Clustering shown for  $k = 6$  clusters, which minimizes cross validation error, and  $k = 9$  clusters, at which island individuals have the majority of their ancestry assigned to an island-specific cluster. (B) Results of the clustering analysis with  $k = 9$  clusters for LVB individuals, split by sampling locality. (C) Ancestry of individuals in Lake Victoria Basin and of Ag1000G reference populations as inferred by clustering into  $k = 9$  clusters. Samples are *A. gambiae* unless noted. (D) PCA plot of study individuals and *A. gambiae* and *A. coluzzii* individuals from reference Ag1000G populations. (E) Plot of first two components of PCA of Lake Victoria Basin individuals showing locality of origin. Mainland individuals are colored red, while island individuals are blue, and point shape indicates sampling locality. Based on these results, the island sample of Bugala was split into mainland- and island-like subpopulations (“Bugala (M)” and “Bugala (I),” respectively) for subsequent analyses (Fig. S3). (F) Heatmap of  $F_{ST}$  between sites (lower triangle) and associated  $z$ -score (upper triangle). “Bugala (M)” and “Bugala (I)” are the mainland- and island-like subpopulations of Bugala. (G) Genome-wide pairwise IBD proportions between individuals, based on the full genome, plotted on a logarithmic scale.

127 mentary Fig. S5). In all comparisons involving islands and some between mainland sites, the  
128 best fitting model as chosen via AIC had zero migration (Supplementary Tables S4, S5, and  
129 S6). Time since population split was much more recent for mainland-mainland comparisons  
130 (excluding Bugala, median: 511 years) than those involving islands (island-island median:  
131 9,080 years; island-mainland median: 5,450 years). Island-island split time confidence inter-  
132 vals typically did not overlap those involving mainland sites.

## 133 Selection

134 We next investigated patterns of selection using genome scans of between- and within-locality  
135 statistics (Supplementary Figs. S8, S7, Supplemental Text), including  $F_{ST}$  [23], Extended  
136 Haplotype Homozygosity (XP-EHH, [24]), and haplotype homozygosity (H12, [25]). Outlier  
137 regions included known selective sweep targets [6], including insecticide resistance-associated  
138 cytochrome P450 *Cyp6P2* which exhibited low diversity ( $\pi$ ), an excess of low frequency

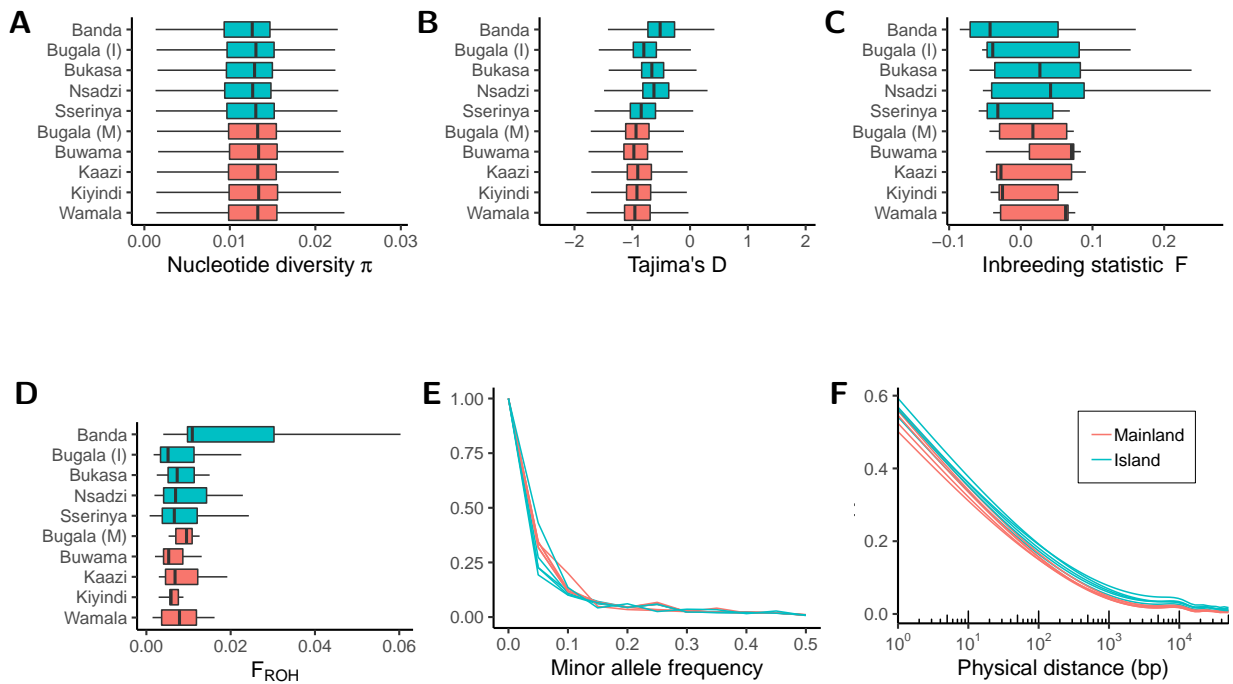


Figure 3: Diversity metrics in the Lake Victoria Basin samples.

Shown are a (A) boxplot of nucleotide diversity ( $\pi$ ; in 10 kilobase windows), (B) boxplot of Tajima's  $D$  (in 10 kilobase windows), (C) boxplot of inbreeding statistic ( $F$ ), (D) boxplot of length of runs of homozygosity ( $F_{ROH}$ ), (E) histogram of Minor Allele Frequency (MAF), and (F) decay in linkage disequilibrium ( $r^2$ ), all grouped by sampling locality. For all boxplots, outlier points are not shown.

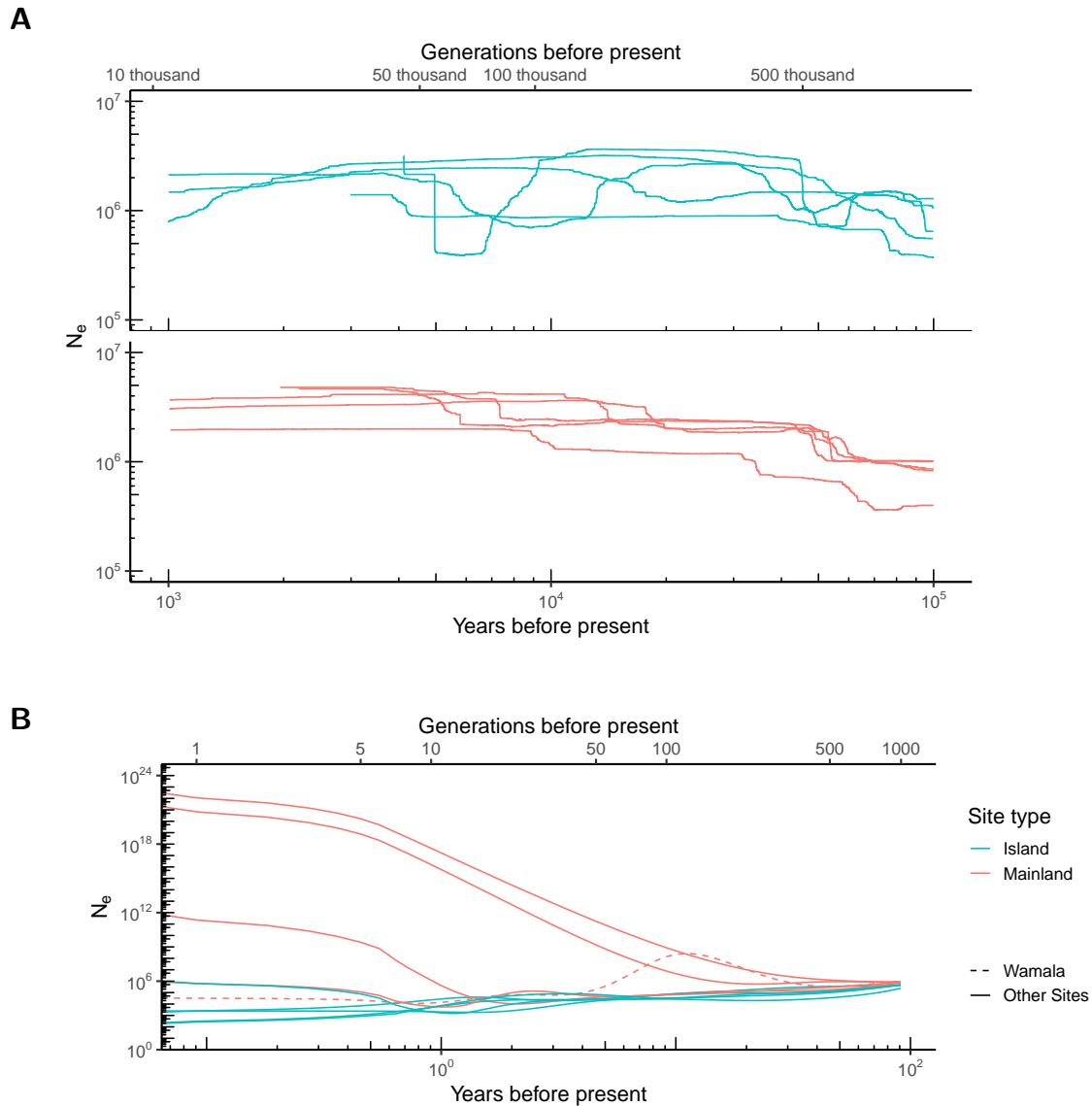


Figure 4: Population history of the Lake Victoria Basin samples.

(A) Long-term evolutionary population histories inferred via stairway plots for island and mainland samples. (B) Contemporary or short-term effective population size ( $N_e$ ) history inferred using sharing of regions that are identical by descent (IBD). Wamala, a mainland locality showing island-like fluctuations in population size, is indicated with a dashed line, and Kaazi shows the most consistently low population size for any mainland site.

139 polymorphisms (Tajima's  $D$ ), and elevated haplotype homozygosity (H12), but low differ-  
140 entiation between LVB localities, as expected for a continent-wide sweep (Supplementary  
141 Fig. S11). Hierarchical clustering of LVB and Ag1000G haplotypes revealed clades with  
142 low inter-individual diversity, expected after selection rapidly increases the frequency of a  
143 haplotype containing adaptive variation (Supplementary Fig. S17).

144 Widespread selective sweeps that are absent or at extremely low frequency on the islands  
145 are strong evidence against contemporary connectivity. To test for such sweeps, we identified  
146 locality-, island-, or LVB mainland-specific sweeps ( $H12 > 99$ th percentile), and intersected  
147 these regions with those under putative selection across the continent ( $H12 > 95$ th percentile  
148 in Ag1000G; [6]). Locality-specific putative sweeps were more prevalent on island than LVB  
149 mainland localities (mean per locality: island = 52.4; mainland = 26.8), concordant with  
150 increased isolation (Supplementary Table S7). Some sweeps targeting insecticide genes with  
151 continent-wide prevalence ([6]; Supplementary Table S10) were found to have colonized the  
152 islands incompletely (*Cyp9K1*: Supplementary Figs. S16 and S12). For instance, the sweep  
153 on the cluster of genes encoding glutathione S-transferases (*Gste1-Gste7*) was present across  
154 the continent but confined largely to the mainland in the LVB (Supplementary Figs. S18,  
155 and S13).

156 Besides known insecticide-related loci, we identified two regions of elevated between-  
157 locality differentiation, low diversity, and extended homozygosity (Supplementary Figs. S8,  
158 S7, S10, and 5). The first, at 2L:34.1 Mb, contains many genes, including a cluster involved in  
159 chorion formation [26] near the signal peak. Haplotype clustering revealed a group of closely-  
160 related Ugandan individuals, consistent with a geographically bounded selective sweep (Sup-  
161 plementary Fig. S14), but the selected variation had not fully colonized the islands. Similar  
162 low-variation clades in distinct genetic backgrounds were also found in, *e.g.*, Cameroon and  
163 Angola, suggesting convergent selection. The second, at X:9.2 Mb, coincided precisely with  
164 eye-specific diacylglycerol kinase (AGAP000519, chrX:9,215,505-9,266,532). Suggestive of

165 a single sweep, low diversity haplotypes formed a single cluster including LVB haplotypes  
166 overwhelmingly from the islands and surprisingly most closely related to haplotypes from  
167 distant locations, primarily Gabon and Burkina Faso rather than Uganda (Supplementary  
168 Fig. S15). Other adaptive variation supports this surprising affinity between the islands  
169 and West Africa: While the LVB mainland-specific sweeps (Supplementary Table S9) were  
170 co-located more often with those of the nearby reference Ugandan population (24%) than  
171 those of Gabon (16%), far more island-specific sweeps (Supplementary Table S8) were also  
172 putative sweeps in the Gabon population (33%) than in the Ugandan population (4.8%).

## 173 Discussion

174 Understanding the population genetics of island *Anopheles gambiae* has both evolutionary  
175 and practical importance. A limited number of genetic investigations have been conducted  
176 on oceanic [27–30] and lacustrine islands [31–34], though the latter have been limited in the  
177 type or count of molecular markers used. In contrast to shallow population structure across  
178 Africa [6, 7], partitioning of genetic variation on islands suggests varying isolation. Using  
179 a genome-wide dataset, we found differentiation between the Ssesse Islands to be relatively  
180 high in the context of continent-wide structure, with the differentiation between Banda Island  
181 (only 30 km offshore) and mainland localities on par with or higher than for populations on  
182 opposite sides of the continent or from different species (*e.g.*, Banda vs. Wamala,  $F_{ST} =$   
183 0.034; mainland Uganda vs. Burkina Faso,  $F_{ST} = 0.007$  [6]; *An. gambiae* vs. *An. coluzzii* in  
184 Burkina Faso:  $F_{ST} = 0.031$  [6]). The Ssesse Islands are approximately as differentiated as all  
185 but the most outlying oceanic islands tested (*e.g.* mainland Tanzania vs. Comoros, 690-830  
186 km apart,  $F_{ST} = 0.199-0.250$  [29]). Patterns of haplotype sharing did include direct evidence  
187 for the recent exchange of migrants between nearby islands, but analyses based on haplotype  
188 sharing, Bayesian clustering, and demographic reconstruction included no evidence of sharing

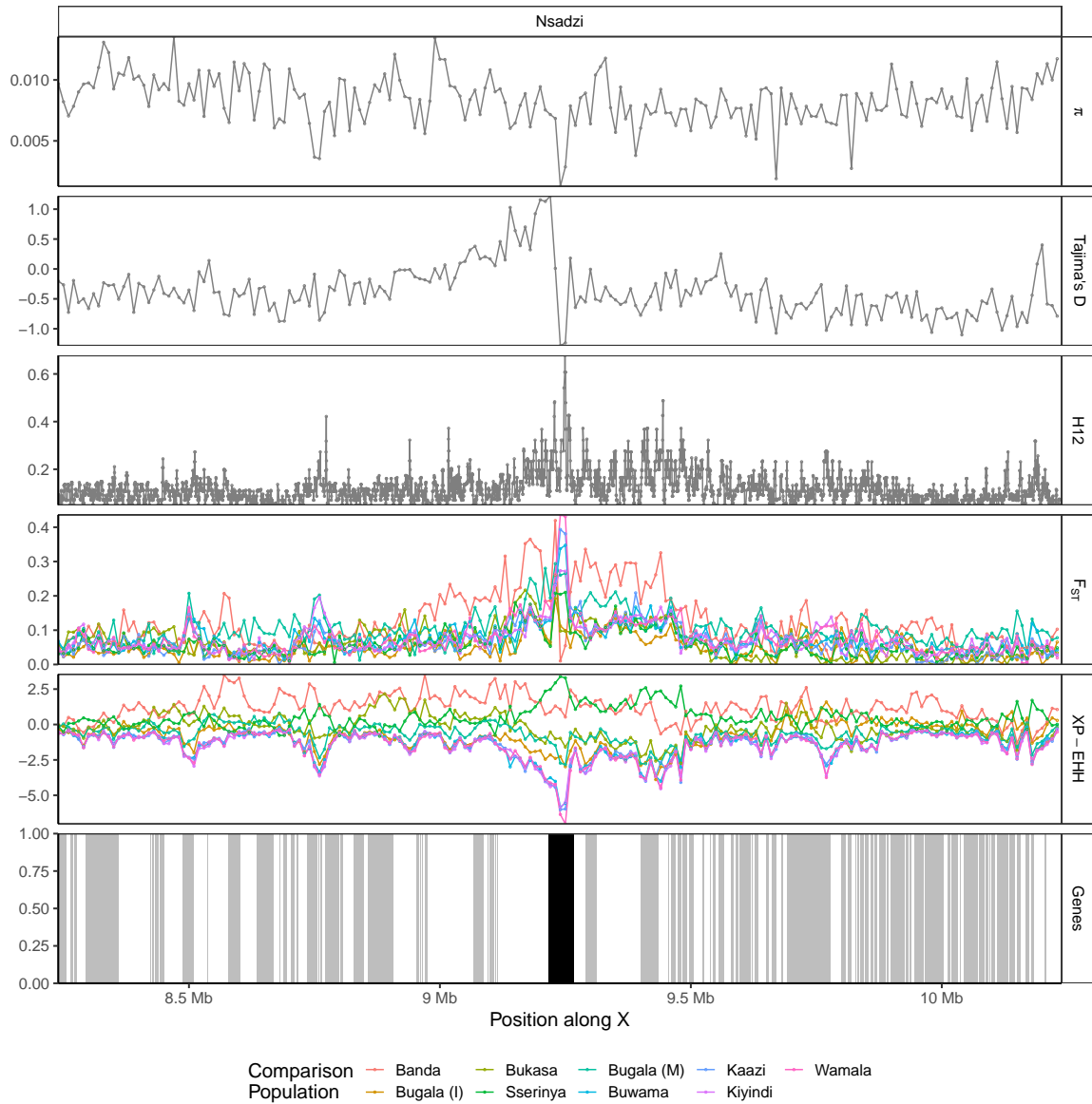


Figure 5: (Caption on next page.)

Figure 5: Selective sweep signal on X-chromosome near *rdgA* ortholog.

Population genetic statistics plotted near putative sweep on X-chromosome. Focus population for all pairwise  $F_{ST}$  and XP-EHH comparisons is island site Nsadzi. Region shown is 1 Mb up- and downstream of sweep target, centered at chrX:9,238,942. The gene eye-specific diacylglycerol kinase (AGAP000519, chrX:9,215,505-9,266,532) is highlighted in black.

189 between Banda and the mainland.

190 The name “Ssese” derives from another arthropod vector, the tsetse fly (*Glossina* spp.)  
191 The tsetse-mediated arrival of sleeping sickness in 1902 brought “enormous mortality” [35,  
192 pp. 332] to the 20 thousand residents, who were evacuated in 1909 [35, 36]. Though en-  
193 couraged to return by 1920, the human population numbered only 4 thousand in 1941 [35]  
194 and took until 1980 to double [37], but has since rapidly risen to over 62 thousand (2015,  
195 projected; [18, 38]). The impacts on mosquito populations of this prolonged depression in  
196 human population size, coupled with water barriers to mosquito migration, are reflected in  
197 the distinctive demographic histories of island *An. gambiae* populations, which were smaller  
198 and fluctuated more than mainland localities, echoing previous results [32, 34]. Two main-  
199 land sites had island-like recent population histories, with Wamala abruptly switching from  
200 a mainland-like to island-like growth pattern around 2005. This coincides precisely with a  
201  $\geq 20\%$  reduction from 2000-2010 in the *Plasmodium falciparum* parasite rate (PfPR<sub>2-10</sub>; a  
202 measure of malaria transmission intensity) in Mityana, the district containing Wamala [39].

203 Though previous *Anopheles* population genetic studies have inferred gene flow even  
204 among species [6, 40], we inferred that no genetic exchange had occurred since the split  
205 between island sites and between islands and the mainland. Island pairs were inferred to  
206 have split far deeper in the past (5,000-14,000 years ago) than mainland sites (typically  
207  $< 500$  years ago), on par with the inferred split time between Uganda and Kenya (approx-  
208 imately 4,000 years ago; [6]). Although bootstrapping-derived confidence intervals permit  
209 some certainty, our model fit is not optimal and additional sampling is necessary to clarify  
210 population history. Our inferred lack of gene flow to the islands appears contradictory to



211 the presence of individuals who share ancestry with the mainland on all islands but Banda.  
212 We cannot dismiss the possibility that this indicates actual migration occurs. If so, effects  
213 of migration would have to be sufficiently countered by local selection to limit its effect on  
214 allele frequency spectra, rendering effective migration (as estimated in population history  
215 inference) zero. The apparent contradiction can also be resolved if shared ancestry between  
216 islands and mainland suggested by the clustering result is interpreted as retention of shared  
217 ancestral polymorphism or the existence of inadequately sampled ancestral variation [41],  
218 rather than recent admixture. This interpretation is consistent with the affinity we observed  
219 between the Ssesse Islands and West Africa in the structure of adaptive variation.

220 As insecticide treated bed net usage is present on the islands [18], variation conferring a  
221 major selective advantage would be expected to spread to and persist on the islands if migra-  
222 tion allows the transfer, and the strongest evidence of a lack of contemporary connectivity  
223 is therefore the absence of a sweep on the islands that is widespread on the continent. We  
224 found two sweeps on insecticide-related genes that are common targets of selection elsewhere  
225 but which have incompletely colonized the Ssesse Islands: one on cytochrome P450 monooxy-  
226 genase *Cyp9K1* [42, 43] present on some islands, and another on glutathione S-transferase  
227 genes (*Gste1-Gste7*; [44–47]) at extremely low frequency on the islands. That the selective  
228 sweeps targeting these loci [6] have not fully colonized the islands despite the advantage in  
229 detoxifying pyrethroids and DDT suggests a lack of contemporary exchange.

230 Our investigation also identified two previously unknown signatures of selection with  
231 similar uneven distributions. The first encompassed many genes, including a cluster involved  
232 in egg shell formation, and the confinement of the signal to Ugandan mosquitoes and limited  
233 distribution on the islands suggests a local origin and spread via short distance migration.  
234 Overlapping signals in distinct backgrounds suggest the region has been affected by multiple  
235 independent convergent sweeps. The putative target of the second sweep is diacylglycerol  
236 kinase on the X-chromosome, a homolog of retinal degeneration A (*rdgA*) in *Drosophila*. The

237 gene is highly pleiotropic, contributing to signal transduction in the fly visual system [48, 49],  
238 but also olfactory [50] and auditory [51] sensory processing. It has been recently implicated  
239 in nutritional homeostasis in *Drosophila* [52] and is known to interact with the TOR pathway  
240 [53], which has been identified as a target of ecological adaptation in *Drosophila* [54, 55] and  
241 *An. gambiae* [56]. The sweep appears largely confined to island individuals in the LVB,  
242 but their most closely related haplotypes are primarily from Gabon, Burkina Faso, and  
243 Kenya. Shared extended haplotypes suggest a single sweep event, not convergence. Possible  
244 explanations include long distance migration of an adaptive variant persisting on only the  
245 islands, possibly due to a local selective advantage resisting the introgression of mainland  
246 haplotypes. We have not found obvious candidate targets of selection, *e.g.* coding changes,  
247 which may be due to imperfect annotation of the genome or the likely possibility that the  
248 target is a non-coding regulator of transcription or was filtered from our dataset. Further  
249 functional studies would be needed to clarify the selective advantage that these haplotypes  
250 confer. Interestingly, the putative sweep coincides with a similar region of low diversity in  
251 a cryptic subgroup of *Anopheles gambiae sensu lato* (GOUNDRY; [40]), suggesting possible  
252 convergence.

253 Population structure investigations are paramount for informing the design and deploy-  
254 ment of control strategies, including field trials of transgenic mosquitoes. We demonstrate  
255 alternatives to simple extrapolation of migration rates from differentiation, which is fraught  
256 [57] particularly given the assumption of equilibrium between the evolutionary forces of  
257 migration and drift [57–59], an unlikely state for huge *An. gambiae* populations [10]. We  
258 suggest that future assessments of connectivity include, as we have, the spatial distribution of  
259 adaptive variation, identification of recent migrants via haplotype sharing, and demographic  
260 history modeling, from which we have inferred the Ssesse Islands to be relatively isolated on  
261 contemporary time scales.

262 Though no island, lacustrine or oceanic, is completely isolated, the probability of contem-

263 porary migration may be sufficiently low to qualify some Ssesse Islands as candidate field sites.  
264 Furthermore, the assessment of the islands' suitability as potential sites for field trials of ge-  
265 netically modified mosquitoes must also consider the logistical ease of access and monitoring  
266 that the bounded geography of a small lacustrine island with low human population density  
267 affords initial field trials. Due consideration should be provided to these characteristics of  
268 small lake islands that may be appealing to regulators, field scientists, local communities,  
269 and other stakeholders. Given such features and the probable rarity of migration, the Ssesse  
270 Islands may be logical and tractable candidates for initial field tests of genetically modified  
271 *An. gambiae* mosquitoes, warranting further study.

## 272 **Acknowledgments**

273 The authors would like to thank the UVRI field entomology team: Christine Babirye, Ronald  
274 Mayanja, Paul Mawejje, Kevin Nakato, and Fred Ssenfuka. This study was supported by  
275 Target Malaria, which receives core funding from the Bill & Melinda Gates Foundation and  
276 from the Open Philanthropy Project Fund, an advised fund of Silicon Valley Community  
277 Foundation, through subcontracts to J.K.K. and N.J.B. N.J.B. also received support from  
278 NIH R01 AI125360 and R21 AI123491. The NYU School of Medicine's Genome Technology  
279 Center is partially supported by the Cancer Center Support Grant P30CA016087 at the  
280 Laura and Isaac Perlmutter Cancer Center. We thank Nicholas Harding and Alistair Miles  
281 for helpful discussion.

## 282 **Author Contributions**

283 C.M.B., J.K.K., and N.J.B. designed the study; C.M.B., M.L., R.M.W., and J.K.K. col-  
284 lected biological samples; C.M.B. analyzed the data; C.M.B., M.C.F., and N.J.B. wrote the

285 manuscript; M.C.F., J.K.K., and N.J.B. supervised the research; C.M.B., M.L., R.M.W.,  
286 M.C.F., J.K.K., and N.J.B. edited the manuscript.

## 287 Conflict of Interest Statement

288 The authors declare no competing financial interests.

## 289 References

- 290 [1] World Health Organization. *World Malaria Report 2017* (Geneva, 2017).
- 291 [2] Burt, A. Heritable strategies for controlling insect vectors of disease. *Philosophical*  
292 *transactions of the Royal Society of London. Series B, Biological sciences* **369**, 20130432  
293 (2014).
- 294 [3] Champer, J., Buchman, A. & Akbari, O. S. Cheating evolution: engineering gene drives  
295 to manipulate the fate of wild populations. *Nature Reviews Genetics* **17**, 146–159 (2016).
- 296 [4] Alphey, L. Genetic control of mosquitoes. *Annual Review of Entomology* **59**, 205–224  
297 (2014).
- 298 [5] Hammond, A. *et al.* A CRISPR-Cas9 gene drive system targeting female reproduction  
299 in the malaria mosquito vector *Anopheles gambiae*. *Nature Biotechnology* **34**, 78–83  
300 (2016).
- 301 [6] Miles, A. *et al.* Genetic diversity of the African malaria vector *Anopheles gambiae*.  
302 *Nature* **552**, 96–100 (2017).
- 303 [7] Lehmann, T. *et al.* Population structure of *Anopheles gambiae* in Africa. *Journal of*  
304 *Heredity* **94**, 133–147 (2003).

- 305 [8] Waples, R. S. Separating the wheat from the chaff: Patterns of genetic differentiation  
306 in high gene flow species. *Journal of Heredity* **89**, 438–450 (1998).
- 307 [9] Waples, R. S. A bias correction for estimates of effective population size based on linkage  
308 disequilibrium at unlinked gene loci. *Conservation Genetics* **7**, 167 (2006).
- 309 [10] Gagnaire, P. A. *et al.* Using neutral, selected, and hitchhiker loci to assess connectivity  
310 of marine populations in the genomic era. *Evolutionary Applications* **8**, 769–786 (2015).
- 311 [11] Dao, A. *et al.* Signatures of aestivation and migration in Sahelian malaria mosquito  
312 populations. *Nature* **516**, 387–390 (2014).
- 313 [12] World Health Organization. *Guidance framework for testing genetically modified*  
314 *mosquitoes* (2014).
- 315 [13] Alphey, L. Malaria control with genetically manipulated insect vectors. *Science* **298**,  
316 119–121 (2002).
- 317 [14] James, A. A. Gene drive systems in mosquitoes: rules of the road. *Trends in Parasitology*  
318 **21**, 64–67 (2005).
- 319 [15] James, S. *et al.* Pathway to deployment of gene drive mosquitoes as a potential bio-  
320 control tool for elimination of malaria in sub-Saharan Africa: Recommendations of a  
321 scientific working group. *American Journal of Tropical Medicine and Hygiene* **98**, 1–49  
322 (2018).
- 323 [16] Uganda Bureau of Statistics (UBOS) & ICF. *Uganda Demographic and Health Survey*  
324 *2016: Key Indicators Report* (Kampala, Uganda: UBOS, and Rockville, Maryland,  
325 USA: UBOS and ICF, 2017).
- 326 [17] Zeemeijer, I. *Who Gets What, When and How?: New Corporate Land Acquisitions and*  
327 *the Impact on Local Livelihoods in Uganda*. Master’s thesis, Utrecht University (2012).

- 328 [18] Kalangala District Local Government District Management Improvement Plan 2012-  
329 2015. Tech. Rep. (2012).
- 330 [19] Coetzee, M. *et al.* *Anopheles coluzzii* and *Anopheles amharicus*, new members of the  
331 *Anopheles gambiae* complex. *Zootaxa* **3619**, 246–274 (2013).
- 332 [20] Alexander, D., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in  
333 unrelated individuals. *Genome Research* **19**, 1655–1664 (2009).
- 334 [21] Liu, X. & Fu, Y.-X. Exploring population size changes using SNP frequency spectra.  
335 *Nature Genetics* **47**, 555–559 (2015).
- 336 [22] Browning, S. R. & Browning, B. L. Accurate non-parametric estimation of recent  
337 effective population size from segments of identity by descent. *American Journal of*  
338 *Human Genetics* **97**, 404–418 (2015).
- 339 [23] Weir, B. & Cockerham, C. Estimating F-statistics for the analysis of population struc-  
340 ture. *Evolution* **38**, 1358–1370 (1984).
- 341 [24] Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive selection in  
342 human populations. *Nature* **449**, 913–918 (2007).
- 343 [25] Garud, N. R., Messer, P. W., Buzbas, E. O. & Petrov, D. A. Recent selective sweeps in  
344 North American *Drosophila melanogaster* show signatures of soft sweeps. *PLoS Genetics*  
345 **11**, 1–32 (2015).
- 346 [26] Ameny, D. A. *et al.* Proteomics reveals novel components of the *Anopheles gambiae*  
347 eggshell. *Journal of Insect Physiology* **56**, 1414–1419 (2010).
- 348 [27] Moreno, M. *et al.* Genetic population structure of *Anopheles gambiae* in Equatorial  
349 Guinea. *Malaria Journal* **6**, 137 (2007).

- 350 [28] Marshall, J. C. *et al.* Exploring the origin and degree of genetic isolation of *Anopheles*  
351 *gambiae* from the islands of São Tomé and Príncipe, potential sites for testing transgenic-  
352 based vector control. *Evolutionary Applications* **1**, 631–644 (2008).
- 353 [29] Marsden, C. D. *et al.* An analysis of two island groups as potential sites for trials of  
354 transgenic mosquitoes for malaria control. *Evolutionary Applications* **6**, 706–720 (2013).
- 355 [30] Maliti, D. *et al.* Islands and stepping-stones: Comparative population structure of  
356 *Anopheles gambiae sensu stricto* and *Anopheles arabiensis* in Tanzania and implications  
357 for the spread of insecticide resistance. *PLoS ONE* **9**, e110910 (2014).
- 358 [31] Chen, H., Minakawa, N., Beier, J. & Yan, G. Population genetic structure of *Anopheles*  
359 *gambiae* mosquitoes on Lake Victoria islands, west Kenya. *Malaria Journal* **3**, 48 (2004).
- 360 [32] Kayondo, J. K. *et al.* Genetic structure of *Anopheles gambiae* populations on islands  
361 in northwestern Lake Victoria, Uganda. *Malaria Journal* **4**, 59 (2005).
- 362 [33] Lukindu, M. *et al.* Spatio-temporal genetic structure of *Anopheles gambiae* in the  
363 Northwestern Lake Victoria Basin, Uganda: Implications for genetic control trials in  
364 malaria endemic regions. *Parasites and Vectors* **11** (2018).
- 365 [34] Wiltshire, R. M. *et al.* Reduced-representation sequencing identifies small effective pop-  
366 ulation sizes of *Anopheles gambiae* in the north-western Lake Victoria basin, Uganda.  
367 *Malaria Journal* **17**, 285 (2018).
- 368 [35] Thomas, A. The vegetation of the Sese Islands, Uganda: An illustration of edaphic  
369 factors in tropical ecology. *Journal of Ecology* **29**, 330–353 (1941).
- 370 [36] Hale Carpenter, G. D. *A Naturalist on Lake Victoria* (E. P. Dutton and Company, New  
371 York, NY, 1920).



- 372 [37] Uganda Bureau of Statistics. *2002 Uganda Population and Housing Census Analytical*  
373 *Report* (2002).
- 374 [38] Uganda Bureau of Statistics. *The National Population and Housing Census 2014 - Main*  
375 *Report* (2016).
- 376 [39] National Malaria Control Programme, Abt Associates & the INFORM Project. *An*  
377 *epidemiological profile of malaria and its control in Uganda* (2013).
- 378 [40] Crawford, J. E. *et al.* Evolution of GOUNDRY, a cryptic subgroup of *Anopheles gambiae*  
379 *s.l.*, and its impact on susceptibility to *Plasmodium* infection. *Molecular Ecology* **25**  
380 (2016).
- 381 [41] Lawson, D. J., van Dorp, L. & Falush, D. A tutorial on how not to over-interpret  
382 STRUCTURE and ADMIXTURE bar plots. *Nature Communications* **9**, 1–11 (2018).
- 383 [42] Vontas, J. *et al.* Rapid selection of a pyrethroid metabolic enzyme CYP9K1 by op-  
384 erational malaria control activities. *Proceedings of the National Academy of Sciences*  
385 201719663 (2018).
- 386 [43] Fossog Tene, B. *et al.* Resistance to DDT in an urban setting: Common mechanisms im-  
387 plicated in both M and S forms of *Anopheles gambiae* in the city of Yaoundé, Cameroon.  
388 *PLoS ONE* **8** (2013).
- 389 [44] Enayati, A. A., Ranson, H. & Hemingway, J. Insect glutathione transferases and insect-  
390 icide resistance. *Insect Molecular Biology* **14**, 3–8 (2005).
- 391 [45] Mitchell, S. N. *et al.* Metabolic and target-site mechanisms combine to confer strong  
392 DDT resistance in *Anopheles gambiae*. *PLoS ONE* **9** (2014).

- 393 [46] Jones, C. M. *et al.* Additional selection for insecticide resistance in urban malaria  
394 vectors: DDT resistance in *Anopheles arabiensis* from Bobo-Dioulasso, Burkina Faso.  
395 *PLoS ONE* **7** (2012).
- 396 [47] Fouet, C., Kamdem, C., Gamez, S. & White, B. J. Genomic insights into adaptive diver-  
397 gence and speciation among malaria vectors of the *Anopheles nili* group. *Evolutionary*  
398 *Applications* **10**, 897–906 (2017).
- 399 [48] Hardie, R. *et al.* Molecular basis of amplification in *Drosophila* phototransduction:  
400 Roles for G protein, phospholipase C, and diacylglycerol kinase. *Neuron* **36**, 689–701  
401 (2002).
- 402 [49] Huang, Y., Xie, J. & Wang, T. A fluorescence-based genetic screen to study retinal  
403 degeneration in *Drosophila*. *PLoS ONE* **10**, 1–19 (2015).
- 404 [50] Kain, P. *et al.* Reduced odor responses from antennal neurons of G(q)alpha, phos-  
405 pholipase Cbeta, and rdgA mutants in *Drosophila* support a role for a phospholipid  
406 intermediate in insect olfactory transduction. *The Journal of Neuroscience* **28**, 4745–  
407 4755 (2008).
- 408 [51] Senthilan, P. R. *et al.* *Drosophila* auditory organ genes and genetic hearing defects. *Cell*  
409 **150**, 1042–1054 (2012).
- 410 [52] Nelson, C. S. *et al.* Cross-phenotype association tests uncover genes mediating nutrient  
411 response in *Drosophila*. *BMC Genomics* **17**, 1–14 (2016).
- 412 [53] Lin, Y. H. *et al.* Diacylglycerol lipase regulates lifespan and oxidative stress response  
413 by inversely modulating TOR signaling in *Drosophila* and *C. elegans*. *Aging Cell* **13**,  
414 755–764 (2014).

- 415 [54] De Jong, G. & Bochdanovits, Z. Latitudinal clines in *Drosophila melanogaster*: Body  
416 size, allozyme frequencies, inversion frequencies, and the insulin-signalling pathway.  
417 *Journal of Genetics* **82**, 207–223 (2003).
- 418 [55] Fabian, D. K. *et al.* Genome-wide patterns of latitudinal differentiation among popula-  
419 tions of *Drosophila melanogaster* from North America. *Molecular Ecology* **21**, 4748–4769  
420 (2012).
- 421 [56] Cheng, C., Tan, J. C., Hahn, M. W. & Besansky, N. J. Systems genetic analysis of  
422 inversion polymorphisms in the malaria mosquito *Anopheles gambiae*. *Proceedings of*  
423 *the National Academy of Sciences* 201806760 (2018).
- 424 [57] Whitlock, M. C. & McCauley, D. E. Indirect measures of gene flow and migration:  $F_{ST}$   
425 not equal to  $1/(4Nm + 1)$ . *Heredity* **82** (Pt. 2), 117–125 (1999).
- 426 [58] Stow, A. J. & Magnusson, W. E. Genetically defining populations is of limited use for  
427 evaluating and managing human impacts on gene flow. *Wildlife Research* **39**, 290–294  
428 (2012).
- 429 [59] Storfer, A., Murphy, M. A., Spear, S. F., Holderegger, R. & Waits, L. P. Landscape  
430 genetics: Where are we now? *Molecular Ecology* **19**, 3496–3514 (2010).

## 431 **Methods**

432 **Sample collection, morphological ID** Mosquitoes were sampled from 5 of the Ssesse  
433 Islands in Lake Victoria, Uganda (Banda, Bukasa, Bugala, Nsadzi, and Sserinya) and 4  
434 mainland sampling localities (Buwama, Kaazi, Kiyindi, and Wamala) at varying distances  
435 from the lake in May and June, 2015. Sampling took place between 4:40 and 8:15 over a 30  
436 day period as follows: Indoor resting mosquitoes were collected from residences via mouth or  
437 mechanical aspirators and subsequently identified morphologically to species group. Female  
438 mosquitoes assigned to the *An. gambiae sensu lato* complex based on morphology ( $N=575$ )  
439 were included in further analyses. All mosquitoes were preserved with silica desiccant and  
440 transported to the University of Notre Dame, Indiana, U.S.A. for analysis.

441 **DNA extraction, Library preparation, and Whole Genome Sequencing** Animals  
442 were assigned to species level via a PCR-based assay [60] using DNA present in a single leg  
443 or wing. DNA from individual *An. gambiae s. s.*  $N=116$  mosquitoes was extracted from  
444 the whole body via phenol-chloroform extraction [61] and then quantified via fluorometry  
445 (PicoGreen). Automated library preparation took place at the NYU Langone Medical Center  
446 with the Biomek SPRIWorks HT system using KAPA Library Preparation Kits, and libraries  
447 were sequenced on the Illumina HiSeq 2500 with 100 paired end cycles.

448 **Mapping and SNP calling, filtering** Software version information is provided in Sup-  
449 plementary Table S11. After quality filtering and trimming using ea-utils' fastq-mcf (-l 15  
450 -q 15 -w 4; [62]), reads were mapped to the *An. gambiae* reference genome (AgamP4 PEST;  
451 [63, 64]) using BWA aln and sampe with default parameters [65].

452 After realignment around indels with GATK's IndelRealigner, variants were called using  
453 GATK's UnifiedGenotyper (with -stand\_call\_conf 50.0 and -stand\_emit\_conf 10.0) and filtered  
454 for quality [66], excluding SNPs with QualByDepth (QD) < 2.0, RMSMappingQuality (MQ)

455 < 40.0, FisherStrand (FS) > 60.0, HaplotypeScore > 13.0, or ReadPosRankSum < -8.0. All  
456 bioinformatic steps for read mapping and variant identification are encapsulated in the NGS-  
457 map pipeline (<https://github.com/bergeycm/NGS-map>). This yielded 33.1 million SNPs.  
458 Individuals and variants with high levels of missingness (> 10%) and variants that were  
459 not biallelic or exhibited values of HWE that were likely due to sequencing error ( $p <$   
460 0.00001) were excluded from further analysis. For use in population structure inference, the  
461 SNP dataset was further pruned for linkage disequilibrium by sliding a window of 50 SNPs  
462 across the genome in 5 SNP increments and recursively removing random SNPs in pairs  
463 with  $r^2 > 0.5$  using PLINK [67, 68]. After filtration, the dataset contained 28,569,621 SNPs  
464 before LD pruning and 115 individuals. SNPs unpruned for linkage disequilibrium were  
465 phased with SHAPEIT2 [69] using an effective population size ( $N_e$ ) of 1,000,000 (consistent  
466 with previous demographic modeling [70]), default MCMC parameters (7 burn-in MCMC  
467 iterations, 8 pruning iterations, and 20 main iterations), conditioning states for haplotype  
468 estimation ( $K = 100$ ), and window size of 2 Mb.

469 **Population structure inference** To explore population structure in a larger, continent-  
470 wide context, we merged our LVB SNP set with a recently published dataset of *Anopheles*  
471 *gambiae* individuals (from the Ag1000G project) [70]. Prior to filtering, biallelic SNPs from  
472 the LVB and Ag1000G datasets were merged using bcftools [71]. We excluded any SNP with  
473 greater than 10% missingness in either dataset, any SNPs that did not pass the accessibility  
474 filter of the Ag1000G dataset, and SNPs with MAF < 1%. After this filtration, our merged  
475 SNP dataset contained 12,537,007 SNPs.

476 After pruning the merged dataset for LD (leaving 9,861,756 SNPs) and excluding labo-  
477 ratory crosses (leaving 881 individuals), we assigned individuals' genomes to ancestry com-  
478 ponents using ADMIXTURE [72]. We created 10 replicate samples of 100,000 SNPs from  
479 chromosome 3 (prior to LD-pruning), including only biallelic SNPs in euchromatic regions

480 with  $MAF > 1\%$ . These replicate datasets were pruned for LD by randomly selecting from  
481 pairs of SNPs with  $r^2 > 0.01$  in sliding windows of size 500 SNPs and with a stepsize of  
482 250 SNPs. For each replicate, we ran ADMIXTURE for 5 iterations in five-fold cross val-  
483 idation mode for values of  $k$  from 2 to 10. This resulted in 50 estimates for each value of  
484  $k$ . We assessed these results using the online version of CLUMPAK with default settings to  
485 ensure the stability of the resulting clustering [73]. CLUMPAK clusters the replicate runs'  
486 Q-matrices to produce a major cluster for each value of  $k$ , which we then visualized. The  
487 lowest cross-validation error was found for  $k = 6$  clusters, but we also display ancestry esti-  
488 mates with  $k = 9$  clusters to further explore patterns of structure with a level of subdivision  
489 at which the Ssesse Island individuals are assigned a unique ancestry component.

490 We visualized population structure via principal components analysis (PCA) with PLINK  
491 [67, 68], using the LVB-Ag1000G merged dataset (excluding the outlier Kenyan population;  
492 [70]) and 3,212,485 chromosome 3 SNPs (to avoid the well-known inversions on chromosome  
493 2 and the X-chromosome) outside of heterochromatic regions (such as centromeric regions;  
494 [64]; Supplementary Table S3). We next performed a PCA on the LVB dataset alone, pruning  
495 for LD and low-MAF ( $< 1\%$ ) SNPs on chromosome 3. Based on the results of this analyses,  
496 we split individuals from the large island of Bugala into two clusters for subsequent analyses:  
497 those that cluster with mainland individuals and those that cluster with individuals from  
498 the smaller islands.

499 We computed the pairwise fixation index ( $F_{ST}$ ) between locality samples for *An. gambiae*  
500 using the unbiased estimator of Hudson [74] as implemented in smartpca [75, 76]. To obtain  
501 overall values between sampling sites, per-SNP values were averaged across the genome  
502 excluding known inversions (*2La*, *2Rb*, and *2Rc*) and heterochromatic regions. We also  
503 computed  $z$ -scores via block jackknife, using 42 blocks of size 5 Mb. We tested for isolation  
504 by distance, or a correlation between genetic and geographic distances, with a Mantel test  
505 [77] as implemented in the R package ade4 [78], using these  $F_{ST}$  estimates and Euclidean

506 geographic distances between localities.

507 To estimate fine-scale structure and relatedness between individuals, we estimated the  
508 proportion of pairs of individuals genomes that are identical by descent (IBD) using PLINK  
509 [67, 68]. We excluded heterochromatic and inversion regions, and retained informative pairs  
510 of SNPs within 500 kb in the pairwise population concordance test.

511 **Diversity estimation** Grouping individuals by site (except for Bugala, which was split  
512 based on the results of the PCA), we calculated nucleotide diversity ( $\pi$ ) and Tajima's  $D$   
513 in nonoverlapping windows of size 10 kb, the inbreeding coefficient ( $F$ ) estimated with the  
514 method of moments, minor allele frequencies (the site frequency spectrum, SFS), and a mea-  
515 sure of linkage disequilibrium ( $r^2$ ) using VCFtools (Danecek2011). For  $r^2$ , we computed  
516 the measure for all SNPs (unpruned for linkage) within 50 kb of a random set of 100 SNPs  
517 with MAF > 10% and corrected for differences in sample size by subtracting  $1/n$ , where  $n$   
518 equaled the number of sampled chromosomes per site. To visualize decay in LD, we plotted  
519  $r^2$  between SNPs against their physical distance in base pairs, first smoothing the data by  
520 fitting a generalized additive model (GAM) to them. We also inferred runs of homozygos-  
521 ity using PLINK [67, 68] to compare their length ( $F_{ROH}$ ), requiring 10 homozygous SNPs  
522 spanning a distance of 100 kb and allowing for 3 heterozygous and 5 missing SNPs in the  
523 window. Runs of homozygosity were inferred using LD-pruned SNPs outside of inversions  
524 or heterochromatic regions.

525 **Demographic history inference** To estimate the long-term evolutionary demographic  
526 history of mosquitoes on and near the Ssesse Islands, including a long-term estimate of  $N_e$   
527 [79], we inferred population demographic history for each site via stairway plots using the  
528 full site frequency spectra from the same dataset [80].

529 To estimate the contemporary or short-term  $N_e$  for each site, we inferred regions of IBD  
530 from unphased data with IBDseq [81] and analyzed them with IBDNe [82]. We restricted



531 our analysis to SNPs from chromosome 3 to avoid inverted regions. We allowed a minimum  
532 IBD tract length of 0.005 cM (or 5 kb), scaling it down from the recommended length for  
533 human genomes due to mosquitoes' high level of heterozygosity [70] and assumed a constant  
534 recombination rate of 2.0 cM/Mb [83].

535 We also inferred a “two-population” isolation-with-migration (IM) demographic model  
536 with  $\delta a \delta i$  [84, 85] in which the ancestral population splits to form two daughter populations  
537 that are allowed to grow exponentially and exchange migrants asymmetrically, as described  
538 in the main text. For  $\delta a \delta i$ -based analyses, we used the full dataset of SNPs on chromosome 3,  
539 not pruned for LD but with heterochromatic regions masked. We polarized the SNPs using  
540 outgroup information from *Anopheles merus* and *An. merus* [86]. We fit this two-population  
541 model and the same model without migration to all pairs of locality samples, choosing the  
542 optimal model using the Godambe Information Matrix and an adjusted likelihood ratio test  
543 to compare the two nested models. We compared the test statistic to a  $\chi^2$  distribution  
544 and rejected the null model if the p-value for the test statistic was  $> 0.05$ . For both,  
545 singletons and doubletons private to one population were masked from the analysis and a  
546 parameter encompassing genotype uncertainty was included in the models and found to be  
547 low (mean = 0.67%). We assessed the goodness-of-fit visually using the residuals of the  
548 comparison between model and data frequency spectra (Supplementary Fig. S6). Using the  
549 site frequency spectrum, we projected down to 2-6 fewer chromosomes than the total for the  
550 smaller population to maximize information given missing data. We set the grid points to  
551  $\{n, n + 10, n + 20\}$ , where  $n$  = the number of chromosomes. Bounds for  $N_e$  scalars were  
552  $\nu \in (0.01, 10, 000)$ , for time were  $T \in (1e-8, 0.1)$ , for migration were  $m \in (1e-8, 10)$ , and for  
553 genotyping uncertainty were  $p_{misid} \in (1e-8, 1)$ . Parameters were perturbed before allowing  
554 up to 1000 iterations for optimization. We estimated parameter uncertainty using the Fisher  
555 information matrix and 100 bootstrap replicates of 1 Mb from the dataset. If the Hessian  
556 was found to be not invertible when computing the Fisher information matrix, the results

557 of that iteration were excluded from the analysis.

558 To translate  $\delta a\delta i$ - and stairway plot-based estimates of  $N_e$  and time to individuals and  
559 years respectively, we assumed a generation time of 11 per year and a mutation rate of  $3.5e-9$   
560 per generation [70].

561 **Selection inference** To infer candidate genes and regions with selection histories that  
562 varied geographically, we compared allele frequencies and haplotype diversity between the  
563 sampling sites. To infer differing selection between sampling sites, we computed  $F_{ST}$  between  
564 all populations in windows of size 10 kb using the estimator of Weir and Cockerham [87]  
565 (as implemented in VCFtools [88]), and H12 (as implemented in SelectionHapStats [89])  
566 and XP-EHH on a per-site basis (as implemented in selscan [90]) to detect long stretches of  
567 homozygosity in a given population considered alone or relative to another population [91].  
568 For XP-EHH, EHH was calculated in windows of size 100 kb in each direction from core  
569 SNPs, allowing EHH decay curves to extend up to 1 Mb from the core, and SNPs with MAF  
570  $< 0.05$  were excluded from consideration as a core SNP. As we lacked a fine-scale genetic  
571 map for *Anopheles*, we assumed a constant recombination rate of 2.0 cM/Mb [83]. Scores  
572 were normalized within chromosomal arms and the X-chromosome. The between-locality  
573 statistics,  $F_{ST}$  and XP-EHH, were summarized using the composite selection score [CSS;  
574 [92, 93]].

575 We plotted these statistics across the genome to identify candidate regions with signa-  
576 tures of selection, including high differentiation between samples from different localities,  
577 reduced variability within a sample, and extended haplotype homozygosity. To identify re-  
578 gions of the genome showing signatures of selection specific to certain geographic areas, we  
579 identified genomic regions with elevated H12 in a subset of localities, and confirmed both ele-  
580 vated differentiation (as inferred from  $F_{ST}$ ) and evidence of differing selective sweep histories  
581 (as inferred from XP-EHH). Excluding the mainland-like portion of Bugala, we identified

582 putative locality-specific sweeps (H12 over 99<sup>th</sup> percentile in one population), island-specific  
583 sweeps (H12 over 99<sup>th</sup> percentile in 4 or more of the 5 island localities but 0 or 1 mainland  
584 localities), or LVB mainland-specific sweeps (H12 over 99<sup>th</sup> percentile in 3 or more of the 4  
585 island localities but 0 or 1 island localities). To place these putative sweeps in their continen-  
586 tal context, for the region of each putative locality-, island-, or LVB mainland-specific sweep,  
587 we determined if the H12 values of each of the Ag1000G populations (excluding Kenya due  
588 to its signatures of admixture and recent population decline; [70]) were in the top 5% for  
589 that population, indicating a possible selective sweep at the same location.

590 We further explored the haplotype structure and putative functional impact of loci for  
591 which we detected signatures of potential selection to determine the count and geographic  
592 distribution of independent selective sweeps. To provide necessary context for the recon-  
593 struction of sweeps and quantify long distance haplotype sharing between populations, we  
594 included data from several other *An. gambiae* populations across Africa (Burkina Faso,  
595 Cameroon, Gabon, Guinea, Guinea-Bissau, Kenya, and other Ugandan individuals; [70]).  
596 We computed the pairwise distance matrix as the raw number of base pairs that differed  
597 and grouped haplotypes via hierarchical clustering analysis (implemented in the hclust R  
598 function) in regions of size 100 kb centered on each peak or the average of peaks, in the case  
599 for multiple nearby spikes. As short terminal branches can result from a beneficial allele and  
600 linked variants rising to fixation during a recent selective sweep, we identified such clusters  
601 by cutting the tree at a height of 0.4 SNP differences per kb.

602 **Script and data availability** All scripts used in the analysis are available at [https://](https://github.com/bergeycm/Anopheles_gambiae_structure_LVB)  
603 [github.com/bergeycm/Anopheles\\_gambiae\\_structure\\_LVB](https://github.com/bergeycm/Anopheles_gambiae_structure_LVB) and released under the GNU  
604 General Public License v3. Sequencing read data for the LVB individuals are deposited in  
605 the NCBI Short Read Archive (SRA) under BioProject accession PRJNA493853.

## References

- [60] Scott, J. A., Brogdon, W. G. & Collins, F. H. Identification of single specimens of the *Anopheles gambiae* complex by the polymerase chain reaction. *The American Journal of Tropical Medicine and Hygiene* **49**, 520–529 (1993).
- [61] Green, M. R. & Sambrook, J. *Molecular Cloning: A Laboratory Manual* (2012).
- [62] Aronesty, E. ea-utils: Command-line tools for processing biological sequencing data (2011).
- [63] Holt, R. *et al.* The genome sequence of the malaria mosquito *Anopheles gambiae*. *Science* **298**, 129–149 (2002).
- [64] Sharakhova, M. V. *et al.* Update of the *Anopheles gambiae* PEST genome assembly. *Genome Biology* **8**, R5 (2007).
- [65] Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
- [66] DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics* **43**, 491–498 (2011).
- [67] Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *American Journal of Human Genetics* **81**, 559–75 (2007).
- [68] Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* **4**, 7 (2015).
- [69] Delaneau, O., Howie, B., Cox, A. J., Zagury, J. F. & Marchini, J. Haplotype estimation using sequencing reads. *American Journal of Human Genetics* **93**, 687–696 (2013).

- 627 [70] Miles, A. *et al.* Genetic diversity of the African malaria vector *Anopheles gambiae*.  
628 *Nature* **552**, 96–100 (2017).
- 629 [71] Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,  
630 2078–9 (2009).
- 631 [72] Alexander, D., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in  
632 unrelated individuals. *Genome Research* **19**, 1655–1664 (2009).
- 633 [73] Kopelman, N. M., Mayzel, J., Jakobsson, M., Rosenberg, N. A. & Mayrose, I.  
634 CLUMPAK: a program for identifying clustering modes and packaging population struc-  
635 ture inferences across K. *Molecular Ecology Resources* **15**, 1179–1191 (2015).
- 636 [74] Hudson, R. R., Slatkin, M. & Maddison, W. P. Estimation of levels of gene flow from  
637 DNA sequence data. *Genetics* **132**, 583–589 (1992).
- 638 [75] Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS*  
639 *Genetics* **2**, 2074–2093 (2006).
- 640 [76] Price, A. *et al.* Principal components analysis corrects for stratification in genome-wide  
641 association studies. *Nature Genetics* **38**, 904–9 (2006).
- 642 [77] Mantel, N. The detection of disease clustering and a generalized regression approach.  
643 *Cancer Research* **27**, 209–220 (1967).
- 644 [78] Dray, S. & Dufour, A. B. The ade4 package: implementing the duality diagram for  
645 ecologists. *Journal of Statistical Software* **22**, 1–20 (2007).
- 646 [79] Hare, M. P. *et al.* Understanding and estimating effective population size for practical  
647 application in marine species management. *Conservation Biology* **25**, 438–449 (2011).

- 648 [80] Liu, X. & Fu, Y.-X. Exploring population size changes using SNP frequency spectra.  
649 *Nature Genetics* **47**, 555–559 (2015).
- 650 [81] Browning, B. L. & Browning, S. R. Detecting identity by descent and estimating  
651 genotype error rates in sequence data. *American Journal of Human Genetics* **93**, 840–  
652 851 (2013).
- 653 [82] Browning, S. R. & Browning, B. L. Accurate non-parametric estimation of recent  
654 effective population size from segments of identity by descent. *American Journal of*  
655 *Human Genetics* **97**, 404–418 (2015).
- 656 [83] Clarkson, C. S. *et al.* The genetic architecture of target-site resistance to pyrethroid  
657 insecticides in the African malaria vectors *Anopheles gambiae* and *Anopheles coluzzii*.  
658 *bioRxiv* Preprint at: <https://www.biorxiv.org/content/early/2018/08/06/323980>  
659 (2018).
- 660 [84] Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H. & Bustamante, C. D. Infer-  
661 ring the joint demographic history of multiple populations from multidimensional SNP  
662 frequency data. *PLoS Genetics* **5**, e1000695 (2009).
- 663 [85] Coffman, A. J., Hsieh, P. H., Gravel, S. & Gutenkunst, R. N. Computationally effi-  
664 cient composite likelihood statistics for demographic inference. *Molecular Biology and*  
665 *Evolution* **33**, 591–593 (2016).
- 666 [86] Fontaine, M. C. *et al.* Extensive introgression in a malaria vector species complex  
667 revealed by phylogenomics. *Science* **347**, 1258524 (2014).
- 668 [87] Weir, B. & Cockerham, C. Estimating F-statistics for the analysis of population struc-  
669 ture. *Evolution* **38**, 1358–1370 (1984).

- 670 [88] Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–8  
671 (2011).
- 672 [89] Garud, N. R., Messer, P. W., Buzbas, E. O. & Petrov, D. A. Recent selective sweeps in  
673 North American *Drosophila melanogaster* show signatures of soft sweeps. *PLoS Genetics*  
674 **11**, 1–32 (2015).
- 675 [90] Szpiech, Z. A. & Hernandez, R. D. selscan: an efficient multithreaded program to  
676 perform EHH-based scans for positive selection. *Molecular Biology and Evolution* **31**,  
677 2824–2827 (2014).
- 678 [91] Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive selection in  
679 human populations. *Nature* **449**, 913–918 (2007).
- 680 [92] Randhawa, I. A., Khatkar, M., Thomson, P. & Raadsma, H. Composite selection signals  
681 can localize the trait specific genomic regions in multi-breed populations of cattle and  
682 sheep. *BMC Genetics* **15**, 34 (2014).
- 683 [93] Wallberg, A., Pirk, C. W., Allsopp, M. H. & Webster, M. T. Identification of multiple  
684 loci associated with social parasitism in honeybees. *PLoS Genetics* **12**, 1–30 (2016).

685 **Supplemental Material for: Assessing connectivity**  
686 **despite high diversity in island populations of the**  
687 **malaria mosquito *Anopheles gambiae***

688 Christina M. Bergey, Martin Lukindu, Rachel M. Wiltshire, Michael C. Fontaine, Jonathan  
689 K. Kayondo, and Nora J. Besansky



690 **Tables**

Table S1: Sampling sites and coordinates.

Location	Latitude	Longitude	Sample Count
Banda	-0.25893	32.39594	11
Bugala - Bugoma	-0.26697	32.07936	11
Bugala - Lutoboka	-0.31624	32.29246	7
Bugala - Mweena	-0.32806	32.31113	5
Bukasa	-0.48609	32.45091	11
Buwama	0.02077	32.10574	11
Kaazi	-0.31831	31.88183	11
Kiyindi	0.27558	33.14699	10
Nsadzi	-0.08632	32.58895	11
Sserinya	-0.26476	32.37228	16
Wamala	0.40811	31.99609	11

Table S2: List of individuals included in study with mean depth of sequencing coverage.

ID	Field ID	Island	Site	Mean depth
LVB2015-1	CM-KSB-J5	Nsadzi	Kansambwe	20.40
LVB2015-2	K-KSB-E1	Nsadzi	Kansambwe	24.50
LVB2015-3	RM-KSB-G1	Nsadzi	Kansambwe	17.90
LVB2015-4	NKG-F-G3	Bukasa	Nakibanga	4.90
LVB2015-6	NKG-F-H1	Bukasa	Nakibanga	15.80
LVB2015-7	NKG-K-I1	Bukasa	Nakibanga	15.80
LVB2015-8	NKG-K-K1	Bukasa	Nakibanga	19.90
LVB2015-9	NKG-M-C1	Bukasa	Nakibanga	22.30
LVB2015-10	NKG-M-D1	Bukasa	Nakibanga	20.30
LVB2015-11	NKG-M-F1	Bukasa	Nakibanga	23.90
LVB2015-14	MWN-K-A1	Bugala	Mweena	21.20
LVB2015-15	MWN-K-C2	Bugala	Mweena	18.30
LVB2015-16	MWN-P-D1	Bugala	Mweena	5.34
LVB2015-17	MWN-R-E1	Bugala	Mweena	22.60
LVB2015-18	MWN-R-F1	Bugala	Mweena	17.10
LVB2015-19	BDA-K-B1	Banda	Banda	14.40
LVB2015-20	BDA-K-B2	Banda	Banda	18.00

LVB2015-21	BBS-C-M1	Sserinya	Bbosa	19.70
LVB2015-22	BBS-F-F1	Sserinya	Bbosa	20.80
LVB2015-24	BBS-K-J3	Sserinya	Bbosa	22.80
LVB2015-25	BBS-K-J8	Sserinya	Bbosa	17.30
LVB2015-26	BBS-K-K2	Sserinya	Bbosa	18.20
LVB2015-27	BBS-M-L1	Sserinya	Bbosa	22.40
LVB2015-28	BBS-P-I4	Sserinya	Bbosa	23.30
LVB2015-29	BBS-R-A2	Sserinya	Bbosa	19.60
LVB2015-30	BBS-R-C1	Sserinya	Bbosa	16.10
LVB2015-32	KSS-F-E2	Sserinya	Kasisa	21.00
LVB2015-33	LBK-C-F1	Bugala	Lutoboka	21.00
LVB2015-34	LBK-C-F6	Bugala	Lutoboka	18.60
LVB2015-35	LBK-C-G6	Bugala	Lutoboka	20.90
LVB2015-36	LBK-K-E2	Bugala	Lutoboka	22.00
LVB2015-37	LBK-M-A1	Bugala	Lutoboka	18.50
LVB2015-39	LBK-R-O1	Bugala	Lutoboka	20.20
LVB2015-42	BGM-F-D1	Bugala	Bugoma	16.70
LVB2015-43	BGM-F-E2	Bugala	Bugoma	24.50
LVB2015-45	BGM-K-M2	Bugala	Bugoma	23.80
LVB2015-46	BGM-M-G1	Bugala	Bugoma	18.40
LVB2015-47	BGM-M-H2	Bugala	Bugoma	14.10
LVB2015-48	BGM-M-J1	Bugala	Bugoma	20.90
LVB2015-50	BGM-P-F9	Bugala	Bugoma	18.90
LVB2015-51	BGM-R-O2	Bugala	Bugoma	16.80
LVB2015-52	KZI-F-F001	Kaazi	Nabugabo	19.10
LVB2015-53	KZI-F-G001	Kaazi	Nabugabo	19.50
LVB2015-54	KZI-F-H001	Kaazi	Nabugabo	19.40
LVB2015-55	KZI-P-A001	Kaazi	Nabugabo	10.30
LVB2015-56	KZI-P-B005	Kaazi	Nabugabo	16.50
LVB2015-59	KZI-R-C003	Kaazi	Nabugabo	18.10
LVB2015-60	KZI-R-D007	Kaazi	Nabugabo	15.90
LVB2015-61	BWM-C-G001	Buwama	Buwama	16.10
LVB2015-62	BWM-C-H001	Buwama	Buwama	11.10
LVB2015-63	BWM-F-A001	Buwama	Buwama	20.40
LVB2015-64	BWM-F-B001	Buwama	Buwama	21.50
LVB2015-65	BWM-P-J001	Buwama	Buwama	14.80
LVB2015-66	BWM-R-C002	Buwama	Buwama	19.10

LVB2015-67	BWM-R-F005	Buwama	Buwama	22.30
LVB2015-68	NMA-C-E003	Wamala	Naama	20.30
LVB2015-69	NMA-C-F002	Wamala	Naama	17.80
LVB2015-70	NMA-F-A001	Wamala	Naama	13.10
LVB2015-71	NMA-K-B001	Wamala	Naama	22.10
LVB2015-72	NMA-K-C002	Wamala	Naama	18.00
LVB2015-73	NMA-P-G001	Wamala	Naama	18.20
LVB2015-74	NMA-P-H003	Wamala	Naama	16.60
LVB2015-76	KYD-C-G001	Kiyindi	Kiyindi	16.10
LVB2015-77	KYD-C-H001	Kiyindi	Kiyindi	11.80
LVB2015-78	KYD-C-I001	Kiyindi	Kiyindi	16.40
LVB2015-79	KYD-C-J002	Kiyindi	Kiyindi	11.50
LVB2015-80	KYD-F-A003	Kiyindi	Kiyindi	10.30
LVB2015-81	KYD-F-B004	Kiyindi	Kiyindi	21.50
LVB2015-82	KYD-K-D002	Kiyindi	Kiyindi	18.40
LVB2015-84	KYD-R-K001	Kiyindi	Kiyindi	16.80
LVB2015-89	BDA-K-E2	Banda	Banda	15.10
LVB2015-90	BDA-K-F1	Banda	Banda	25.10
LVB2015-91	BDA-M-N1	Banda	Banda	25.60
LVB2015-92	BDA-M-O4	Banda	Banda	17.60
LVB2015-93	BDA-M-Q1	Banda	Banda	39.20
LVB2015-96	CM-KSB-J2	Nsadzi	Kansambwe	9.22
LVB2015-97	CM-KSB-J3	Nsadzi	Kansambwe	10.10
LVB2015-98	CM-KSB-J6	Nsadzi	Kansambwe	16.90
LVB2015-100	K-KSB-D1	Nsadzi	Kansambwe	6.05
LVB2015-101	ML-KSB-M1	Nsadzi	Kansambwe	4.27
LVB2015-102	ML-KSB-M2	Nsadzi	Kansambwe	19.90
LVB2015-103	RM-KSB-G2	Nsadzi	Kansambwe	14.20
LVB2015-104	RM-KSB-G3	Nsadzi	Kansambwe	17.50
LVB2015-105	NKG-R-A12	Bukasa	Nakibanga	15.30
LVB2015-106	NKG-C-E1	Bukasa	Nakibanga	16.20
LVB2015-108	NKG-K-C5	Bukasa	Nakibanga	18.50
LVB2015-109	NKG-M-A1	Bukasa	Nakibanga	12.80
LVB2015-112	BDA-K-D4	Banda	Banda	12.70
LVB2015-113	BDA-K-E3	Banda	Banda	12.20
LVB2015-114	BDA-M-N5	Banda	Banda	15.00
LVB2015-115	BDA-M-P1	Banda	Banda	16.80

LVB2015-116	BBS-C-M3	Sserinya	Bbosa	16.60
LVB2015-117	BBS-K-J1	Sserinya	Bbosa	18.80
LVB2015-118	BBS-K-J11	Sserinya	Bbosa	14.60
LVB2015-120	BBS-K-K6	Sserinya	Bbosa	18.10
LVB2015-121	BBS-P-I8	Sserinya	Bbosa	15.00
LVB2015-122	BBS-R-A19	Sserinya	Bbosa	15.50
LVB2015-125	LBK-R-A5	Bugala	Lutoboka	18.10
LVB2015-126	BGM-K-K1	Bugala	Bugoma	15.20
LVB2015-128	BGM-M-H4	Bugala	Bugoma	20.30
LVB2015-129	BGM-P-F4	Bugala	Bugoma	19.00
LVB2015-130	KZI-F-G005	Kaazi	Nabugabo	18.60
LVB2015-131	KZI-P-A007	Kaazi	Nabugabo	15.80
LVB2015-132	KZI-R-C012	Kaazi	Nabugabo	15.10
LVB2015-133	KZI-R-E011	Kaazi	Nabugabo	16.30
LVB2015-134	BWM-P-I001	Buwama	Buwama	18.20
LVB2015-135	BWM-P-K002	Buwama	Buwama	19.30
LVB2015-136	BWM-R-D001	Buwama	Buwama	14.40
LVB2015-137	BWM-R-F002	Buwama	Buwama	19.90
LVB2015-138	NMA-C-E006	Wamala	Naama	21.90
LVB2015-139	NMA-C-F003	Wamala	Naama	20.40
LVB2015-140	NMA-P-G003	Wamala	Naama	18.90
LVB2015-141	NMA-R-I001	Wamala	Naama	14.10
LVB2015-142	KYD-F-B006	Kiyindi	Kiyindi	18.10
LVB2015-143	KYD-K-E003	Kiyindi	Kiyindi	14.20

Table S3: Genomic coordinates of heterochromatic and inverted regions.

Chromosome arm	Start	End	Information
2L	20,524,058	42,165,532	2La inversion [94]
2R	18,575,300	26,767,588	2Rb inversion [94]
2L	1	2,431,617	Heterochromatic region [94]
2L	5,078,962	5,788,875	Heterochromatic region [94]
2R	58,984,778	61,545,105	Heterochromatic region [94]
3L	1	1,815,119	Heterochromatic region [94]
3L	4,264,713	5,031,692	Heterochromatic region [94]
3R	38,988,757	41,860,198	Heterochromatic region [94]

3R	52,161,877	53,200,684	Heterochromatic region [94]
X	20,009,764	24,393,108	Heterochromatic region [94]

---

---

Table S4: Results of two population demographic inference with IM model in  $\delta a\delta i$  when comparing island to island localities. Numbers in parentheses are bounds of 95% confidence interval computed using Fisher information matrix and 100 bootstrap replicates of 1 Mb from the dataset.

Localities	$N_a$	% Pop. 1 in Split	Pop. 1 $\nu_F$	Pop. 2 $\nu_F$	Time since split	$m_{12}$	$m_{21}$
Banda - Bugala (I)	762,000 (762,000, 763,000)	0.538 (0.522, 0.555)	2.79 (2.55, 3.03)	9,570 (8,000, 11,100)	5,050 (4,900, 5,200)	None	None
Banda - Bukasa	755,000 (754,000, 756,000)	0.595 (0.59, 0.599)	8.53 (8.22, 8.84)	9,850 (9,030, 10,700)	9,760 (9,660, 9,850)	None	None
Banda - Nsadzi	760,000 (759,000, 761,000)	0.557 (0.551, 0.562)	40.3 (37.6, 43)	8,550 (7,500, 9,590)	13,800 (13,600, 14,000)	None	None
Banda - Sserinya	764,000 (763,000, 765,000)	0.497 (0.489, 0.506)	4.74 (4.23, 5.25)	8,460 (6,590, 10,300)	5,100 (4,930, 5,280)	None	None
Bugala (I) - Bukasa	759,000 (758,000, 760,000)	0.588 (0.575, 0.6)	9,330 (7,690, 11,000)	2,960 (2,600, 3,320)	8,400 (8,210, 8,580)	None	None
Bugala (I) - Nsadzi	759,000 (758,000, 760,000)	0.499 (0.492, 0.505)	9,350 (7,740, 11,000)	30.2 (26.2, 34.2)	8,380 (8,170, 8,590)	None	None
Bugala (I) - Sserinya	763,000 (762,000, 764,000)	0.592 (0.574, 0.61)	7,940 (6,530, 9,350)	593 (513, 673)	4,860 (4,720, 5,010)	None	None
Bukasa - Nsadzi	759,000 (758,000, 760,000)	0.436 (0.427, 0.446)	9,720 (8,270, 11,200)	77.4 (67.2, 87.6)	12,200 (11,900, 12,400)	None	None
Bukasa - Sserinya	755,000 (754,000, 756,000)	0.493 (0.488, 0.497)	9,960 (8,190, 11,700)	5,090 (4,470, 5,710)	12,700 (12,600, 12,900)	None	None
Nsadzi - Sserinya	777,000 (776,000, 778,000)	0.615 (0.594, 0.635)	53.9 (33.2, 74.5)	9,210 (7,420, 11,000)	10,500 (9,940, 11,100)	None	None

Table S5: Results of two population demographic inference with IM model in  $\delta a\delta i$  when comparing island to mainland localities. Numbers in parentheses are bounds of 95% confidence interval computed using Fisher information matrix and 100 bootstrap replicates of 1 Mb from the dataset.

Localities	$N_a$	% Pop. 1 in Split	Pop. 1 $\nu_F$	Pop. 2 $\nu_F$	Time since split	$m_{12}$	$m_{21}$
Banda - Bugala (M)	751,000 (750,000, 752,000)	0.522 (0.511, 0.532)	4.38 (4.15, 4.61)	8,580 (7,470, 9,690)	7,610 (7,460, 7,750)	None	None
Banda - Buwama	751,000 (751,000, 752,000)	0.457 (0.439, 0.476)	1.92 (1.68, 2.16)	7,470 (6,460, 8,480)	5,160 (4,960, 5,350)	None	None
Banda - Kaazi	752,000 (751,000, 753,000)	0.477 (0.466, 0.488)	3.75 (3.52, 3.99)	9,040 (7,680, 10,400)	7,550 (7,400, 7,690)	None	None
Banda - Kiyindi	735,000 (734,000, 736,000)	0.511 (0.501, 0.52)	1.82 (1.76, 1.87)	9,360 (8,160, 10,600)	5,450 (5,410, 5,490)	None	None
Banda - Wamala	750,000 (749,000, 751,000)	0.596 (0.586, 0.606)	2.2 (2.09, 2.3)	8,740 (7,610, 9,860)	6,600 (6,480, 6,720)	None	None
Bugala (I) - Bugala (M)	752,000 (751,000, 753,000)	0.496 (0.484, 0.508)	3,090 (2,490, 3,690)	8,940 (7,420, 10,500)	6,580 (6,450, 6,710)	None	None
Bugala (I) - Buwama	753,000 (753,000, 754,000)	0.5 (0.497, 0.502)	0.198 (0.194, 0.202)	9,330 (-63,800, 82,500)	274 (271, 276)	None	None
Bugala (I) - Kaazi	753,000 (752,000, 754,000)	0.401 (0.388, 0.414)	1,500 (1,350, 1,640)	9,240 (7,950, 10,500)	7,780 (7,630, 7,940)	None	None
Bugala (I) - Kiyindi	735,000 (734,000, 736,000)	0.478 (0.47, 0.487)	0.14 (0.137, 0.143)	8,260 (-92,500, 109,000)	254 (250, 259)	None	None
Bugala (I) - Wamala	748,000 (748,000, 749,000)	0.479 (0.462, 0.496)	7,230 (5,220, 9,240)	6,590 (5,120, 8,050)	5,880 (5,730, 6,040)	None	None
Bugala (M) - Bukasa	751,000 (750,000, 751,000)	0.497 (0.49, 0.504)	9,080 (7,630, 10,500)	47.2 (42.5, 52)	6,870 (6,750, 7,000)	None	None
Bugala (M) - Nsadzi	771,000	0.381	9,030	3.95	5,610	None	None

	(770,000, 772,000)	(0.368, 0.394)	(7,690, 10,400)	(3.64, 4.26)	(5,480, 5,750)		
Bugala (M) - Sserinya	768,000	0.483	9,070	22.7	5,080	None	None
	(767,000, 769,000)	(0.471, 0.495)	(7,620, 10,500)	(21.5, 24)	(5,020, 5,140)		
Bukasa - Buwama	751,000	0.536	1.74	9,680	2,650	None	None
	(750,000, 752,000)	(0.527, 0.545)	(1.71, 1.77)	(7,860, 11,500)	(2,630, 2,670)		
Bukasa - Kaazi	751,000	0.533	16.2	7,730	6,690	None	None
	(750,000, 752,000)	(0.524, 0.543)	(15.7, 16.6)	(6,260, 9,200)	(6,620, 6,760)		
Bukasa - Kiyindi	733,000	0.549	4.64	9,170	4,120	None	None
	(732,000, 734,000)	(0.523, 0.575)	(3.93, 5.36)	(7,280, 11,100)	(3,900, 4,350)		
Bukasa - Wamala	748,000	0.361	281	6,790	7,320	None	None
	(747,000, 749,000)	(0.345, 0.377)	(222, 340)	(5,510, 8,070)	(7,110, 7,530)		
Buwama - Nsadzi	756,000	0.608	9,500	3.2	3,960	None	None
	(755,000, 756,000)	(0.593, 0.624)	(8,020, 11,000)	(2.98, 3.42)	(3,820, 4,090)		
Buwama - Sserinya	753,000	0.498	5,090	0.134	273	None	None
	(752,000, 754,000)	(0.495, 0.501)	(-12,800, 23,000)	(0.132, 0.136)	(271, 274)		
Kaazi - Nsadzi	756,000	0.516	9,050	11.6	7,350	None	None
	(755,000, 756,000)	(0.493, 0.539)	(7,510, 10,600)	(8.26, 14.9)	(6,970, 7,730)		
Kaazi - Sserinya	752,000	0.529	7,950	15	4,920	None	None
	(752,000, 753,000)	(0.511, 0.547)	(6,600, 9,300)	(13.4, 16.6)	(4,780, 5,060)		
Kiyindi - Nsadzi	738,000	0.534	9,990	3.79	4,800	None	None
	(738,000, 739,000)	(0.527, 0.541)	(8,580, 11,400)	(3.59, 3.99)	(4,750, 4,860)		
Kiyindi - Sserinya	735,000	0.5	5,400	0.129	267	None	None
	(734,000, 736,000)	(0.497, 0.502)	(-18,000, 28,800)	(0.127, 0.131)	(265, 269)		
Wamala - Nsadzi	755,000	0.485	9,160	5.71	5,890	None	None
	(754,000, 755,000)	(0.479, 0.491)	(7,870, 10,400)	(5.6, 5.81)	(5,860, 5,930)		
Wamala - Sserinya	751,000	0.636	9,240	152	5,290	None	None
	(750,000, 752,000)	(0.622, 0.65)	(7,720, 10,700)	(135, 169)	(5,140, 5,440)		



Table S6: Results of two population demographic inference with IM model in  $\delta a \delta i$  when comparing mainland to mainland localities. Numbers in parentheses are bounds of 95% confidence interval computed using Fisher information matrix and 100 bootstrap replicates of 1 Mb from the dataset.

Localities	$N_a$	% Pop. 1 in Split	Pop. 1 $\nu_F$	Pop. 2 $\nu_F$	Time since split	$m_{12}$	$m_{21}$
Bugala (M) - Buwama	750,000 (749,000, 751,000)	0.172 (0.164, 0.18)	1.58 (1.52, 1.63)	4.63 (4.35, 4.9)	341 (337, 344)	0 (-5.81, 5.81)	0.339 (-15.3, 16)
Bugala (M) - Kaazi	750,000 (749,000, 751,000)	0.517 (0.484, 0.55)	1,840 (1,430, 2,250)	6,270 (4,610, 7,930)	3,330 (3,230, 3,430)	None	None
Bugala (M) - Kiyindi	735,000 (734,000, 736,000)	0.552 (0.412, 0.693)	9,800 (2,120, 17,500)	267 (160, 373)	1,340 (733, 1,960)	None	None
Bugala (M) - Wamala	748,000 (747,000, 749,000)	0.325 (0.308, 0.342)	9.38 (8.92, 9.85)	18.8 (17.4, 20.2)	911 (899, 922)	0 (-16.9, 16.9)	1.89 (-30.7, 34.5)
Buwama - Kaazi	753,000 (752,000, 754,000)	0.352 (-10.2, 10.9)	1,230 (-1,800, 4,270)	614 (-900, 2,130)	180 (176, 183)	0 (-10,600, 10,600)	236 (204, 267)
Buwama - Kiyindi	738,000 (737,000, 739,000)	0.677 (-1,560, 1,560)	12.5 (9.48, 15.5)	4.56 (3.48, 5.64)	91 (74.2, 108)	0 (-524, 524)	0 (-375, 375)
Buwama - Wamala	753,000 (752,000, 754,000)	0.709 (0.621, 0.798)	9,280 (-510, 19,100)	6,750 (315, 13,200)	255 (177, 333)	718 (-36,500, 37,900)	1,660 (-74,700, 78,000)
Kaazi - Kiyindi	736,000 (735,000, 737,000)	0.447 (0.371, 0.523)	549 (376, 722)	7,600 (4,230, 11,000)	764 (661, 868)	143 (83.5, 202)	603 (-25,400, 26,600)
Kaazi - Wamala	749,000 (749,000, 750,000)	0.499 (0.493, 0.504)	5,890 (4,040, 7,730)	7,880 (5,300, 10,500)	1,720 (1,610, 1,830)	None	None
Kiyindi - Wamala	736,000 (735,000, 737,000)	0.5 (0.494, 0.505)	7,750 (3,650, 11,900)	6,520 (3,050, 10,000)	511 (463, 560)	3.98 (-1,690, 1,690)	47 (-20,600, 20,700)

Table S7: Locality-specific (in LVB) putative sweeps based on H12 statistic.

Site	Count	Chr.	Putative Sweeps	Other sites <sup>1</sup>
Banda	44	2L	28.6 Mb; 36 Mb; 36.4 Mb; 36.9 Mb; 37.6 Mb; 38.1 Mb; 39.1 Mb; 42.2 Mb; 43.4 Mb; 43.8 Mb; 44.3 Mb; 44.9 Mb; 45.4 Mb	1 also found in BFS, GNS
		2R	4.2 Mb; 12.3 Mb; 18.3 Mb; 23.6 Mb; 29.4 Mb; 30.3 Mb; 33.7 Mb; 34.8 Mb; 35.8 Mb; 36.5 Mb; 44.1 Mb; 44.6 Mb; 49.7 Mb	1 also found in BFM, BFS, CMS, GNS, GWA; 1 also found in BFM, GWA; 5 also found in GWA
		3L	18.5 Mb; 21.6 Mb; 23.4 Mb; 23.9 Mb; 32.8 Mb	
		3R	2.6 Mb; 7.9 Mb; 29.2 Mb; 30.5 Mb; 31.3 Mb; 32.1 Mb; 33.2 Mb; 45.3 Mb; 46.4 Mb; 47 Mb	1 also found in GNS
		X	0.5 Mb; 2.1 Mb; 4.3 Mb	1 also found in AOM
Bugala (I)	24	2L	2.5 Mb; 5.5 Mb; 7.1 Mb; 19 Mb; 31.1 Mb; 43 Mb; 45.7 Mb	1 also found in AOM, BFM, BFS, CMS, GAS, GNS, UGS; 1 also found in AOM, UGS
		2R	6.7 Mb; 21.1 Mb; 24 Mb; 24.6 Mb; 35.6 Mb; 37.1 Mb; 38.6 Mb; 39 Mb; 55.9 Mb	1 also found in BFM, GWA; 2 also found in GWA
		3L	17.2 Mb; 29.5 Mb	
		3R	26 Mb; 35.8 Mb; 37.5 Mb	
		X	3.5 Mb; 5.7 Mb; 10.8 Mb	
Bukasa	112	2L	12.6 Mb; 13.6 Mb; 17.7 Mb; 20.1 Mb; 20.9 Mb; 21.6 Mb; 22.7 Mb; 23.6 Mb; 24.7 Mb; 25.4 Mb; 26.2 Mb; 26.9 Mb; 27.3 Mb; 27.8 Mb; 28.4 Mb; 29.1 Mb; 30.1 Mb; 31.5 Mb; 32.3 Mb; 33.3 Mb; 35.8 Mb; 39.4 Mb; 39.8 Mb; 40.6 Mb; 41.4 Mb; 43.1 Mb; 45.6 Mb; 48.1 Mb; 49.3 Mb	1 also found in AOM, BFM, BFS, CMS, GAS, GNS; 1 also found in BFM, GAS; 1 also found in BFS, GAS, GNS; 1 also found in BFS, GNS; 2 also found in CMS; 2 also found in GAS; 1 also found in GWA

	2R	1.3 Mb; 4.7 Mb; 5.3 Mb; 7.2 Mb; 7.6 Mb; 8 Mb; 9.7 Mb; 10.5 Mb; 12 Mb; 12.4 Mb; 13.5 Mb; 14 Mb; 15.7 Mb; 16.9 Mb; 17.5 Mb; 19.4 Mb; 22.9 Mb; 24.9 Mb; 25.8 Mb; 26.6 Mb; 29.9 Mb; 30.8 Mb; 32.4 Mb; 33.4 Mb; 35.5 Mb; 37.6 Mb; 43 Mb; 45.6 Mb; 47 Mb; 49.5 Mb; 54.8 Mb	1 also found in AOM, GAS, GWA; 2 also found in BFM; 1 also found in BFM, GWA; 3 also found in GAS, GWA; 6 also found in GWA	
	3L	7.3 Mb; 11.6 Mb; 13.1 Mb; 15.6 Mb; 18.1 Mb; 19.1 Mb; 19.8 Mb; 20.6 Mb; 24.2 Mb; 25.2 Mb; 27.3 Mb; 28 Mb; 28.7 Mb; 29.7 Mb; 30.6 Mb; 33.7 Mb; 34.7 Mb; 35.3 Mb; 36.1 Mb; 38.7 Mb; 39.9 Mb; 40.3 Mb; 41.2 Mb	2 also found in BFM; 1 also found in GAS	
	3R	5.1 Mb; 5.9 Mb; 7.2 Mb; 8.9 Mb; 12.7 Mb; 13.3 Mb; 14.1 Mb; 14.9 Mb; 15.9 Mb; 17.2 Mb; 22.3 Mb; 23.3 Mb; 23.8 Mb; 24.9 Mb; 26.8 Mb; 27.9 Mb; 31.4 Mb; 33 Mb; 35.9 Mb; 36.9 Mb	1 also found in GWA	
	X	1.7 Mb; 2.8 Mb; 4.9 Mb; 6 Mb; 7 Mb; 11.5 Mb; 12.5 Mb; 13.6 Mb; 16.7 Mb	1 also found in BFM, GAS, GWA; 4 also found in GAS	
Buwama	27	2L	14.9 Mb; 15.9 Mb; 25.1 Mb; 26.5 Mb; 31.6 Mb	1 also found in BFM, GAS; 1 also found in GWA
		2R	24.4 Mb; 39.5 Mb; 44.5 Mb; 46.3 Mb; 49.1 Mb; 53.7 Mb; 55.3 Mb	1 also found in AOM, BFM, CMS; 1 also found in BFS, CMS, GNS; 1 also found in CMS; 1 also found in GWA
		3L	2.4 Mb; 3.1 Mb; 3.6 Mb; 4.1 Mb; 10.6 Mb; 16.1 Mb; 21.7 Mb; 29.8 Mb	
		3R	18 Mb; 29.1 Mb; 35.5 Mb; 37.7 Mb; 38.4 Mb; 38.9 Mb; 40.6 Mb	1 also found in GNS

Kaazi	15	2L	8.5 Mb; 34.6 Mb	1 also found in AOM
		2R	8.3 Mb; 23 Mb	1 also found in AOM; 1 also found in GAS
		3L	3.5 Mb; 4.8 Mb; 8.6 Mb; 11.8 Mb; 13 Mb; 15.8 Mb; 25 Mb; 26.8 Mb	1 also found in BFM; 1 also found in BFM, GNS
		3R	14.7 Mb; 15.6 Mb; 46 Mb	1 also found in GNS
Kiyindi	40	2L	2 Mb; 10.6 Mb; 17.8 Mb; 22.1 Mb; 23.9 Mb; 26 Mb; 28.7 Mb; 29.9 Mb; 34.8 Mb	1 also found in AOM, BFM, BFS, CMS, GNS, UGS; 1 also found in BFS, GNS; 3 also found in GAS
		2R	19.1 Mb; 20.2 Mb; 25.9 Mb; 35.3 Mb; 36.6 Mb; 38.1 Mb; 40 Mb; 41.7 Mb; 42.4 Mb; 45.3 Mb; 48.2 Mb; 48.6 Mb; 50.1 Mb; 52.2 Mb; 53.6 Mb; 54.7 Mb; 55.1 Mb	1 also found in AOM; 1 also found in AOM, BFS, CMS, GNS, GWA; 1 also found in BFM; 2 also found in GWA
		3L	1.2 Mb; 8.9 Mb; 12.1 Mb; 12.6 Mb; 13.5 Mb; 14.8 Mb; 15.4 Mb; 16 Mb; 16.8 Mb; 19.7 Mb; 26.7 Mb	1 also found in BFM; 1 also found in GWA
		3R	38 Mb; 41.9 Mb; 48.3 Mb	
Nsadzi	47	2L	23.2 Mb; 27 Mb; 45.5 Mb	
		2R	1.6 Mb; 2.3 Mb; 3.2 Mb; 4 Mb; 8.8 Mb; 10.2 Mb; 13.2 Mb; 16.1 Mb; 20 Mb; 21.3 Mb; 24.7 Mb; 30.5 Mb; 34.2 Mb; 37.3 Mb; 41.2 Mb; 43.5 Mb; 52 Mb	1 also found in BFM, GAS; 1 also found in BFM, GWA; 1 also found in BFS, CMS, GNS; 1 also found in CMS; 2 also found in GWA
		3L	10.5 Mb; 11 Mb; 24.3 Mb; 35 Mb; 35.4 Mb; 36.8 Mb; 37.6 Mb	1 also found in GAS
		3R	3.8 Mb; 6 Mb; 7.4 Mb; 19.9 Mb; 20.5 Mb; 21.4 Mb; 23 Mb; 24.2 Mb; 27.7 Mb; 41.6 Mb; 42.2 Mb; 48.2 Mb; 49.8 Mb; 50.4 Mb	1 also found in BFS, GNS
		X	0.7 Mb; 2.3 Mb; 5.2 Mb; 7.7 Mb; 11.9 Mb; 17.9 Mb	1 also found in BFM, GAS, GWA; 1 also found in GAS

Sserinya	35	2L	22.2 Mb; 24.2 Mb; 25.7 Mb; 33.2 Mb; 34.9 Mb; 35.4 Mb; 40.2 Mb; 41.1 Mb; 45.1 Mb; 45.9 Mb; 46.8 Mb	1 also found in BFM, GNS; 1 also found in CMS, GAS; 1 also found in GAS, GNS
		2R	0.4 Mb; 7.7 Mb; 21.5 Mb; 30 Mb; 32 Mb; 36.1 Mb	3 also found in GWA
		3L	10.1 Mb; 10.9 Mb; 14.6 Mb; 34.5 Mb; 41.8 Mb	1 also found in BFS, CMS, GNS, GWA, UGS
		3R	1.9 Mb; 10 Mb; 15 Mb; 24.8 Mb; 26.2 Mb; 27 Mb; 29 Mb	1 also found in GAS
		X	5.8 Mb; 12.7 Mb; 13.1 Mb; 18.1 Mb; 18.8 Mb; 21.3 Mb	1 also found in BFM, CMS, GAS, GWA; 1 also found in BFM, GAS, GWA; 1 also found in CMS, GNS, GWA; 2 also found in GAS
Wamala	25	2L	13.4 Mb; 15.5 Mb; 17.1 Mb; 19.1 Mb; 20 Mb	2 also found in GAS
		2R	21.2 Mb; 22.2 Mb; 29.6 Mb; 38.8 Mb; 39.7 Mb; 47.6 Mb; 48.3 Mb; 48.9 Mb	2 also found in AOM; 2 also found in GWA
		3L	3.3 Mb; 7.6 Mb; 8.2 Mb	
		3R	5 Mb; 39.2 Mb; 43.2 Mb; 46.5 Mb; 47.5 Mb; 48.5 Mb; 50.5 Mb; 50.9 Mb; 51.8 Mb	

50

691 <sup>1</sup> Ag1000G site codes: AOM: Angola [*coluzzii*]; BFM: Burkina Faso [*coluzzii*]; BFS: Burkina Faso [*gambiae*]; CMS: Cameroon [*gambiae*]; GAS: Gabon

692 [*gambiae*]; GNS: Guinea [*gambiae*]; GWA: Guinea-Bissau; UGS: Uganda [*gambiae*]

Table S8: Putative sweeps based on H12 statistic present on islands but rare or absent on LVB mainland.

Chr.	Region Start	Region End	Island Sites with Putative Sweep	Mainland Sites with Putative Sweep	Outlier Island Localities	Outlier Mainland Localities	Ag1000G Populations with Putative Sweep
2R	16,200,000	16,300,000	4 / 5	0 / 4	Banda; Bukasa; Nsadzi; Sserinya	None	Guinea-Bissau
2R	17,300,000	17,500,000	4 / 5	1 / 4	Banda; Bugala (I); Bukasa; Sserinya	Buwama	Guinea-Bissau
2R	21,000,000	21,100,000	5 / 5	1 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	Buwama	None
2R	40,400,000	40,800,000	4 / 5	1 / 4	Bugala (I); Bukasa; Nsadzi; Sserinya	Wamala	Burkina Faso [ <i>gambiae</i> ], Cameroon [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Guinea-Bissau
2R	41,100,000	41,200,000	4 / 5	1 / 4	Banda; Bukasa; Nsadzi; Sserinya	Wamala	Cameroon [ <i>gambiae</i> ]
2R	55,800,000	55,900,000	4 / 5	1 / 4	Banda; Bugala (I); Bukasa; Sserinya	Kiyindi	Angola [ <i>coluzzii</i> ]
2L	7,700,000	7,800,000	4 / 5	1 / 4	Banda; Bugala (I); Bukasa; Nsadzi	Buwama	Guinea-Bissau, Uganda [ <i>gambiae</i> ]
2L	8,100,000	8,200,000	4 / 5	0 / 4	Banda; Bugala (I); Nsadzi; Sserinya	None	None
2L	42,400,000	42,500,000	4 / 5	0 / 4	Banda; Bugala (I); Nsadzi; Sserinya	None	None
2L	43,500,000	43,600,000	5 / 5	1 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	Buwama	None
2L	49,000,000	49,100,000	4 / 5	1 / 4	Banda; Bugala (I); Bukasa; Sserinya	Wamala	None

3R	26,600,000	26,700,000	5 / 5	0 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	None	None
3R	36,700,000	36,800,000	4 / 5	0 / 4	Banda; Bugala (I); Bukasa; Nsadzi	None	None
3R	44,200,000	44,300,000	4 / 5	1 / 4	Banda; Bukasa; Nsadzi; Sserinya	Kiyindi	Angola [ <i>coluzzii</i> ]
3R	46,200,000	46,300,000	4 / 5	0 / 4	Banda; Bukasa; Nsadzi; Sserinya	None	None
X	6,600,000	7,000,000	4 / 5	0 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	None	Burkina Faso [ <i>coluzzii</i> ], Gabon [ <i>gambiae</i> ], Guinea-Bissau
X	8,100,000	10,700,000	4 / 5	0 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	Kiyindi	Burkina Faso [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ]
X	11,300,000	11,800,000	5 / 5	0 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	None	Gabon [ <i>gambiae</i> ]
X	12,900,000	13,000,000	4 / 5	0 / 4	Banda; Bugala (I); Bukasa; Sserinya	None	Gabon [ <i>gambiae</i> ]
X	14,300,000	14,400,000	5 / 5	1 / 4	Banda; Bugala (I); Bukasa; Nsadzi; Sserinya	Kaazi	Gabon [ <i>gambiae</i> ]
X	16,200,000	16,300,000	4 / 5	1 / 4	Banda; Bugala (I); Bukasa; Sserinya	Kaazi	Burkina Faso [ <i>coluzzii</i> ], Gabon [ <i>gambiae</i> ]

Table S9: Putative sweeps based on H12 statistic present on LVB mainland but rare or absent on islands.

Chr.	Region Start	Region End	Island Sites with Putative Sweep	Mainland Sites with Putative Sweep	Outlier Island Localities	Outlier Mainland Localities	Ag1000G Populations with Putative Sweep
2R	27,600,000	27,700,000	1 / 5	3 / 4	Nsadzi	Buwama; Kiyindi; Wamala	None
2R	38,000,000	38,100,000	1 / 5	3 / 4	Bugala (I)	Buwama; Kiyindi; Wamala	None
2R	42,700,000	42,800,000	0 / 5	3 / 4	None	Buwama; Kiyindi; Wamala	None
2R	45,400,000	45,500,000	1 / 5	3 / 4	Sserinya	Buwama; Kiyindi; Wamala	None
2R	46,800,000	46,900,000	1 / 5	3 / 4	Banda	Buwama; Kiyindi; Wamala	Cameroon [ <i>gambiae</i> ]
2R	48,000,000	48,100,000	1 / 5	3 / 4	Bukasa	Buwama; Kaazi; Wamala	Angola [ <i>coluzzii</i> ], Cameroon [ <i>gambiae</i> ]
2R	48,800,000	48,900,000	1 / 5	3 / 4	Nsadzi	Buwama; Kaazi; Wamala	None
2R	50,900,000	51,000,000	1 / 5	3 / 4	Bukasa	Kaazi; Kiyindi; Wamala	Burkina Faso [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ]
2R	51,500,000	51,600,000	0 / 5	3 / 4	None	Kaazi; Kiyindi; Wamala	None
2R	57,500,000	57,600,000	1 / 5	3 / 4	Banda	Buwama; Kaazi; Kiyindi	Angola [ <i>coluzzii</i> ], Guinea-Bissau
2L	2,900,000	3,000,000	1 / 5	4 / 4	Sserinya	Buwama; Kaazi; Kiyindi; Wamala	Angola [ <i>coluzzii</i> ], Burkina Faso [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Cameroon [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
2L	4,200,000	4,300,000	1 / 5	4 / 4	Bugala (I)	Buwama; Kaazi; Kiyindi; Wamala	Angola [ <i>coluzzii</i> ], Cameroon [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]



2L	5,700,000	5,800,000	1 / 5	3 / 4	Bugala (I)	Buwama; Kaazi; Kiyindi	Angola [ <i>coluzzii</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
2L	6,200,000	6,300,000	1 / 5	3 / 4	Bugala (I)	Kaazi; Kiyindi; Wamala	Uganda [ <i>gambiae</i> ]
2L	6,600,000	6,800,000	1 / 5	3 / 4	Bugala (I)	Kaazi; Kiyindi; Wamala	Angola [ <i>coluzzii</i> ], Cameroon [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
2L	10,000,000	10,100,000	1 / 5	3 / 4	Sserinya	Kaazi; Kiyindi; Wamala	None
2L	10,800,000	10,900,000	0 / 5	3 / 4	None	Kaazi; Kiyindi; Wamala	None
2L	11,300,000	11,400,000	1 / 5	3 / 4	Bugala (I)	Kaazi; Kiyindi; Wamala	Guinea-Bissau
2L	12,000,000	12,100,000	1 / 5	3 / 4	Bugala (I)	Kaazi; Kiyindi; Wamala	None
2L	12,400,000	13,000,000	0 / 5	3 / 4	Bukasa	Buwama; Kaazi; Kiyindi; Wamala	None
2L	14,500,000	14,900,000	1 / 5	3 / 4	Sserinya	Buwama; Kiyindi; Wamala	Gabon [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
2L	16,000,000	16,300,000	1 / 5	3 / 4	Bukasa	Buwama; Kaazi; Wamala	Gabon [ <i>gambiae</i> ]
2L	16,600,000	16,700,000	1 / 5	4 / 4	Bugala (I)	Buwama; Kaazi; Kiyindi; Wamala	None
2L	18,700,000	18,800,000	1 / 5	3 / 4	Nsadzi	Kaazi; Kiyindi; Wamala	None
2L	33,600,000	33,700,000	1 / 5	3 / 4	Bugala (I)	Buwama; Kaazi; Kiyindi	Angola [ <i>coluzzii</i> ]
2L	34,400,000	34,500,000	1 / 5	3 / 4	Sserinya	Buwama; Kaazi; Wamala	None

3R	28,500,000	28,700,000	1 / 5	4 / 4	Sserinya	Buwama; Kaazi; Kiyindi; Wamala	Burkina Faso [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Cameroon [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
3R	36,500,000	36,900,000	0 / 5	3 / 4	Nsadzi	Buwama; Kiyindi; Wamala	None
3R	43,000,000	43,100,000	0 / 5	3 / 4	None	Buwama; Kiyindi; Wamala	None
3R	43,700,000	44,100,000	0 / 5	3 / 4	Nsadzi	Buwama; Kiyindi; Wamala	Angola [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
3R	48,800,000	48,900,000	0 / 5	3 / 4	None	Buwama; Kiyindi; Wamala	None
3R	50,000,000	50,100,000	1 / 5	3 / 4	Sserinya	Kaazi; Kiyindi; Wamala	None
3L	7,000,000	7,100,000	1 / 5	4 / 4	Sserinya	Buwama; Kaazi; Kiyindi; Wamala	None
3L	11,500,000	11,600,000	1 / 5	3 / 4	Sserinya	Buwama; Kiyindi; Wamala	Burkina Faso [ <i>coluzzii</i> ]
3L	12,200,000	12,300,000	0 / 5	3 / 4	None	Kaazi; Kiyindi; Wamala	None
3L	13,400,000	13,500,000	0 / 5	3 / 4	None	Kaazi; Kiyindi; Wamala	None
3L	16,300,000	16,400,000	1 / 5	3 / 4	Sserinya	Buwama; Kiyindi; Wamala	Uganda [ <i>gambiae</i> ]

Table S10: Signatures of selective sweeps on known insecticide genes by site based on H12 statistic.

Chr.	Location	Insecticide Gene	Island Sites with Putative Sweep	Mainland Sites with Putative Sweep	Outlier		Ag1000G Populations with Putative Sweep
					Island Localities	Mainland Localities	
2R	28,497,407	Cyp6p	5 / 5	4 / 4	Banda; (I); Nsadzi; Sserinya	Bugala Buwama; Kaazi; Kiyindi; Wamala	Angola [ <i>coluzzii</i> ], Burkina Faso [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Cameroon [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
3R	28,598,038	Gste	1 / 5	4 / 4	Sserinya	Buwama; Kaazi; Kiyindi; Wamala	Burkina Faso [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Cameroon [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ], Uganda [ <i>gambiae</i> ]
X	15,241,718	Cyp9k1	3 / 5	4 / 4	Banda; Bukasa; Sserinya	Buwama; Kaazi; Kiyindi; Wamala	Burkina Faso [ <i>coluzzii</i> ], Burkina Faso [ <i>gambiae</i> ], Gabon [ <i>gambiae</i> ], Guinea [ <i>gambiae</i> ]

Table S11: Software and versions used for major parts of analysis.

Software	Version	Citation
ea-utils	-	[95]
BWA	0.7.16a	[96]
GATK	3.8	[97]
PLINK	1.90b4.6	[98, 99]
SHAPEIT2	2.837	[100]
SAMtools/BCFtools	1.5	[101, 102]
ADMIXTURE	1.3.0	[103]
CLUMPAK	-	[104]
VCFtools	0.1.15	[105]
$\delta\text{a}\delta\text{i}$ (python package)	1.7.0	[106, 107]
Stairway plot - Jpopgen	2-beta	[108]
selscan	1.2.0a	[109]
adegenet (R package)	2.1.0	[110]
ape (R package)	5.0	[111]
RColorBrewer (R package)	1.1-2	[112]
dendextend (R package)	1.6.0	[113]
rehh (R package)	2.0.2	[114]
eigensoft	7.2.1	[115, 116]
GNU parallel	20170422	[117]
tabix	1.5	[101]
bedtools	2.26.0	[118]

## 693 **Figures**

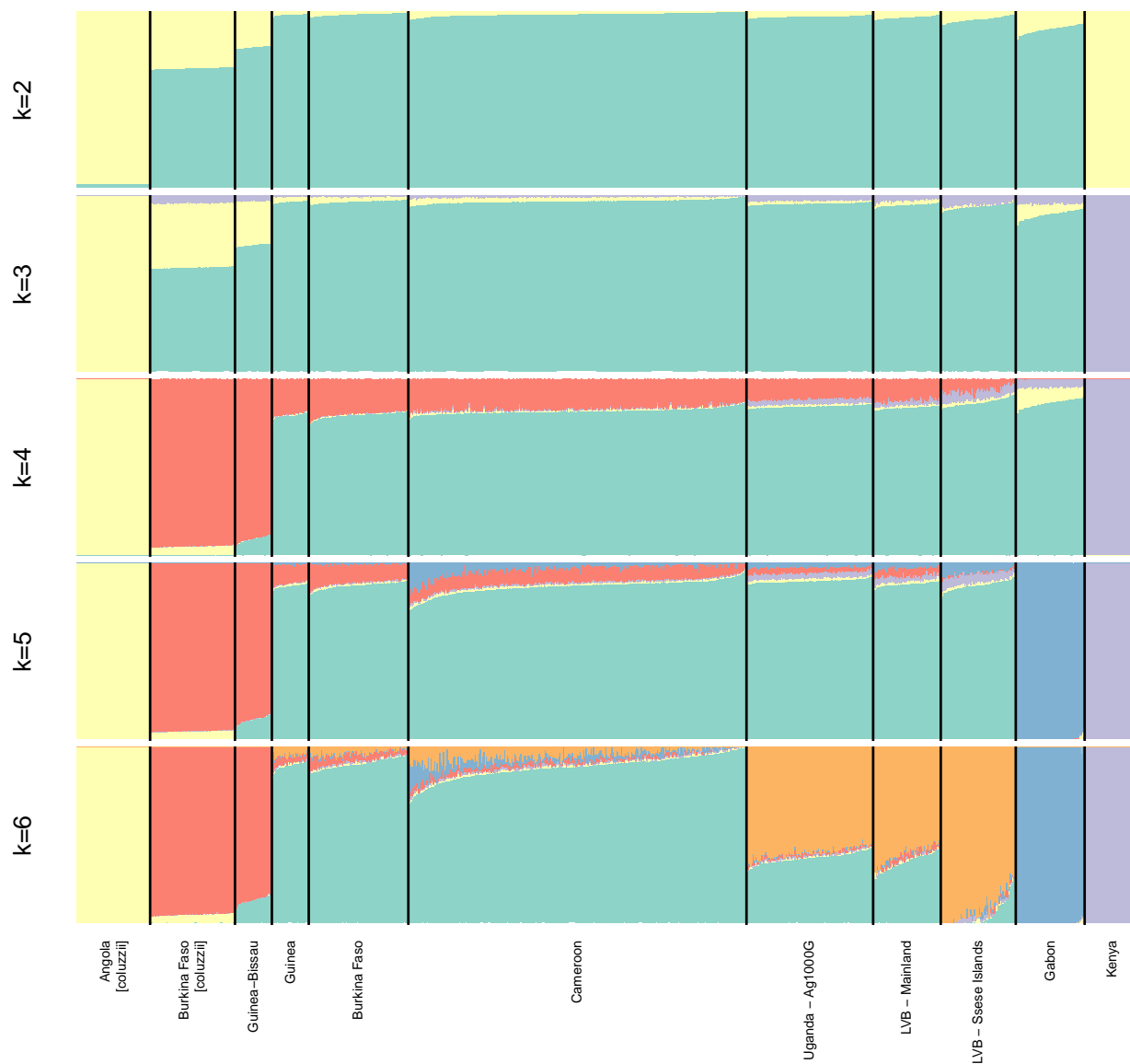


Figure S1: (Caption on next page.)

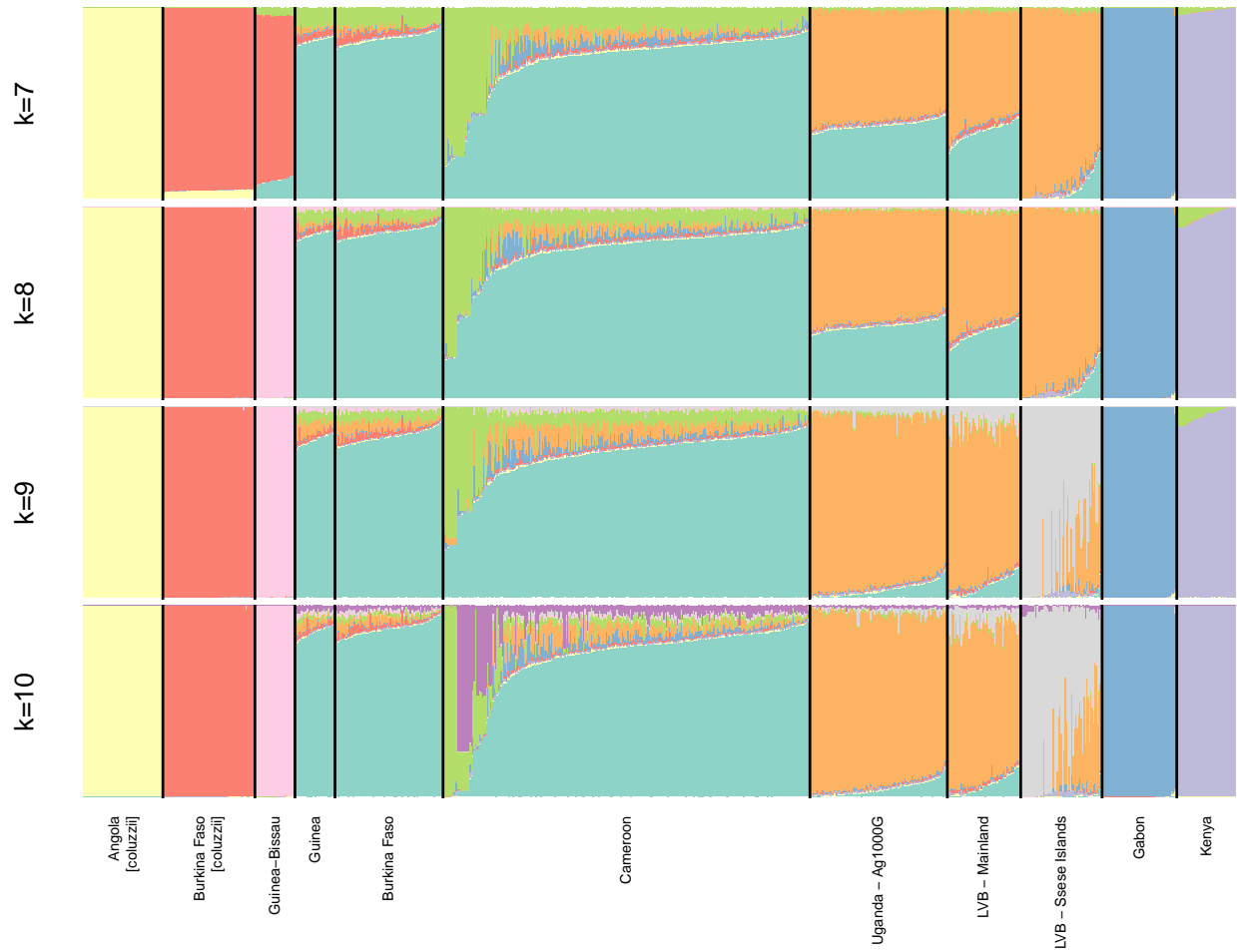


Figure S1: ADMIXTURE-inferred ancestry.

Ancestry of individuals in Lake Victoria Basin and of Ag1000G reference populations as inferred by ADMIXTURE clustering method. Samples are *A. gambiae* unless noted, and analysis is based on chromosome 3. Using  $k = 6$  clusters minimizes cross validation error (Fig. S2).

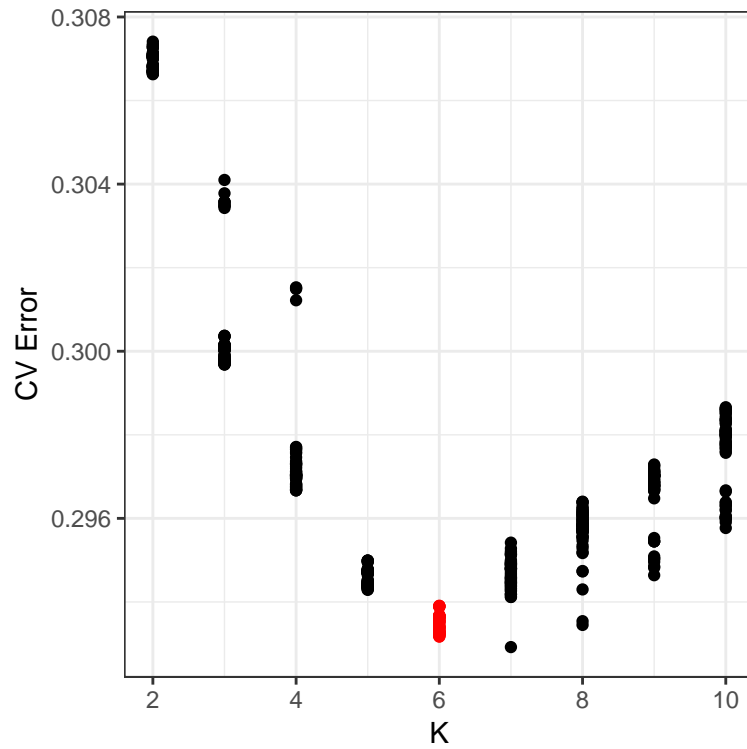


Figure S2: ADMIXTURE cross-validation error.

Cross-validation error for range of  $k$  values for ADMIXTURE analysis of Lake Victoria Basin individuals and *A. gambiae* and *A. coluzzii* Ag1000G reference populations.



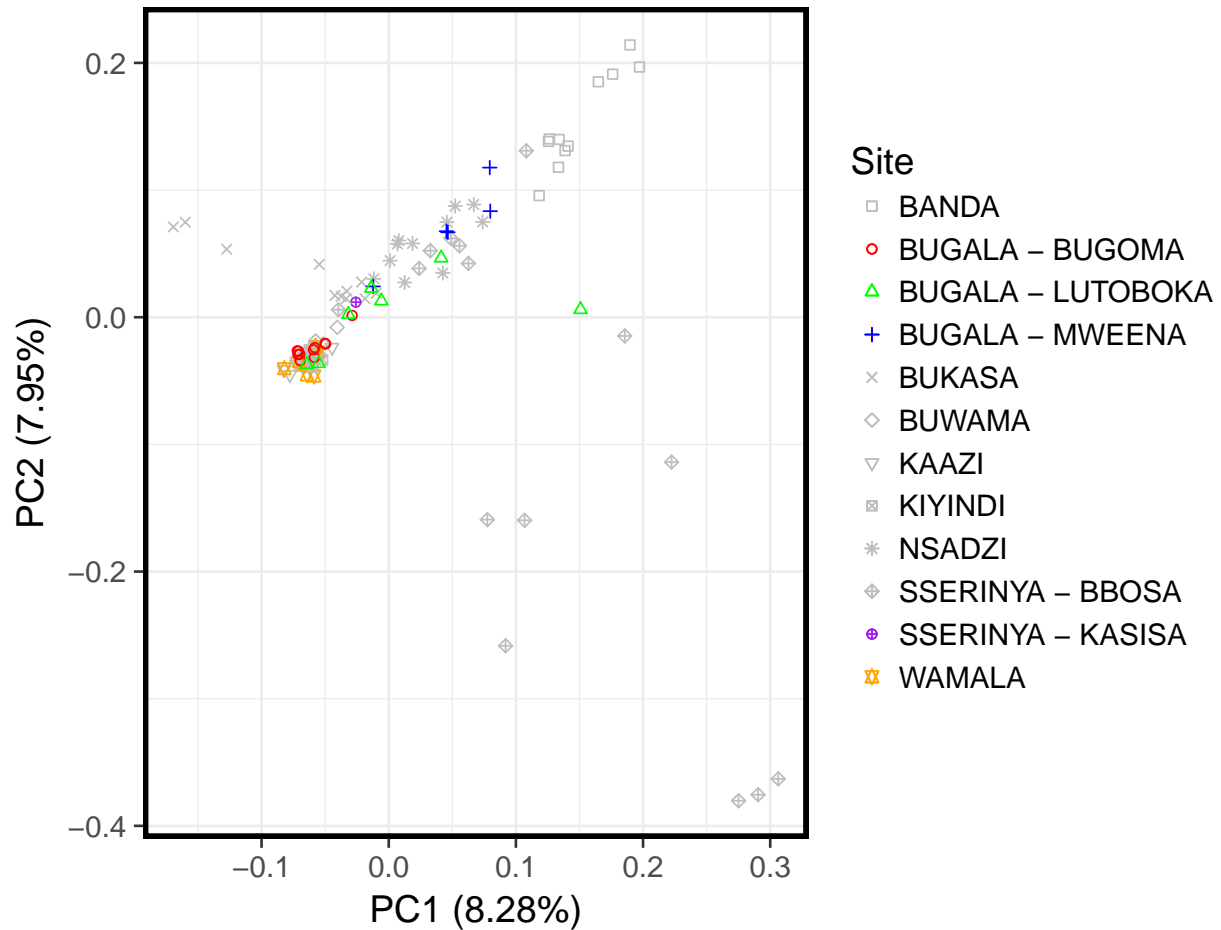


Figure S3: PCA showing Bugala subdivision.

PCA colored by sampling locations. Based on this analysis, individuals from Bugala were split into mainland- and island-like subpopulations. Samples from Sserinya Island, though sampled from two localities, were not split.

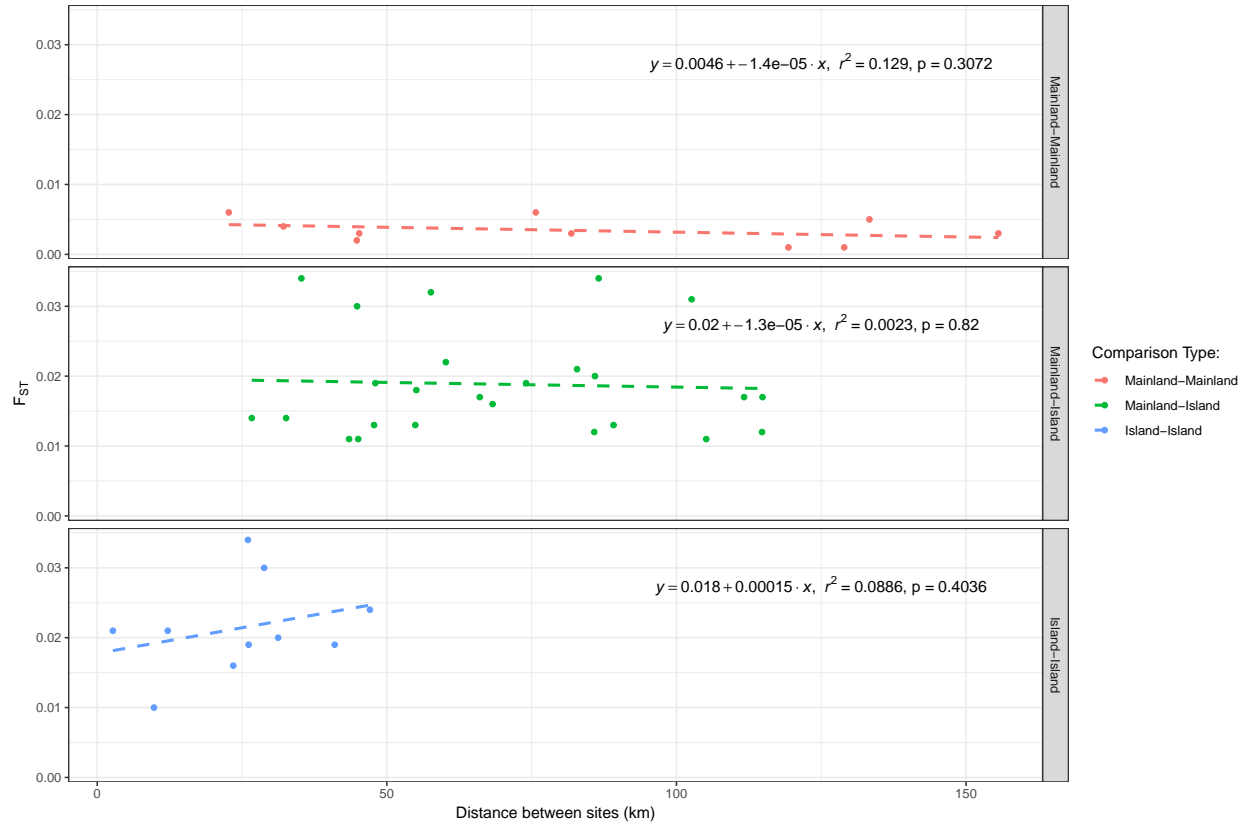


Figure S4: Correlations between genetic distance ( $F_{ST}$ ) and geographic distance between localities. The  $p$ -values are for the test that the slope is significantly different from zero.

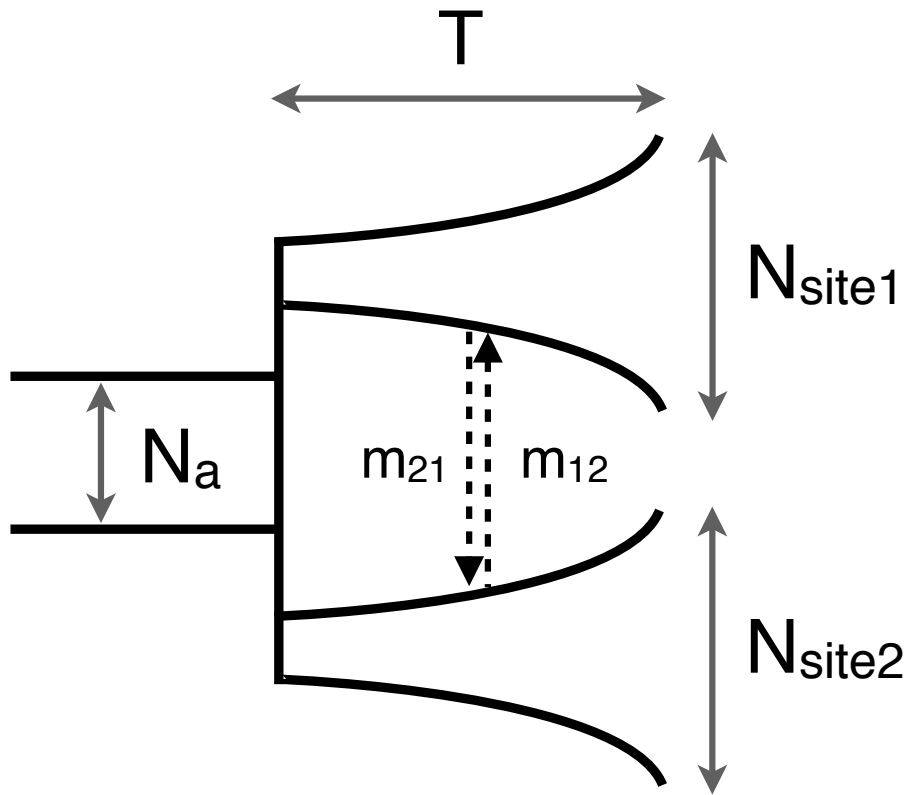


Figure S5: IM model schematics.

Schematic of model fit to data with  $\delta a \delta i$  for population history inference between all pairs of sampled sites using IM model.

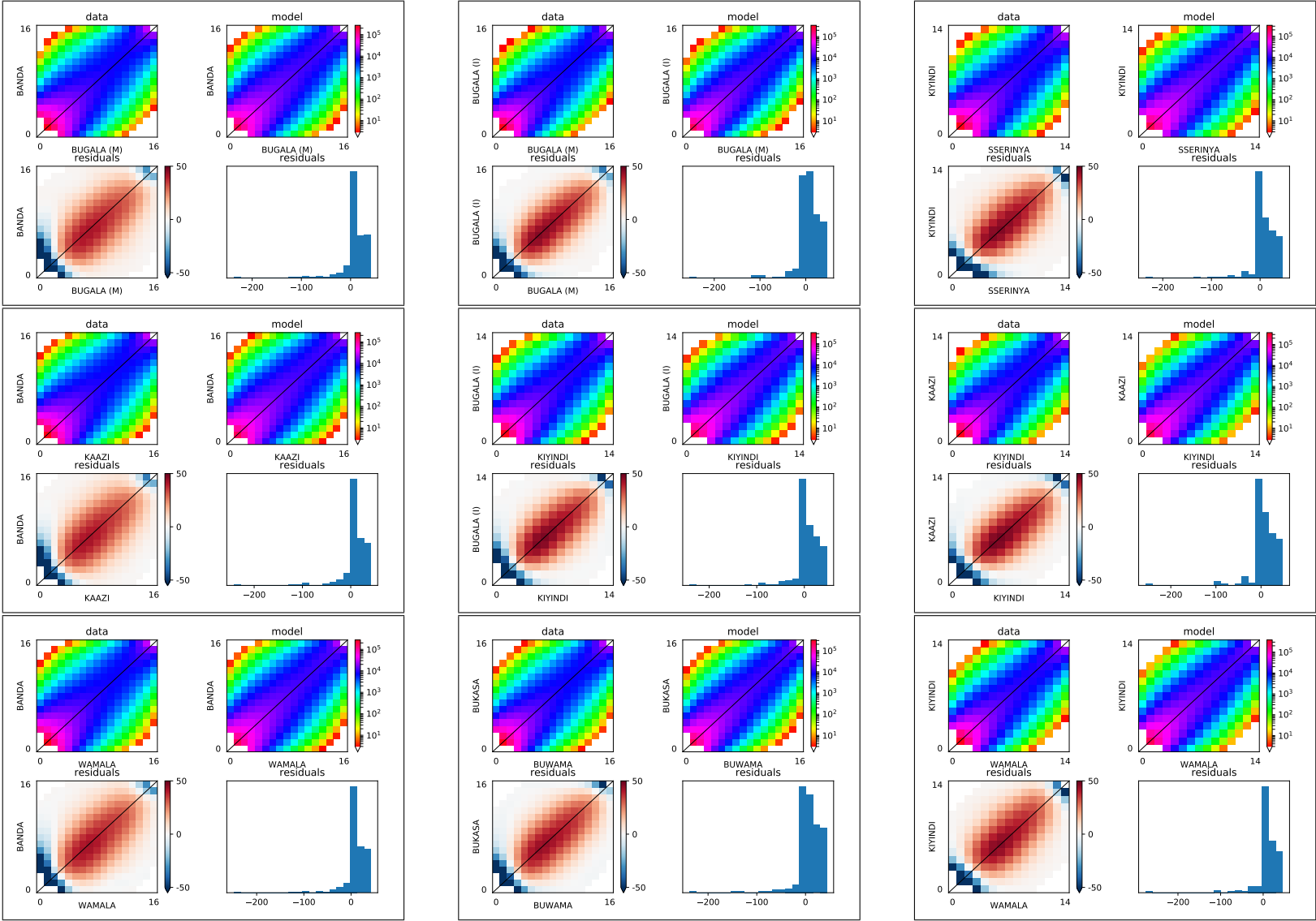


Figure S6: (Caption on next page.)

Figure S6: Two population  $\delta a \delta i$  optimization results.

Comparison between best fitting model and data frequency spectra for two population  $\delta a \delta i$  inference. Of the pairwise comparisons for which the best model included migration, a randomly selected set of nine are shown here. Two-dimensional frequency spectra are plotted as logarithmic colormaps for the data (upper left) and model (upper right), and the bottom row plots show the residuals between model and data. Positive residuals in red indicate the model predicts too many SNPs in that entry while negative residuals in blue indicate the model predicts too few.

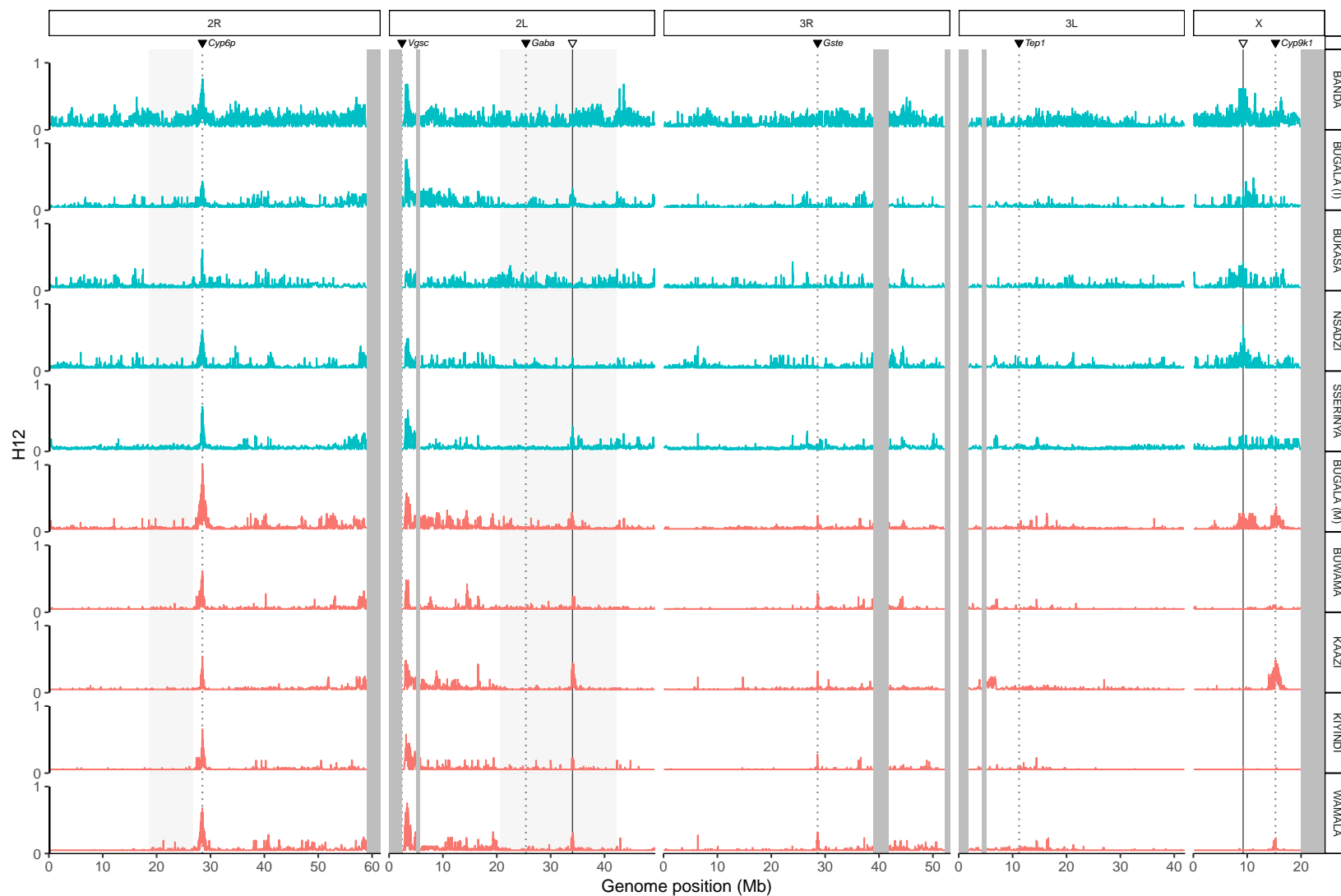


Figure S7: H12 across genome.

Values of H12, a measure of haplotype homozygosity, plotted across genome. Shaded regions indicate inversions or heterochromatic regions (excluded from analysis) and dotted lines indicate known insecticide genes while dashed lines indicate the two putative sweeps identified in the present study.

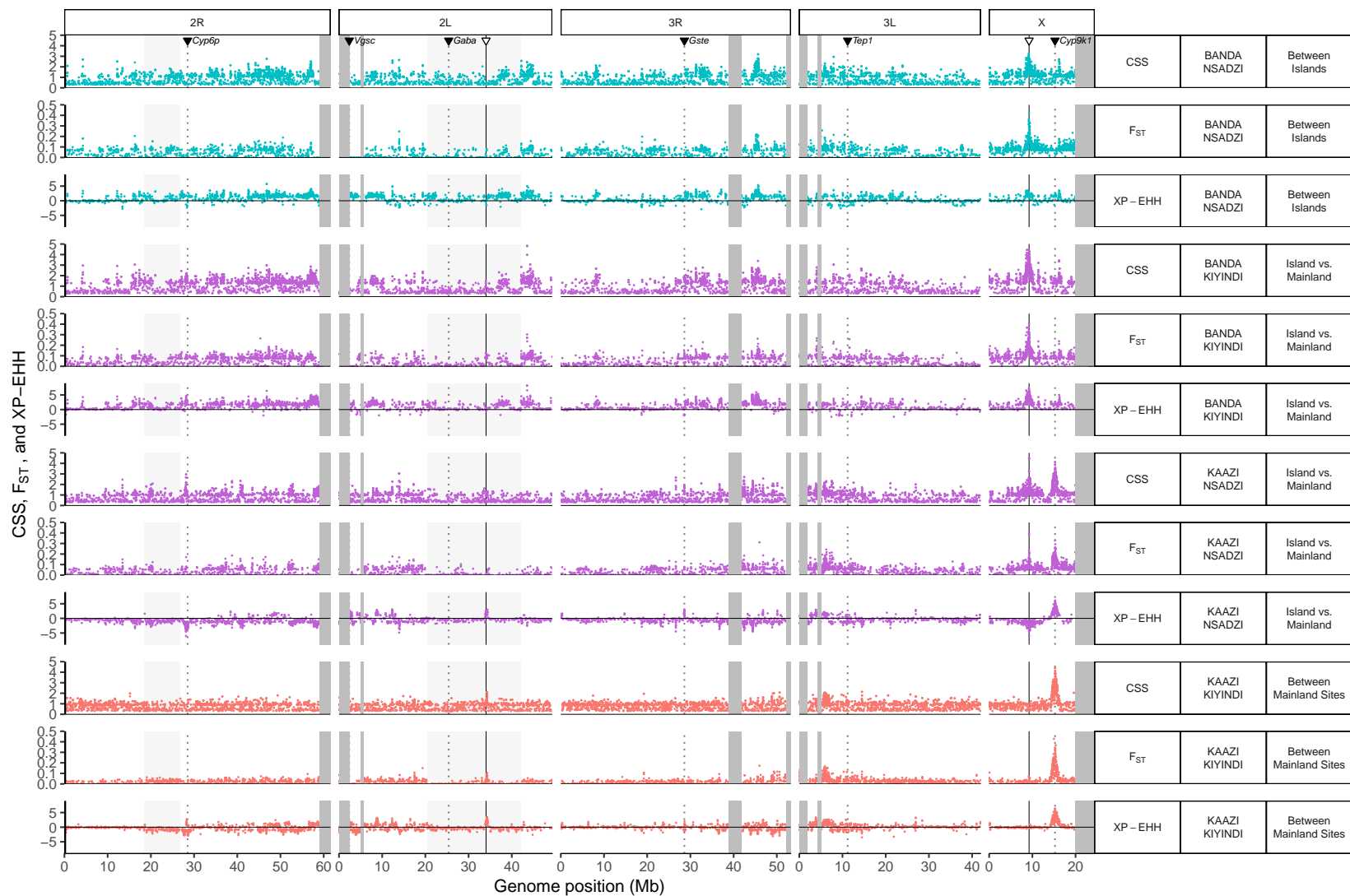
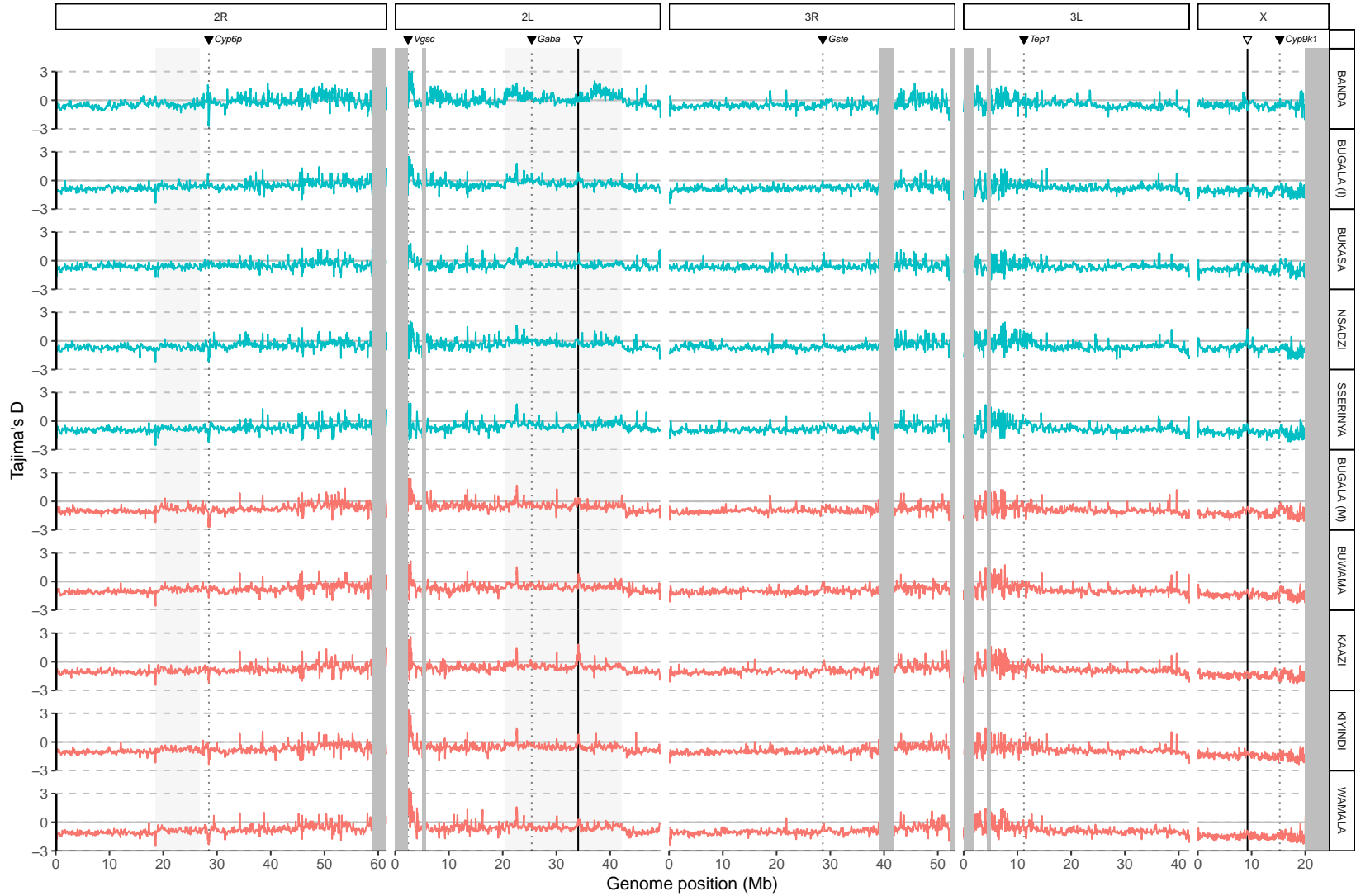


Figure S8:  $F_{ST}$  and XP-EHH across genome.

$F_{ST}$ , XP-EHH, and CSS averaged in windows of size 10 kb plotted across genome for pairwise comparisons of island and mainland localities. Shaded regions indicate inversions or heterochromatic regions (excluded from analysis) and dotted lines indicate known insecticide genes while dashed lines indicate the two putative sweeps identified in the present study. Only several exemplar pairs of populations shown.



69

Figure S9: Tajima's  $D$  across genome.

Tajima's  $D$  plotted across genome. Shaded regions indicate inversions or heterochromatic regions (excluded from analysis) and dotted lines indicate known insecticide genes while dashed lines indicate the two putative sweeps identified in the present study.



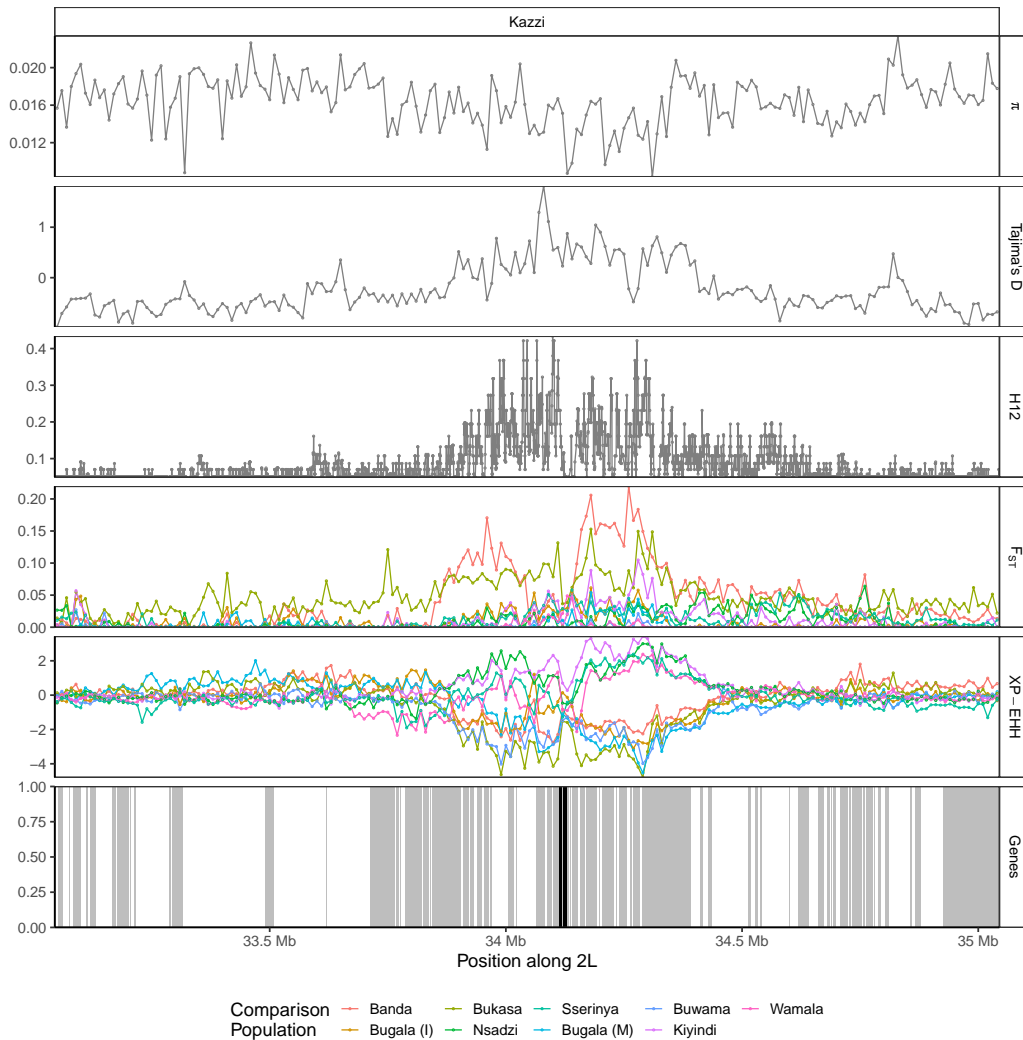


Figure S10: Selective sweep signal on chromosome 2L.

Population genetic statistics plotted near putative sweep on chromosome 2L. Focus population for all pairwise  $F_{ST}$  and XP-EHH comparisons is mainland site Kaazi. Region shown is 1 Mb up- and downstream of sweep target, centered at chr2L:34,044,820. Several genes involved in chorion formation (AGAP006549, AGAP006550, AGAP006551, AGAP006553, AGAP006554, AGAP006555 and AGAP006556) are shown in black.

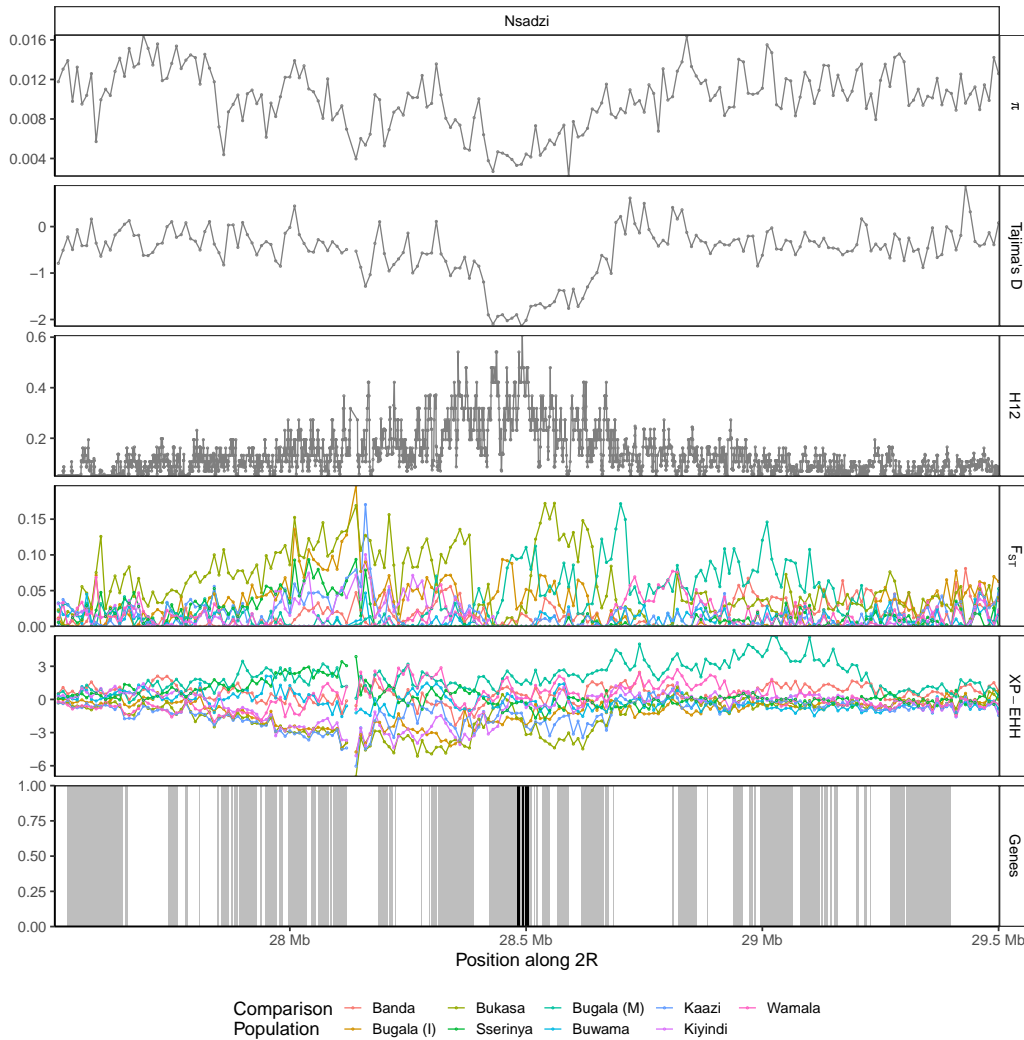


Figure S11: Selective sweep signal at *Cyp6* gene cluster on chromosome 2R. Population genetic statistics plotted near *Cyp6* gene cluster on chromosome 2R. Focus population for all pairwise  $F_{ST}$  and XP-EHH comparisons is island site Nsadzi. Region shown is 1 Mb up- and downstream of gene cluster, centered at chr2R:28,501,972, and *Cyp6* genes are highlighted in black.

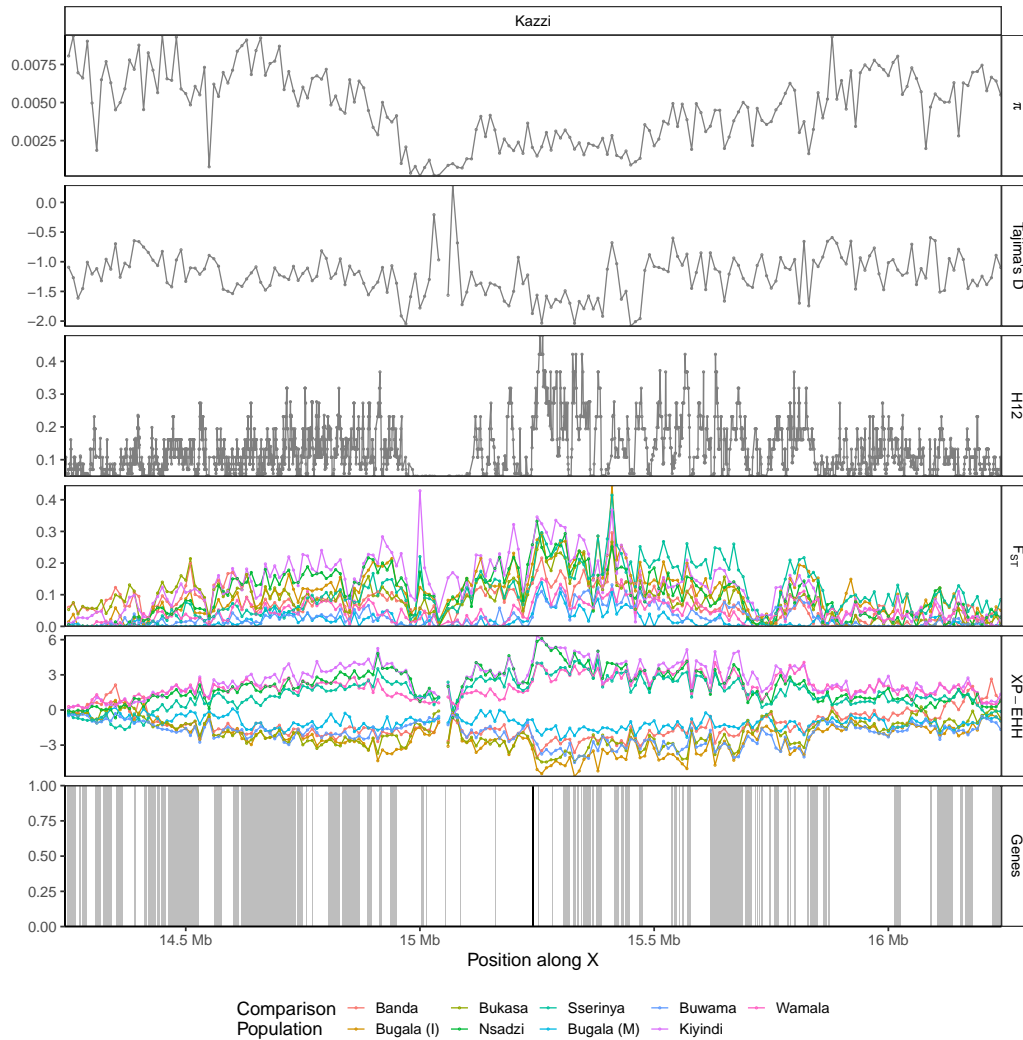


Figure S12: Selective sweep signal at gene *Cyp9K1* on X-chromosome. Population genetic statistics plotted near gene *Cyp9K1* on X-chromosome. Focus population for all pairwise  $F_{ST}$  and XP-EHH comparisons is mainland site Kaazi. Region shown is 1 Mb up- and downstream of gene cluster, centered at chrX:15,241,718, and *Cyp9K1* gene is highlighted in black.

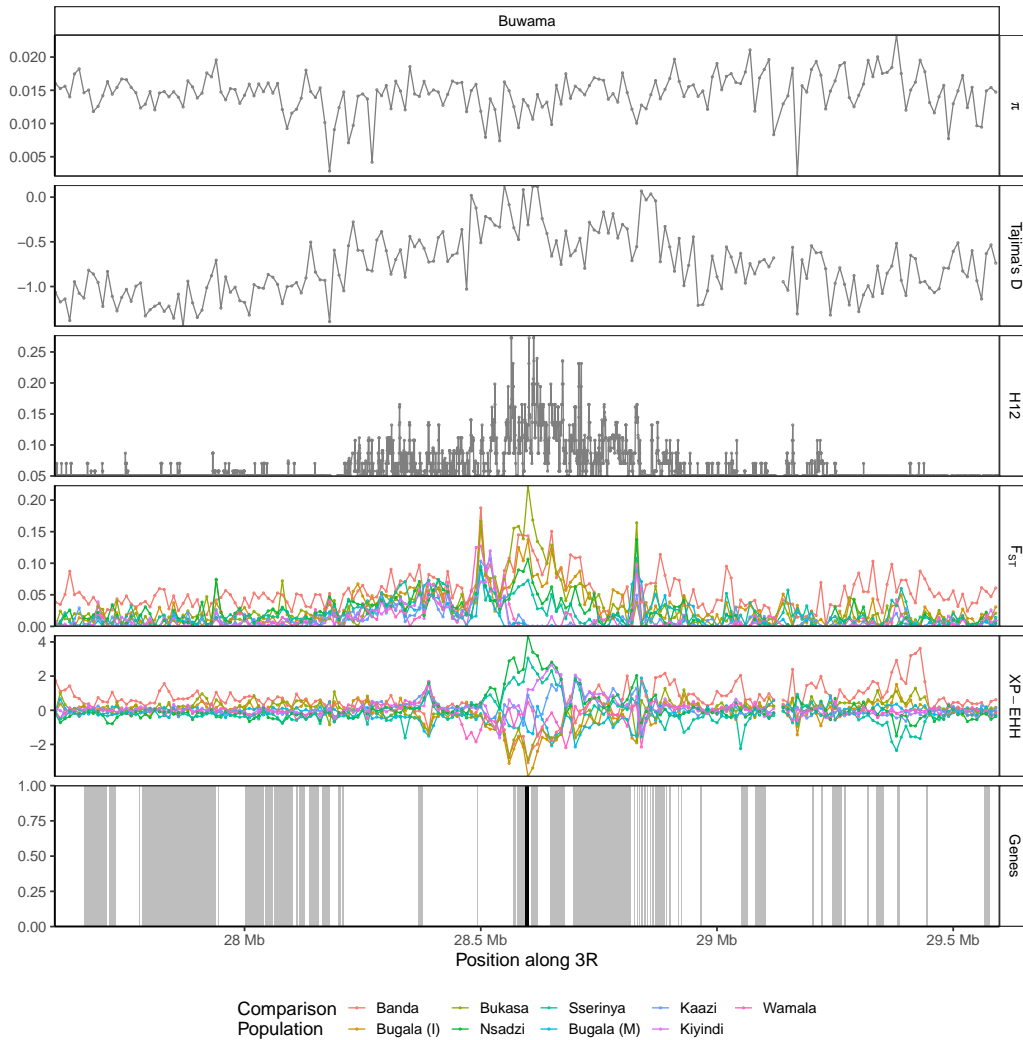


Figure S13: Selective sweep signal at *Gste* gene cluster on chromosome 3R. Population genetic statistics plotted near *Gste* gene cluster on chromosome 3R. Focus population for all pairwise  $F_{ST}$  and XP-EHH comparisons is mainland site Buwama. Region shown is 1 Mb up- and downstream of gene cluster, centered at chr3R:28,598,038, and *Gste* genes are highlighted in black.

74

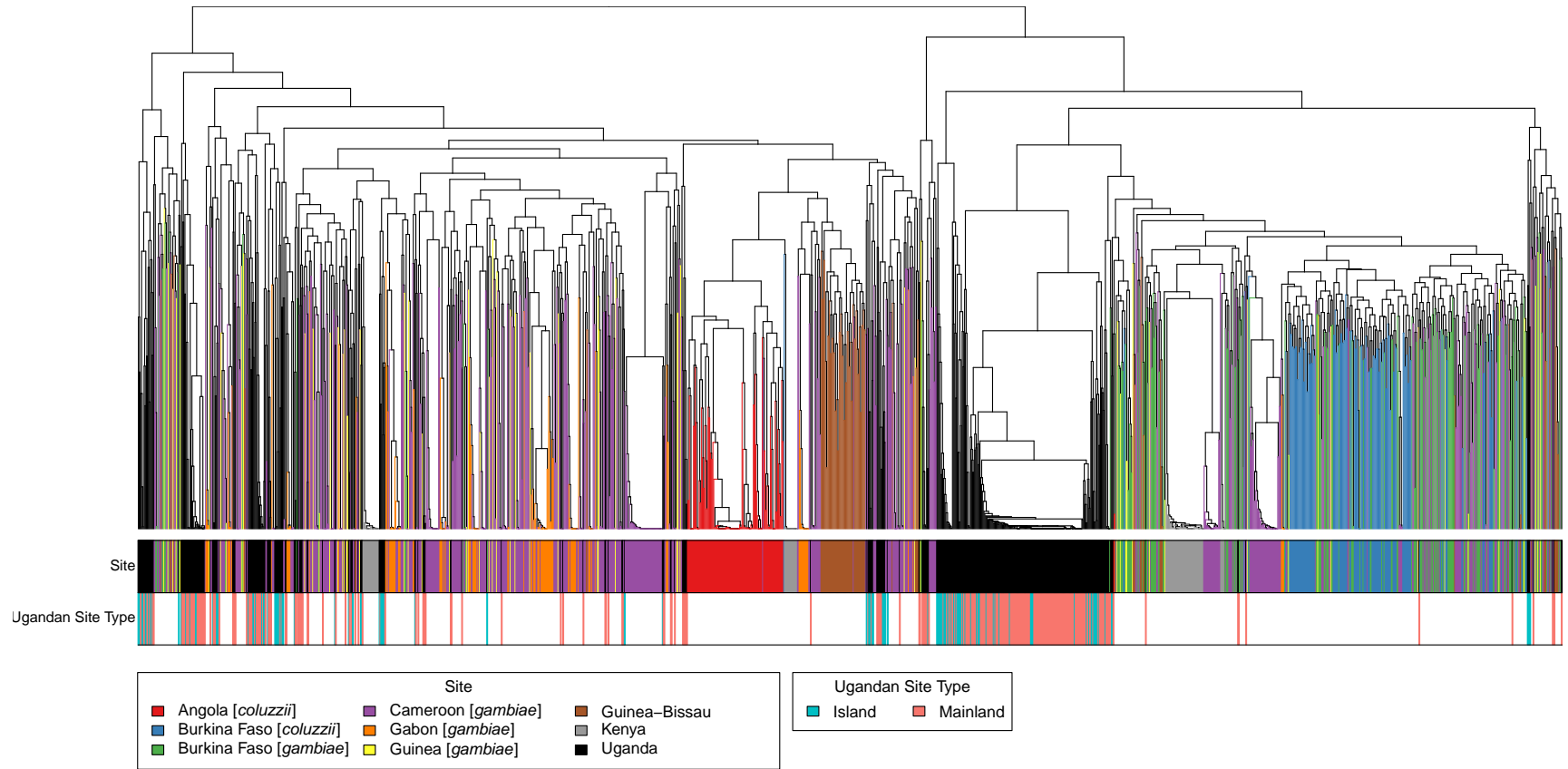


Figure S14: Tree for putative sweep on chromosome 2L.

Distance-based tree of haplotypes near putative sweep on chromosome 2L. Region shown is 100 kb up- and downstream of sweep target, centered at chr2L:34,044,820. Top color bar indicates locality, with all Ugandan individuals, from both the Ag1000G reference population and the LVB, in black. The bottom color bar differentiates the Ugandan individuals into mainland (red) and island (blue) individuals.

75

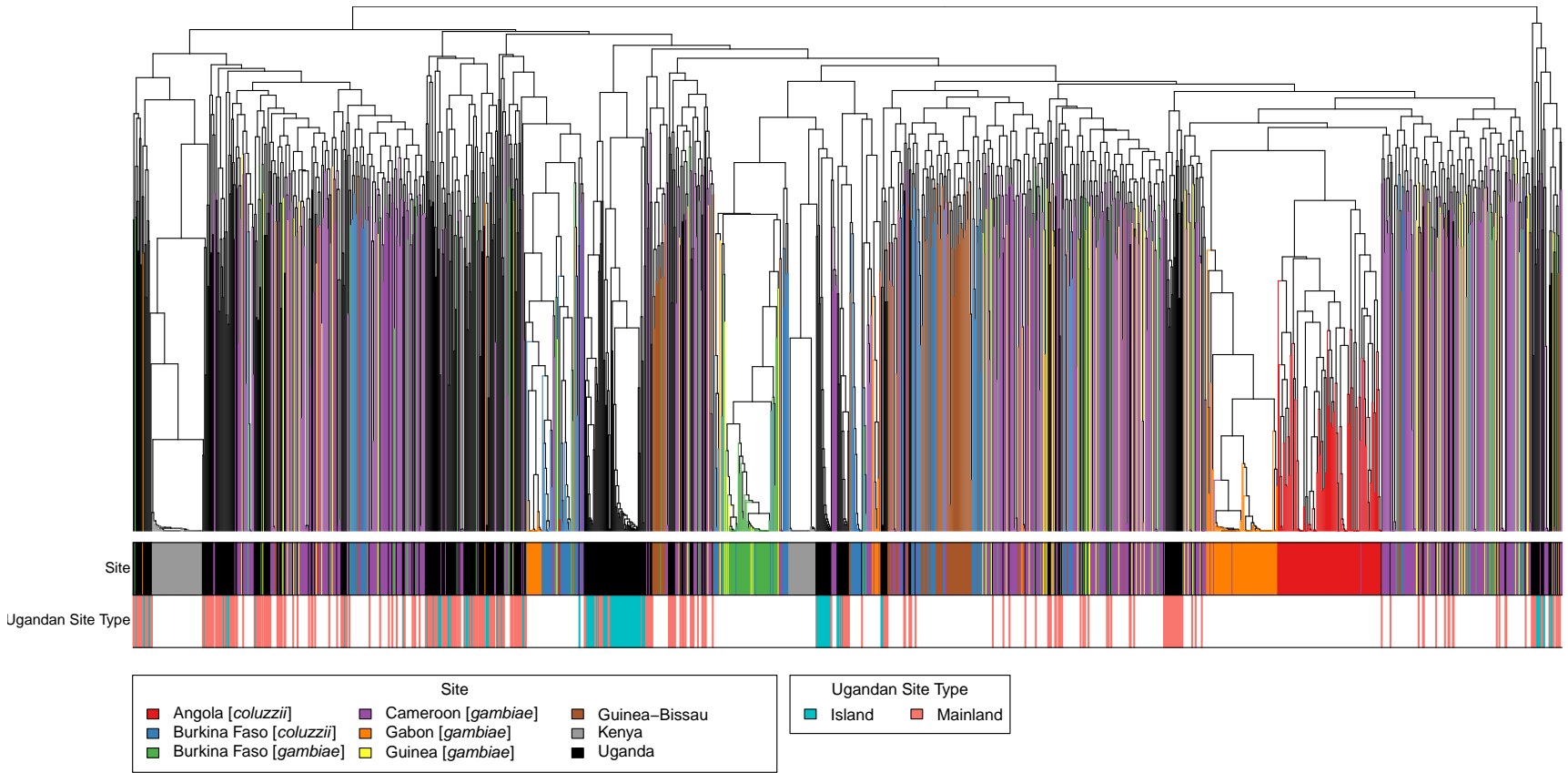


Figure S15: Tree for putative sweep on X-chromosome near *rdgA* ortholog.

Distance-based tree of haplotypes near putative sweep on X-chromosome. Region shown is 100 kb up- and downstream of sweep target, centered at chrX:9,238,942. Top color bar indicates locality, with all Ugandan individuals, from both the Ag1000G reference population and the LVB, in black. The bottom color bar differentiates the Ugandan individuals into mainland (red) and island (blue) individuals.

76

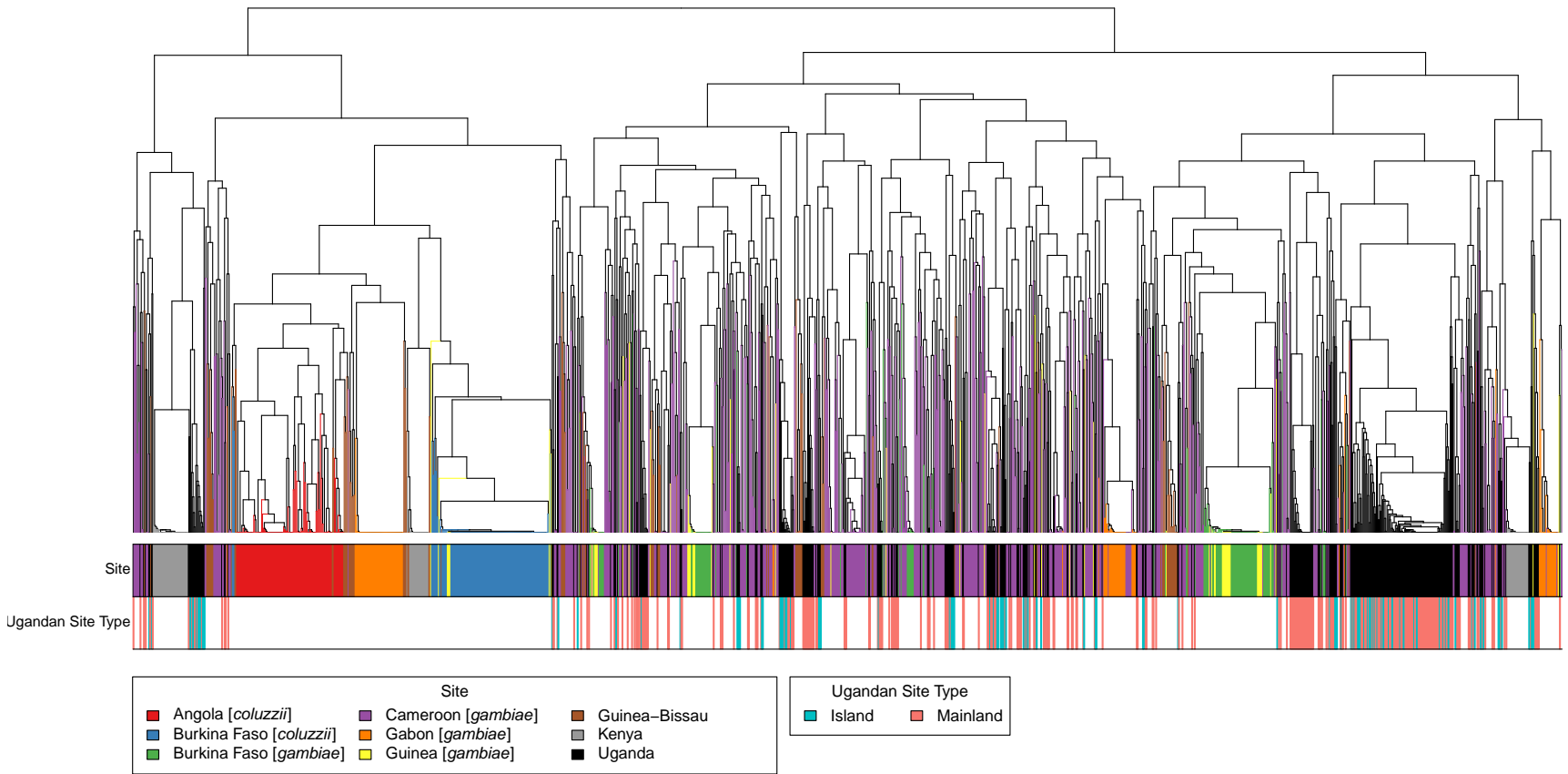


Figure S16: Tree for sweep at gene *Cyp9K1* on X-chromosome.

Distance-based tree of haplotypes near sweep at gene *Cyp9K1* on X-chromosome. Region shown is 100 kb up- and downstream of sweep target, centered at chrX:15,241,718. (Insufficient variants preclude inferring tree for region of width 20 kb.) Top color bar indicates locality, with all Ugandan individuals, from both the Ag1000G reference population and the LVB, in black. The bottom color bar differentiates the Ugandan individuals into mainland (red) and island (blue) individuals.

77

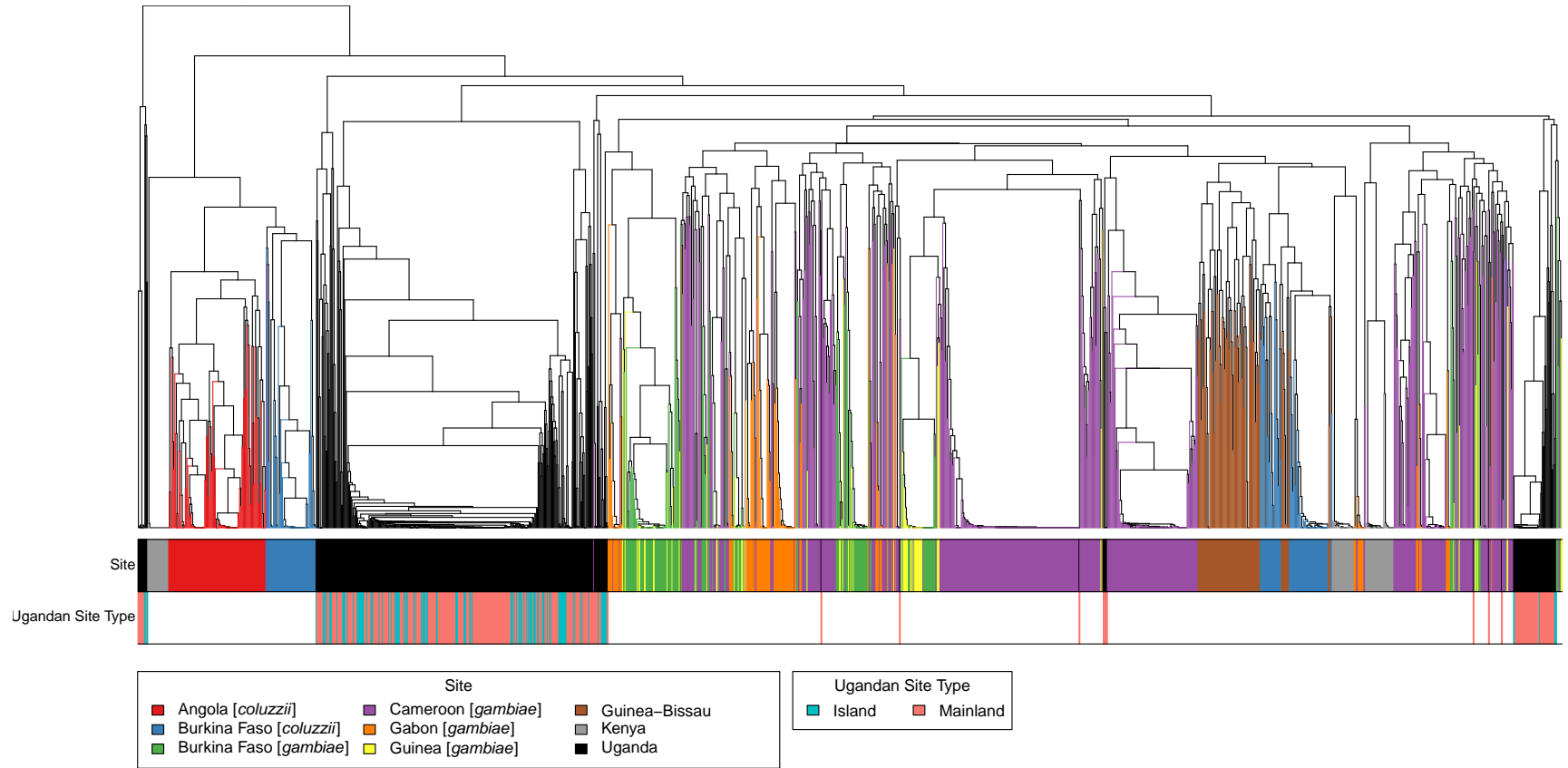


Figure S17: Tree for sweep on *Cyp6* gene cluster on chromosome 2R.

Distance-based tree of haplotypes near sweep at *Cyp6* gene cluster on chromosome 2R. Region shown is 100 kb up- and downstream of sweep target, centered at chr2R:28,501,972. Top color bar indicates locality, with all Ugandan individuals, from both the Ag1000G reference population and the LVB, in black. The bottom color bar differentiates the Ugandan individuals into mainland (red) and island (blue) individuals.



78

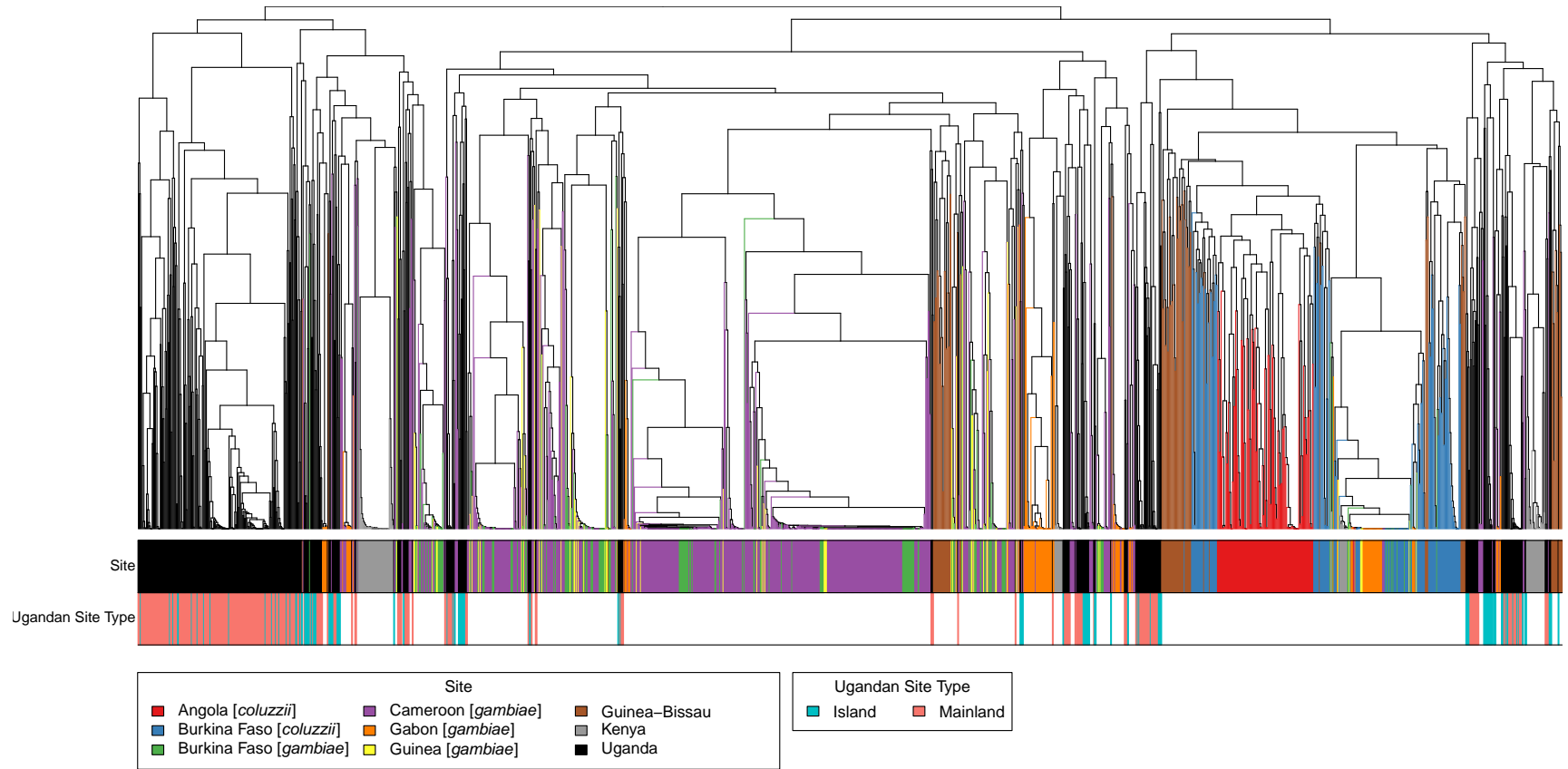


Figure S18: Tree for sweep on *Cyp6* gene cluster on chromosome 3R.

Distance-based tree of haplotypes near sweep at *Cyp6* gene cluster on chromosome 3R. Region shown is 100 kb up- and downstream of sweep target, centered at chr3R:28,598,038. Top color bar indicates locality, with all Ugandan individuals, from both the Ag1000G reference population and the LVB, in black. The bottom color bar differentiates the Ugandan individuals into mainland (red) and island (blue) individuals.

## References

- [94] Sharakhova, M. V. *et al.* Update of the *Anopheles gambiae* PEST genome assembly. *Genome Biology* **8**, R5 (2007).
- [95] Aronesty, E. ea-utils: Command-line tools for processing biological sequencing data (2011).
- [96] Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
- [97] DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics* **43**, 491–498 (2011).
- [98] Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *American Journal of Human Genetics* **81**, 559–75 (2007).
- [99] Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* **4**, 7 (2015).
- [100] Delaneau, O., Howie, B., Cox, A. J., Zagury, J. F. & Marchini, J. Haplotype estimation using sequencing reads. *American Journal of Human Genetics* **93**, 687–696 (2013).
- [101] Li, H. Tabix: Fast retrieval of sequence features from generic TAB-delimited files. *Bioinformatics* **27**, 718–719 (2011).
- [102] Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–9 (2009).
- [103] Alexander, D., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Research* **19**, 1655–1664 (2009).

- 715 [104] Kopelman, N. M., Mayzel, J., Jakobsson, M., Rosenberg, N. A. & Mayrose, I.  
716 CLUMPAK: a program for identifying clustering modes and packaging population struc-  
717 ture inferences across K. *Molecular Ecology Resources* **15**, 1179–1191 (2015).
- 718 [105] Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–8  
719 (2011).
- 720 [106] Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H. & Bustamante, C. D. Infer-  
721 ring the joint demographic history of multiple populations from multidimensional SNP  
722 frequency data. *PLoS Genetics* **5**, e1000695 (2009).
- 723 [107] Coffman, A. J., Hsieh, P. H., Gravel, S. & Gutenkunst, R. N. Computationally effi-  
724 cient composite likelihood statistics for demographic inference. *Molecular Biology and*  
725 *Evolution* **33**, 591–593 (2016).
- 726 [108] Liu, X. & Fu, Y.-X. Exploring population size changes using SNP frequency spectra.  
727 *Nature Genetics* **47**, 555–559 (2015).
- 728 [109] Szpiech, Z. A. & Hernandez, R. D. selscan: an efficient multithreaded program to  
729 perform EHH-based scans for positive selection. *Molecular Biology and Evolution* **31**,  
730 2824–2827 (2014).
- 731 [110] Jombart, T. & Ahmed, I. adegenet 1.3-1: new tools for the analysis of genome-wide  
732 SNP data. *Bioinformatics* (2011).
- 733 [111] Paradis, E., Claude, J. & Strimmer, K. A{PE}: analyses of phylogenetics and evolution  
734 in {R} language. *Bioinformatics* **20**, 289–290 (2004).
- 735 [112] Neuwirth, E. *RColorBrewer: ColorBrewer Palettes* (2014).
- 736 [113] Galili, T. dendextend: an R package for visualizing, adjusting, and comparing trees of  
737 hierarchical clustering. *Bioinformatics* (2015).

- 738 [114] Gautier, M. & Vitalis, R. rehh: An R package to detect footprints of selection in  
739 genome-wide SNP data from haplotype structure. *Bioinformatics* **28**, 1176–1177 (2012).
- 740 [115] Patterson, N., Price, A. L. & Reich, D. Population structure and eigenanalysis. *PLoS*  
741 *Genetics* **2**, 2074–2093 (2006).
- 742 [116] Price, A. *et al.* Principal components analysis corrects for stratification in genome-wide  
743 association studies. *Nature Genetics* **38**, 904–9 (2006).
- 744 [117] Tange, O. GNU Parallel - The Command-Line Power Tool. *login: The USENIX*  
745 *Magazine* **36**, 42–47 (2011).
- 746 [118] Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing  
747 genomic features. *Bioinformatics* **26**, 841–2 (2010).