1   # A chromosome-scale assembly of the major African

2   # malaria vector *Anopheles funestus*

3   Jay Ghurye[1,2], Sergey Koren[2], Scott T. Small[3], Seth Redmond[4,5], Paul Howell[6,†], Adam M.

4   Phillippy[2,*], and Nora J. Besansky[3,*]

5

6   [1] Department of Computer Science, University of Maryland, College Park, MD

7   [2] Genome Informatics Section, Computational and Statistical Genomics Branch, National

8   Human Genome Research Institute, National Institute of Health, Bethesda, MD

9   [3] Department of Biological Sciences, University of Notre Dame, South Bend, IN

10  [4] Infectious Disease and Microbiome Program, Broad Institute, Cambridge, MA

11  [5] Department of Immunology and Infectious Disease, Harvard TH Chan School of Public Health,

12  Boston, MA

13  [6] Centers for Disease Control, Atlanta, GA

14  [†] Current affiliation: Verily Life Sciences, San Francisco, CA

15  [*] Corresponding Authors

16

17

18

19

20    # Abstract

21

22    **Background:** *Anopheles funestus* is one of the three most consequential and widespread

23    vectors of human malaria in tropical Africa. However, the lack of a high-quality reference

24    genome has hindered the association of phenotypic traits with their genetic basis in this

25    important mosquito.

26

27    **Findings:** Here we present a new high-quality *An. funestus* reference genome (AfunF3)

28    assembled using 240x coverage of long-read single-molecule sequencing for contigging,

29    combined with 100x coverage of short-read Hi-C data for chromosome scaffolding. The

30    assembled contigs total 446 Mbp of sequence and contain substantial duplication due to

31    alternative alleles present in the sequenced pool of mosquitos from the FUMOZ colony. Using

32    alignment and depth-of-coverage information, these contigs were deduplicated to a 211 Mbp

33    primary assembly, which is closer to the expected haploid genome size of 250 Mbp. This

34    primary assembly consists of 1,053 contigs organized into 3 chromosome-scale scaffolds with

35    an N50 contig size of 632 kbp and an N50 scaffold size of 93.811 Mbp, representing a 100-fold

36    improvement in continuity versus the current reference assembly, AfunF1.

37

38    **Conclusion:** This highly contiguous and complete *An. funestus* reference genome assembly

39    will serve as an improved basis for future studies of genomic variation and organization in this

40    important disease vector.

41

42 # Data Description

43 ## Introduction and Background

44

45 Many insect genomes remain a challenge to assemble, and mosquito genomes have proven

46 particularly difficult due to their repeat content and structurally dynamic genomes. These issues

47 are compounded by the requirements of long-read sequencing technologies that typically

48 require >10 µg of DNA for library construction. As a result, it is often impossible to construct a

49 sequencing library from a single individual. Instead, sequencing a pool of individuals from an

50 inbred population has been required [1]. For species that are amenable to extensive inbreeding,

51 this approach has led to reference-grade genomes directly from the assembler [2]. However,

52 when inbreeding is not possible, the sequenced pool of individuals can carry population

53 variation that fragments the resulting assembly. In this case, instead of assembling a single

54 genome, the assembler must reconstruct some unknown number of variant haplotypes.

55

56 Motivated by the goal of genome-enabled malaria control, a large international consortium

57 previously sequenced and assembled the genomes of 16 *Anopheles* species using short-read

58 Illumina sequencing [3,4]. Although these draft assemblies represented a crucial first step, their

59 potential for 1) understanding and manipulating vectorial capacity traits, 2) inferring how key

60 vector adaptations to hosts and habitats have arisen and are maintained, and 3) accurately

61 defining vector breeding units and migration between them is constrained by two major

62 limitations. First, many of these *Anopheles* assemblies are highly fragmented collections of

63 relatively short scaffolds, causing gene annotation problems such as missing genes, missing

64 exons, and genes split between scaffolds or sequencing gaps. Thus, one of the consequences

65 of fragmented assemblies is that it is difficult to estimate gene copy number, which may be

66    linked to important phenotypic traits (e.g. insecticide resistance) [5,6]. Genes of particular

67    interest with respect to arthropod disease vectors (e.g., cytochrome P450s and

68    odorant/gustatory receptors) may be especially prone to annotation errors, as many belong to

69    gene families whose members are often physically clustered into tandem arrays.

70

71    A second major limitation of fragmented insect assemblies is that they are rarely scaffolded into

72    chromosomes, owing to difficulty and lack of funding for physical or linkage mapping. Among

73    other consequences, the unknown placement of scaffolds along chromosome arms means that

74    their position within or outside of chromosomal inversions is difficult or impossible to determine.

75    Many anopheline species are highly polymorphic for chromosomal inversions, which tend to

76    occur disproportionately on particular chromosome arms [7–9]. In a heterozygote carrying one

77    inverted and one uninverted chromosome, recombination between the reversed chromosomal

78    segments is greatly reduced [10], creating cryptic population structure that can cause spurious

79    associations in GWAS [11] and mislead recombination-based inference of selection and gene

80    flow [12,13].  Importantly, chromosomal inversions also directly or indirectly influence traits

81    affecting malaria transmission intensity—anopheline biting and resting behavior [14,15],

82    seasonality [16], aridity tolerance [14,17–21], ecological plasticity [22,23] morphometric variation

83    [24], and *Plasmodium* infection rates [25,26]. Thus, correct population genomic and GWAS

84    inferences depend upon knowing the location of a marker in the genome.

85

86    *Anopheles funestus* is one of the three most important and widespread vectors of human

87    malaria in tropical Africa [27–30], and unlike *Anopheles gambiae* with which it broadly co-

88    occurs, it is a relatively neglected species. It is considered even more highly anthropophilic and

89    endophilic than *An. gambiae* and amenable to conventional indoor-based vector control such as

90    bed nets and indoor spraying of houses with residual insecticides. Indeed, historical house

91    spraying campaigns in eastern and southern Africa not only locally eliminated this species, but

4

92    the effect was maintained for several years following the cessation of spraying, due to the

93    apparent inability of *An. funestus* to recolonize some areas. Likewise, *An. funestus* was

94    eliminated from a humid forest and degraded forest areas in West Africa where malaria is meso-

95    or hypoendemic [31]. However, in the savanna environment of West Africa where malaria is

96    holo- or hyperendemic, similar historical indoor spraying campaigns failed to eliminate the

97    species. Exophilic populations persisted which—despite marked anthropophily—continued to

98    feed outdoors on cattle but also entered sprayed houses to bite humans. Today, the situation is

99    worsened by the emergence and spread of insecticide resistance in this species [29,32–34].

100

101    Mastery over malaria will require tackling *An. funestus*, but it remains understudied; information

102    on its behavior and genetics lags far behind *An. gambiae.* At least part of the reason for its

103    neglect may be the historical lack of laboratory colonies, a problem solved with the

104    establishment of the FUMOZ colony and its registration with the Anopheles program of BEI

105    Resources (https://www.beiresources.org/AnophelesProgram.aspx*). An. funestus* shares with

106    *An. gambiae* not only a broad sub-Saharan distribution and major vector status but also

107    abundant chromosomal inversion polymorphism and shallow range-wide population structure

108    [35]. However, there are behavioral and genetic heterogeneities relevant to malaria transmission

109    that remain poorly understood. In West Africa, strong cytogenetic evidence points to cryptic,

110    temporally stable assortatively mating populations co-occurring in the same villages [36–39].

111    These chromosomally recognized forms of *An. funestus*, named Kiribina and Folonzo, seem to

112    differ in larval ecology and—importantly—they also differ in adult behaviors affecting vectorial

113    capacity, most notably indoor resting behavior. Mechanistic understanding of the genomic

114    determinants of these and other epidemiologically important phenotypic and behavioral traits

115    ultimately depends on upgrading the *An. funestus* reference to a chromosome-based assembly

116    in which the unanchored scaffolds are united, ordered and oriented on chromosome arms.

5

## Chromosome-scale assembly of *Anopheles funestus*

To achieve a complete and highly contiguous assembly of the *An. funestus* genome (AfunF3), we first assembled contigs from long, single-molecule reads, and then scaffolded these contigs into chromosome-scale scaffolds using Hi-C proximity ligation data. A similar strategy was recently used to improve the genome of *Aedes aegypti* [40]. An initial assembly of the long-read data alone (AfunF3 contigs) yielded a contig N50 size of 94.05 kbp (N50 such that 50% of assembled bases are in contigs of this size or greater) and extensive haplotype separation as evidenced by an inflated assembly size of 446.04 Mbp and a high rate of core gene duplications (48%) as measured by BUSCO [41]. These alternative alleles likely derive from natural variation circulating within the sequenced FUMOZ colony, as the DNA from a pool of adult mosquitoes was required for PacBio library preparation. Identifying and removing duplicate contigs via an all-vs-all alignment reduced the primary assembly size to 211.75 Mbp and improved the N50 size to 631.72 kbp (Table 1).

The primary set of contigs (excluding alternative alleles) was then scaffolded using Hi-C Illumina reads to first bin the contigs into 3 chromosomes, followed by ordering and orientation of the contigs using the Proximo method (Phase Genomics, Seattle WA). The final scaffolded assembly (AfunF3 primary) contains 210.82 Mbp of sequence and a scaffold N50 of 93.81 Mbp. The resulting scaffolds represent the entirety of the three *An. funestus* chromosomes: 2, 3, and X (Figure 1).

Because single-molecule PacBio data is prone to insertion and deletion errors, all AfunF3 contigs were polished twice with Arrow [42] using the signal-level PacBio data and once with Pilon [43] using paired-end Illumina data from the same FUMOZ colony. Because Illumina-based polishing tools typically do not correct bases that appear heterozygous in the read set,

142    we anticipated that variation in the FUMOZ colony would prevent the correction of variant

143    bases. To help address this issue, we finally polished the assembly using 10X Genomics

144    Illumina data obtained from an individual mosquito. As an independent test of base accuracy,

145    we compared our new assembly (AfunF3 primary) and the prior assembly (AfunF1) to a 10X

146    Genomics dataset from a different individual mosquito. The average Phred-scaled quality value

147    [44] of the new assembly was estimated as QV28 versus QV23 for the Illumina-based AfunF1

148    assembly. This independent data indicates a higher average accuracy for the new assembly,

149    but also revealed significant diversity within the colony. For example, calling variants using 10X

150    Genomics data for two different mosquitos yielded widely different SNP counts (92,759 vs.

151    177,428).

152

153    We next evaluated the structural accuracy of the AfunF1 and AfunF3 assemblies by measuring

154    their agreement with the raw PacBio reads. The intermediate assembly AfunF2 [45] was

155    assembled before collection of all PacBio and Hi-C data, and so was deemed redundant and

156    excluded from these analyses. When compared to the raw data, the AfunF3 primary assembly

157    had fewer called structural differences (insertions, deletions, duplications, and inversions) than

158    AfunF1 (Table 2). Despite the substantial single-nucleotide polymorphism observed within the

159    FUMOZ colony, no large polymorphic inversions could be identified from the combined PacBio,

160    Hi-C, and 10X Genomics data. Comparison of the chromosome-scale AfunF3 primary assembly

161    versus the An. gambiae reference genome (AgamP4) confirmed a known reciprocal whole-arm

162    translocation between 2L and 3R, as well as substantial intra-chromosomal shuffling (Figure 2).

163    AfunF3 contigs also had fewer fragmented BUSCO core genes and a similar number of

164    complete BUSCOs compared to AfunF1 (Table 2), but also a high rate of duplication. The

165    AfunF3 primary scaffolds reduce duplication at the expense of lower BUSCO completeness.

166

167   To further evaluate AfunF3's suitability as an updated reference for *An. funestus*, we mapped

168   RNA-Seq expression data to the assemblies and computed the number of concordant paired-

169   end reads. A better assembly is expected to have both a higher fraction of mapped reads

170   (completeness) as well as a higher fraction of correctly spaced and oriented pairs (structural

171   accuracy). Both AfunF3 assemblies have better agreement of mapped read pairs as well as a

172   higher overall mapping rate versus the AfunF1 assembly (Table 2). The AfunF3 contigs do have

173   a higher rate of multi-mapping RNA-Seq reads, but this is reduced in the primary assembly

174   while preserving the high mapping rate. In addition to a higher mapping rate, more complete

175   transcripts were mapped to single contigs within the long-read assemblies. The average number

176   of complete transcripts contained per contig was 67.38 for AfunF3 primary versus 5.28 for the

177   AfunF1 assembly. These results demonstrate the greater continuity of the updated assembly,

178   which provides sequence-resolved reconstructions of many *An. funestus* intergenic regions for

179   the first time.


180   ## Discussion


181   *Anopheles funestus* is one of the leading vectors of malaria and understanding the organization

182   and function of its genome is key to controlling this deadly disease. Here we described a

183   chromosome-scale assembly of the *An. funestus* genome using multiple sequencing

184   technologies and assembly methods. The tremendous improvement in the completeness and

185   contiguity of its genome will provide a valuable resource for future genomic analyses and

186   functional characterization of this important species and enable a mechanistic understanding of

187   the genomic determinants of epidemiologically important phenotypic and behavioral traits.

188

# 189  Materials and Methods

## 190  Library preparation and sequencing

191  A gravid female mosquito of the FUMOZ colony was allowed to lay eggs, and her offspring were

192  inbred for a single generation. From this, an isofemale line was grown and DNA extracted from

193  the adult females for sequencing with PacBio and Hi-C. 46 SMRT cells of PacBio RSII

194  sequencing using the P6-C4 chemistry were run by the core facility at the Icahn School of

195  Medicine at Mount Sinai (New York, NY), resulting in 173X coverage (assuming a 250 Mbp

196  genome size). A previous study generated 70X coverage of the same colony using the older

197  PacBio P5-C3 chemistry sequencing [45]. This older data was combined with the additional

198  173X coverage, totaling 60.95 Gb of long-read data in 10.93 million sequences (average length

199  5.6 kb, N50 read length 8.4 kb) and an estimated total coverage of 234X. Two Hi-C libraries

200  were prepared and sequenced (one from mixed-sex larvae, the second from adult females) by

201  Phase Genomics (Seattle, WA), resulting in ~100X coverage of Illumina Hi-C data containing

202  ~187 million 80 bp paired-end Illumina reads.

## 203  Assembly and scaffolding

204  PacBio contig assembly was performed with Canu v1.3 [46] using parameters:

205  corOutCoverage=100 genomeSize=250m errorRate=0.013 batOptions="-dg 3 -db 3 -dr 1 -ca

206  500 -cp 50". The resulting contigs were then polished with Arrow [42] using default parameters

207  and the P6-C4 PacBio signal data (because Arrow does not support the older P5-C3 data). After

208  polishing, the assembly was separated into primary and alternative contigs to remove

209  unnecessarily duplicated alleles from the AfunF3 contigs. This was performed using two

210  different approaches. First, contigs containing at least one complete BUSCO gene were

211  identified. For each BUSCO gene, if it was found contained in two or more contigs, the contig

9

212     with the highest alignment score was kept as the primary. Next, all contigs not containing a

213     BUSCO gene but assembled with high coverage (>40X) were added to the primary set.

214

215     To order and orient the primary contigs along the chromosomes, Hi-C reads were aligned using

216     Bowtie2 [47] and scaffolding using Proximo (Phase Genomics, Seattle WA). Scaffold gaps

217     spanned by PacBio reads were filled using PBJelly [48]. This assembly was again run through

218     Arrow to polish the sequences inserted by PBJelly and fill any remaining short gaps. The Hi-C

219     assembled scaffolds were then aligned using NUCmer [49] to the AfunF1 contigs for validation

220     and the alignments visualized using Circos [50] and mummerplot. This identified a mis-join of

221     chromosomes 3R and X, which was manually corrected. Additional manual curation using

222     mapped transcripts, FISH probes [45] , and comparison to AfunF1 scaffolds identified a few

223     additional inversion errors in the scaffolds, mainly on distal 2L. Visual inspection of the Hi-C

224     data showed clear signatures of scaffolding error. These errors were corrected by manually

225     extracting the region and placing the sequence at the correct locus, as indicated by the Hi-C

226     interactions. After these corrections, the scaffolded chromosomes (AfunF3 primary) show good

227     agreement with the Hi-C data (Figure 3).

228

229     As diploid and population variation introduces indels in the Arrow polishing process [51], the

230     final assemblies were also polished by Pilon using paired-end Illumina data (NCBI SRA

231     accession numbers: SRX209628 and SRX209387) and 10X Genomics Illumina data from a

232     single individual (NCBI SRA accession number: SRX4819916). The paired-end Illumina data

233     was mapped using BWA-MEM [52] and the 10X Genomics data mapped using Lariat [53] in a

234     barcode-aware manner, as to improve the mapping quality. Consensus quality of the final

235     assemblies was then estimated using an independent 10X Genomics dataset (NCBI SRA

236     accession number: SRX4819903) of a different mosquito of the same FUMOZ colony. Based on

237     the alignment of reads to the assembly, variants were called using freebayes (parameters: -C 2

10

238   -0 -O -q 20 -z 0.10 -E 0 -X -u -p 2 -F 0.5), and the assembly QV was estimated using called

239   homozygous variants (i.e. positions where nearly all Illumina reads agreed with each other yet

240   disagreed with the assembly).

## Validation

242   To check for the presence of contamination, assembled contigs were classified using Kraken

243   [54] using a custom database including all microbial RefSeq genomes and all available

244   mosquito genomes. Most of the assembled sequence (96.00%) was classified as *An. funestus*

245   or Culicidae. The remaining sequences were primarily unannotated or annotated at a higher

246   taxonomic level (3.76%), from possible bacterial/human sources (0.24%, 32 contigs), and had

247   slightly lower GC content (Figure 4). However, none of these contigs were called contaminants

248   by NCBI's independent contamination check and so all contigs were included in the submitted

249   assembly to avoid excluding novel mosquito sequence missing from the prior draft assemblies.

250

251   The structural accuracy of the assemblies was evaluated by mapping raw PacBio reads and

252   calling structural variants. PacBio reads were aligned to each assembly using NGMLR [55] with

253   parameters: -t 16 -x pacbio --skip-write. Using these alignments, variants were called using

254   Sniffles [55] with parameters: -t 32 -s 10 -f 0.25. Variants were then filtered to avoid capturing

255   heterozygous population variants such that variants for which the alternate variant had ≥45

256   supporting reads and the assembly variant had <10 supporting reads were called as assembly

257   errors.

258

259   Paired-end RNA-Seq for the *An. funestus* FUMOZ colony were downloaded from NCBI under

260   accession SRR826832. These reads were aligned to all assemblies using the HISAT2 aligner

261   [56] and assembled into transcripts using Trinity [57] with default parameters. The assembled

262   transcripts were then mapped to all assemblies using GMAP [58]. Transcripts were required to

263    be aligned over 90% of their length to a single contig to be considered "complete" in the

264    assembly.

265

# Availability of supporting data

267    Raw genomic sequence reads are available in the NCBI Sequence Read Archive under project

268    accession PRJNA494870. This Whole Genome Shotgun project has been deposited at

269    DDBJ/ENA/GenBank under the accession RCWQ00000000. The version described in this

270    paper is version RCWQ01000000.

271

272

# Declarations

## List of abbreviations

275    BUSCO: Benchmarking Universal Single-Copy Ortholog; PacBio: Pacific Biosciences; RNA-

276    Seq: RNA-sequencing; NCBI: National Center for Biotechnology Information; SRA: Sequence

277    Read Archive

## Ethics approval and consent to participate

279    Not applicable.

## Consent for publication

281    Not applicable.

12

## Competing interests

The author(s) declare that they have no competing interests.

## Funding
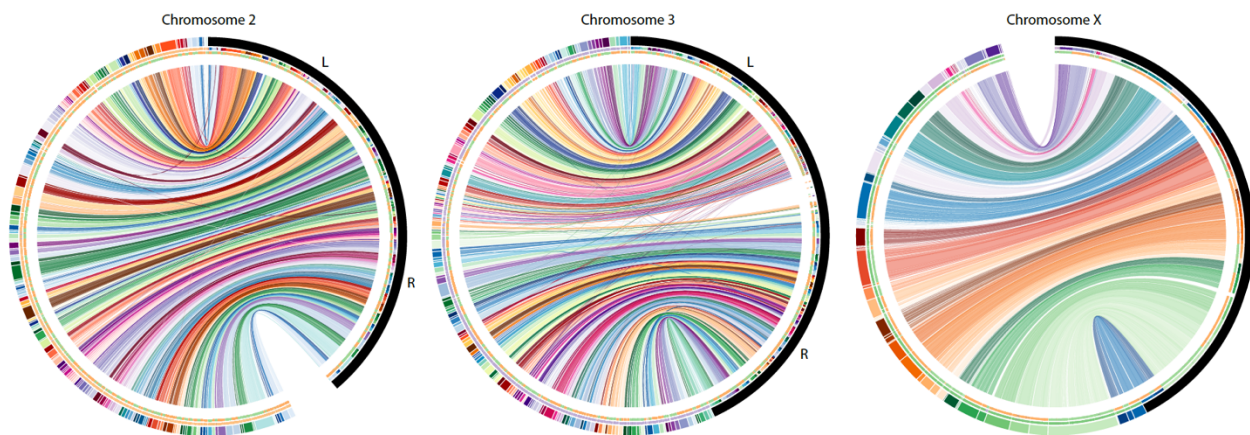
## Authors' contributions

AMP and NJB conceived and coordinated the project. JG, SK, STS, and AMP performed the genome assembly, validation, and comparative analyses. SR provided the 10X Genomics data and analysis. PH provided FUMOZ samples for sequencing. JG, AMP, and NJB drafted the manuscript. All the authors have read and approved the manuscript.

## Acknowledgments

13

# Figures

304

305     Figure 1: Circos plot comparing the AfunF1 assembly of *An. funestus* to the updated

306     AfunF3 assembly. AfunF1 scaffolds (colored half of the outer ring) are ordered by

307     majority alignment location onto AfunF3 (black half of the outer ring). Connecting lines

308     indicate pairwise alignments between the two assemblies, and crossing lines indicate

309     that part of the AfunF1 scaffold aligns to discordant regions on the AfunF3

310     chromosome. The first internal ring color correspond to the AfunF1 scaffold color. The

311     second internal ring represents the orientation of the AfunF1 scaffolds onto AfunF3,

312     where orange is forward and green is reverse.

313



314

315  Figure 2: Hi-C interaction map for assembled *An. funestus* scaffolds generated using

316  the Juicebox Hi-C visualization program [59]. Darker colors indicate a higher frequency

317  of chromatin interaction. The plot shows clear separation of chromosome boundaries

318  and limited off-diagonal interactions, supporting the global structure of the chromosome-
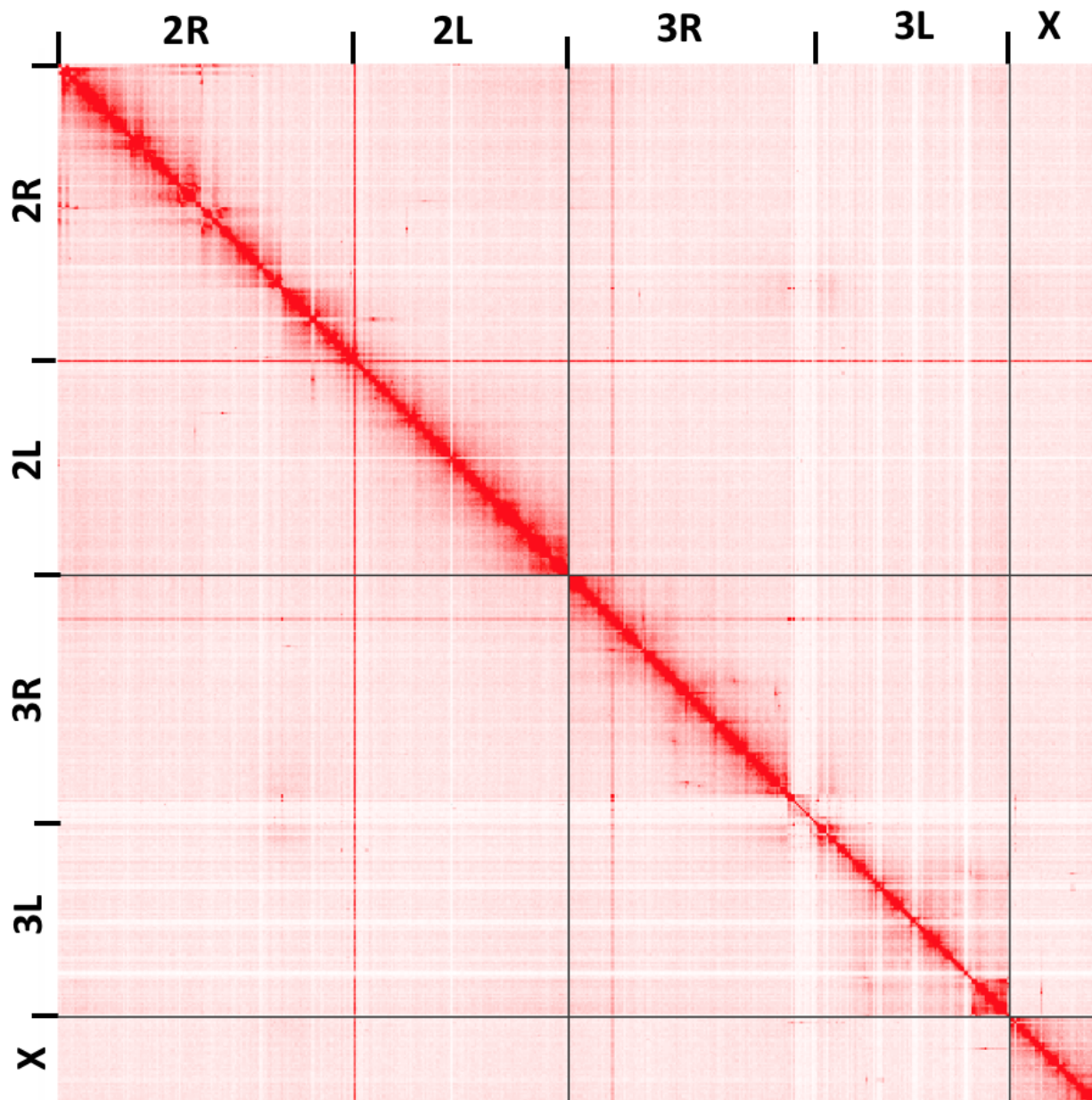
319  scale scaffolds.

320



321

15

322   Figure 3: Whole genome alignment dotplot for *Anopheles funestus* and *Anopheles*

323   *gambiae* genomes generated using D-GENIES [60]. A dot in the plot corresponds to a

324   match between the corresponding genomic positions indicated on the axes. The *An.*

325   *gambiae* reference genome is displayed on the x-axis, and the *An. funestus* AfunF3

326   primary assembly on the y-axis. A reciprocal whole-arm translocation between 2L and

327   3R is apparent, as well as substantial intra-chromosomal shuffling between these
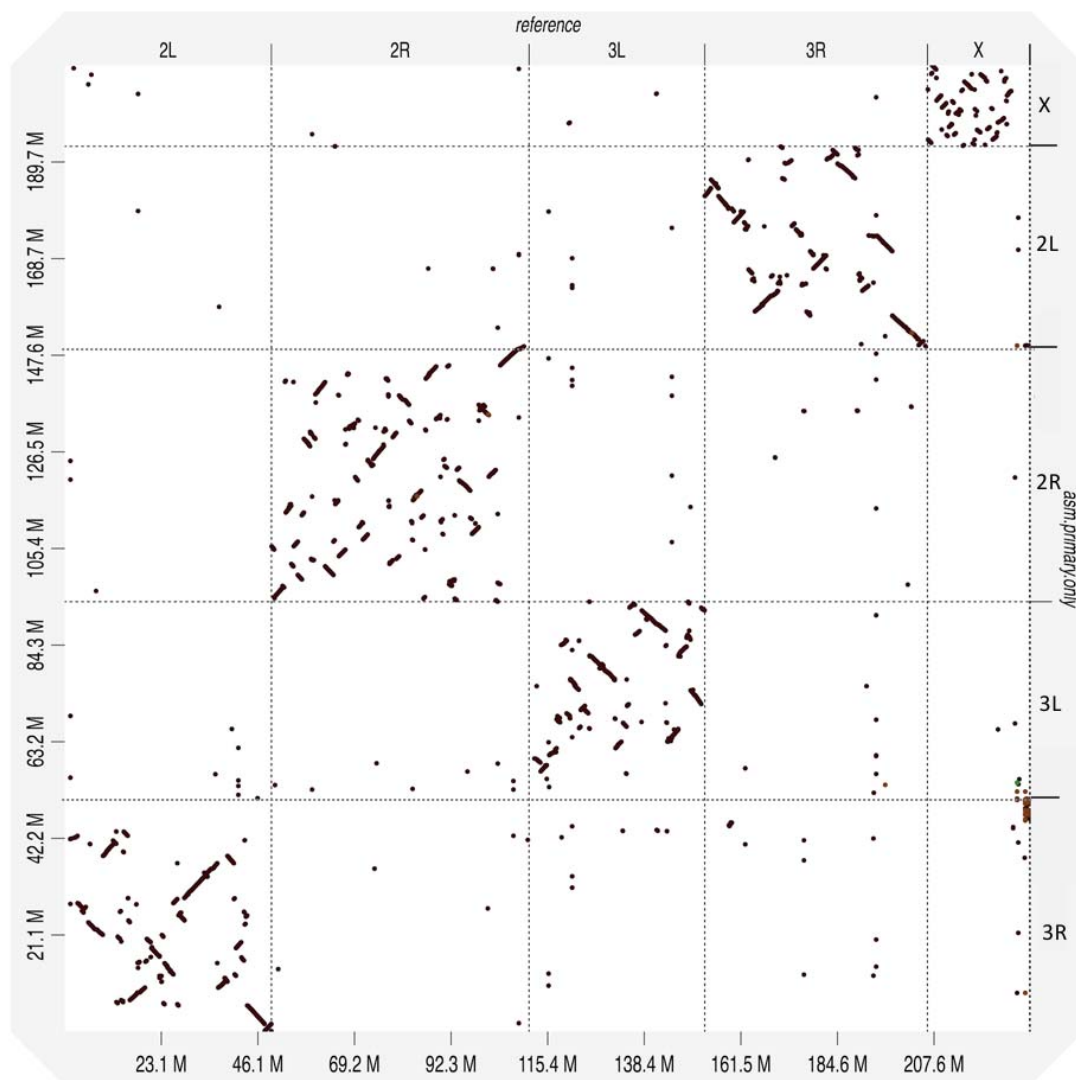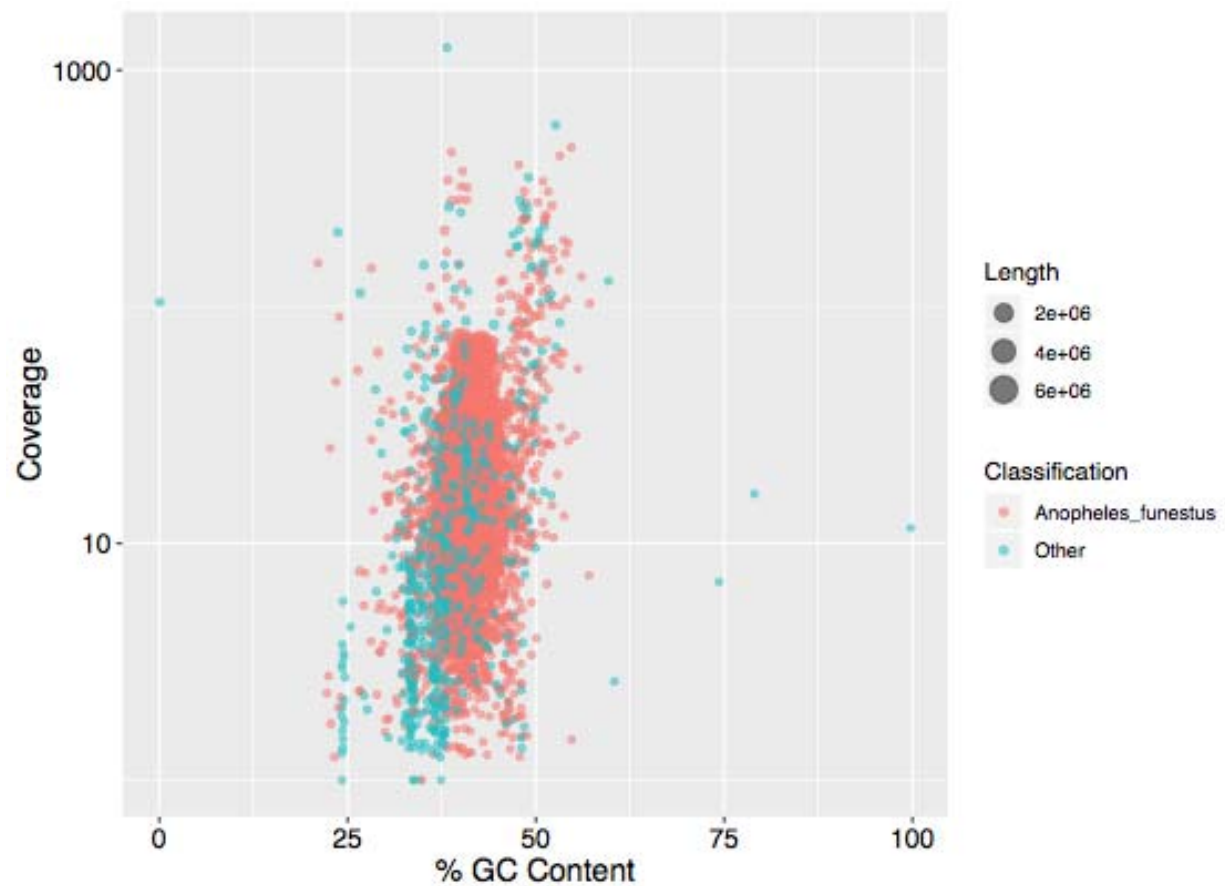
328   genomes.



329

330    Figure 4: GC content versus coverage plot for all assembled *An. funestus* contigs. The

331    orange points denote the contigs classified by Kraken as *An. funestus* and green points

332    denote everything else. A majority of the contigs are classified as *An. funestus* by

333    Kraken and there is no indication of extensive contamination.

334



335

336

337 **Tables**

338 Table 1: Assembly statistics for the *An. funestus* genome. *AfunF1* represents the prior

339 reference assembly, *AfunF3 contigs* denotes the complete long-read assembly with all

340 contigs included and *AfunF3 primary* denotes the assembly after deduplication and

341 scaffolding. QV(Illumina) denotes the assembly QV estimated using Illumina data and

342 QV(10X) denotes the 10X Genomics data. QV(Illumina) is highest for the AfunF1

343 assembly, because it is the same data used to generate that assembly, whereas

344 QV(10X) is based on data from a single mosquito of the same FUMOZ colony.

345

| Assembly | Number of Contigs | Contig N50 | Max Contig Size | Number of Scaffolds | Scaffold N50 | Max Scaffold size | Total Assembly Size | QV (Illumina) | QV (10X) |
|---|---|---|---|---|---|---|---|---|---|
| AfunF1 | 9,880 | 60,925 | 563,645 | 1,392 | 671,960 | 3,832,769 | 225,223,604 | 38.93 | 22.69 |
| AfunF3 contigs | 10,245 | 94,259 | 7,564,979 | 9,175 | 238,902 | 99,362,816 | 446,039,041 | 29.82 | 28.18 |
| AfunF3 primary | 1,053 | 631,722 | 7,564,979 | 3 | 93,811,348 | 99,362,816 | 210,827,327 | 24.94 | 25.82 |

346

18

347 Table 2: Validation of *An. funestus* genome assemblies using BUSCO gene set

348 completeness, agreement of the assemblies with RNA-Seq transcriptome data, and

349 structural accuracy inferred using PacBio long read data. *AfunF1* represents the prior

350 reference assembly, *AfunF3 contigs* denotes the complete long-read assembly with all

351 contigs included and *AfunF3 primary* denotes the assembly after deduplication and

352 scaffolding. For BUSCO categories C denotes "Complete Genes", S denotes "Single

353 Copy Genes", D denotes "Duplicated Genes", F denotes "Fragmented Genes", and M

354 denotes "Missing Genes". For long reads based structural variation, DEL denotes

355 deletions, DUP denotes duplications, INV denotes inversions, and INS denotes

356 insertions.

| Assembly | BUSCO statistics | | | | Transciptome data statistics | | | Structural variants called with long reads | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | C/S | C/D | F | M | Alignment Rate | Multi-mapped reads | % Transcripts in a single contig | DEL | DUP | INV | INS |
| AfunF1 | 2,756 | 16 | 27 | 16 | 81.79% | 23.92% | 84.96% | 9,036 | 455 | 152 | 3,798 |
| AfunF3 contigs | 2,765 | 1,068 | 18 | 17 | 84.34% | 36.97% | 91.16% | NA | NA | NA | NA |
| AfunF3 primary | 2,685 | 54 | 30 | 81 | 84.86% | 27.03% | 89.40% | 571 | 6 | 10 | 702 |

357

# References

358

359    1. Kim KE, Peluso P, Babayan P, Yeadon PJ, Yu C, Fisher WW, et al. Long-read, whole-

360    genome shotgun sequence data for five model organisms. Sci Data. 2014;1:140045.

361    2. Berlin K, Koren S, Chin C-S, Drake JP, Landolin JM, Phillippy AM. Assembling large

362    genomes with single-molecule sequencing and locality-sensitive hashing. Nat Biotechnol.

363    2015;33:623–30.

364    3. Neafsey DE, Christophides GK, Collins FH, Emrich SJ, Fontaine MC, Gelbart W, et al. The

365    evolution of the Anopheles 16 genomes project. G3 . 2013;3:1191–4.

366    4. Neafsey DE, Waterhouse RM, Abai MR, Aganezov SS, Alekseyev MA, Allen JE, et al.

367    Mosquito genomics. Highly evolvable malaria vectors: the genomes of 16 Anopheles

368    mosquitoes. Science. 2015;347:1258522.

369    5. Assogba BS, Milesi P, Djogbénou LS, Berthomieu A, Makoundou P, Baba-Moussa LS, et al.

370    The ace-1 Locus Is Amplified in All Resistant Anopheles gambiae Mosquitoes: Fitness

371    Consequences of Homogeneous and Heterogeneous Duplications. PLoS Biol.

372    2016;14:e2000618.

373    6. Weetman D, Djogbenou LS, Lucas E. Copy number variation (CNV) and insecticide

374    resistance in mosquitoes: evolving knowledge or an evolving problem? Curr Opin Insect Sci.

375    2018;27:82–8.

376    7. Coluzzi M. A Polytene Chromosome Analysis of the Anopheles gambiae Species Complex.

377    Science. 2002;298:1415–8.

378    8. Pombi M, Caputo B, Simard F, Di Deco MA, Coluzzi M, della Torre A, et al. Chromosomal

379    plasticity and evolutionary potential in the malaria vector Anopheles gambiae sensu stricto:

380    insights from three decades of rare paracentric inversions. BMC Evol Biol. 2008;8:309.

381    9. Sharakhov I. A Microsatellite Map of the African Human Malaria Vector Anopheles funestus. J

382    Hered. 2004;95:29–34.

383    10. Kirkpatrick M. How and why chromosome inversions evolve. PLoS Biol [Internet]. 2010;8.

384    Available from: http://dx.doi.org/10.1371/journal.pbio.1000501

385    11. Ma J, Amos CI. Investigation of inversion polymorphisms in the human genome using

386    principal components analysis. PLoS One. 2012;7:e40224.

387    12. Seich Al Basatena N-K, Hoggart CJ, Coin LJ, O'Reilly PF. The effect of genomic inversions

388    on estimation of population genetic parameters from SNP data. Genetics. 2013;193:243–53.

389    13. Houle D, Márquez EJ. Linkage Disequilibrium and Inversion-Typing of the Drosophila

390    melanogaster Genome Reference Panel. G3 . 2015;5:1695–701.

391    14. Coluzzi M, Sabatini A, Petrarca V, Di Deco MA. Chromosomal differentiation and adaptation

392    to human environments in the Anopheles gambiae complex. Trans R Soc Trop Med Hyg.

393    1979;73:483–97.

394    15. Main BJ, Lee Y, Ferguson HM, Kreppel KS, Kihonda A, Govella NJ, et al. The Genetic Basis

395    of Host Preference and Resting Behavior in the Major African Malaria Vector, Anopheles

396    arabiensis. PLoS Genet. 2016;12:e1006303.

397    16. Rishikesh N, Di Deco MA, Petrarca V, Coluzzi M. Seasonal variations in indoor resting

398    Anopheles gambiae and Anopheles arabiensis in Kaduna, Nigeria. Acta Trop. 1985;42:165–70.

399    17. Ayala D, Zhang S, Chateau M, Fouet C, Morlais I, Costantini C, et al. Association mapping

400    desiccation resistance within chromosomal inversions in the African malaria vector Anopheles

401    gambiae. Mol Ecol [Internet]. 2018; Available from: http://dx.doi.org/10.1111/mec.14880

402    18. Petrarca V, Nugud AD, Elkarim Ahmed MA, Haridi AM, Di Deco MA, Coluzzi M.

403    Cytogenetics of the Anopheles gambiae complex in Sudan, with special reference to An.

404    arabiensis: relationships with East and West African populations. Med Vet Entomol.

405    2000;14:149–64.

406    19. Gray EM, Rocca KAC, Costantini C, Besansky NJ. Inversion 2La is associated with

407    enhanced desiccation resistance in Anopheles gambiae. Malar J. 2009;8:215.

408    20. Rocca KAC, Gray EM, Costantini C, Besansky NJ. 2La chromosomal inversion enhances

409    thermal tolerance of Anopheles gambiae larvae. Malar J. 2009;8:147.

410    21. Fouet C, Gray E, Besansky NJ, Costantini C. Adaptation to aridity in the malaria mosquito

411    Anopheles gambiae: chromosomal inversion polymorphism and body size influence resistance

412    to desiccation. PLoS One. 2012;7:e34841.

413    22. Ayala D, Acevedo P, Pombi M, Dia I, Boccolini D, Costantini C, et al. Chromosome

414    inversions and ecological plasticity in the main African malaria mosquitoes. Evolution.

415    2017;71:686–701.

416    23. Cheng C, Tan JC, Hahn MW, Besansky NJ. Systems genetic analysis of inversion

417    polymorphisms in the malaria mosquito. Proc Natl Acad Sci U S A. 2018;115:E7005–14.

418    24. Ayala D, Caro-Riaño H, Dujardin J-P, Rahola N, Simard F, Fontenille D. Chromosomal and

419    environmental determinants of morphometric variation in natural populations of the malaria

420    vector Anopheles funestus in Cameroon. Infect Genet Evol. 2011;11:940–7.

421    25. Riehle MM, Bukhari T, Gneme A, Guelbeogo WM, Coulibaly B, Fofana A, et al. The

422    Anopheles gambiae 2La chromosome inversion is associated with susceptibility to Plasmodium

423    falciparum in Africa. Elife [Internet]. 2017;6. Available from: http://dx.doi.org/10.7554/elife.25813

424   26. Petrarca V, Beier JC. Intraspecific chromosomal polymorphism in the Anopheles gambiae

425   complex as a factor affecting malaria transmission in the Kisumu area of Kenya. Am J Trop Med

426   Hyg. 1992;46:229–37.

427   27. Gillies MT, De Meillon B. The Anophelinae of Africa South of the Sahara: (Ethiopian

428   Zoogeographical Region). 1968.

429   28. Coetzee M, Fontenille D. Advances in the study of Anopheles funestus, a major vector of

430   malaria in Africa. Insect Biochem Mol Biol. 2004;34:599–605.

431   29. Coetzee M, Koekemoer LL. Molecular systematics and insecticide resistance in the major

432   African malaria vector Anopheles funestus. Annu Rev Entomol. 2013;58:393–412.

433   30. Dia I, Guelbeogo MW, Ayala D. Advances and Perspectives in the Study of the Malaria

434   Mosquito Anopheles funestus. Anopheles mosquitoes - New insights into malaria vectors. 2013.

435   31. Zahar AR, World Health Organization. Vector Bionomics in the Epidemiology and Control of

436   Malaria: The WHO African region & the southern WHO eastern Mediterranean region. 1984.

437   32. Menze BD, Riveron JM, Ibrahim SS, Irving H, Antonio-Nkondjio C, Awono-Ambene PH, et

438   al. Multiple Insecticide Resistance in the Malaria Vector Anopheles funestus from Northern

439   Cameroon Is Mediated by Metabolic Resistance Alongside Potential Target Site Insensitivity

440   Mutations. PLoS One. 2016;11:e0163261.

441   33. Riveron JM, Ibrahim SS, Mulamba C, Djouaka R, Irving H, Wondji MJ, et al. Genome-Wide

442   Transcription and Functional Analyses Reveal Heterogeneous Molecular Mechanisms Driving

443   Pyrethroids Resistance in the Major Malaria Vector Anopheles funestus Across Africa. G3:

444   Genes|Genomes|Genetics. 2017;g3.117.040147.

445   34. Ndo C, Kopya E, Donbou MA, Njiokou F, Awono-Ambene P, Wondji C. Elevated

446    Plasmodium infection rates and high pyrethroid resistance in major malaria vectors in a forested

447    area of Cameroon highlight challenges of malaria control. Parasit Vectors [Internet]. 2018;11.

448    Available from: http://dx.doi.org/10.1186/s13071-018-2759-y

449    35. Michel AP, Ingrasci MJ, Schemerhorn BJ, Kern M, Le Goff G, Coetzee M, et al. Rangewide

450    population genetic structure of the African malaria vector Anopheles funestus. Mol Ecol.

451    2005;14:4235–48.

452    36. Michel AP, Guelbeogo WM, Grushko O, Schemerhorn BJ, Kern M, Willard MB, et al.

453    Molecular differentiation between chromosomally defined incipient species of Anopheles

454    funestus. Insect Mol Biol. 2005;14:375–87.

455    37. Guelbeogo WM, Grushko O, Boccolini D, Ouédraogo PA, Besansky NJ, Sagnon NF, et al.

456    Chromosomal evidence of incipient speciation in the Afrotropical malaria mosquito Anopheles

457    funestus. Med Vet Entomol. 2005;19:458–69.

458    38. Costantini C, Sagnon N, Ilboudo-Sanogo E, Coluzzi M, Boccolini D. Chromosomal and

459    bionomic heterogeneities suggest incipient speciation in Anopheles funestus from Burkina Faso.

460    Parassitologia. 1999;41:595–611.

461    39. Guelbeogo WM, Sagnon N 'fale, Grushko O, Yameogo MA, Boccolini D, Besansky NJ, et al.

462    Seasonal distribution of Anopheles funestus chromosomal forms from Burkina Faso. Malar J.

463    2009;8:239.

464    40. Matthews BJ, Dudchenko O, Kingan SB, Koren S, Antoshechkin I, Crawford JE, et al.

465    Improved reference genome of Aedes aegypti informs arbovirus vector control. Nature.

466    2018;563:501–7.

467    41. Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, et al.

468    BUSCO applications from quality assessments to gene prediction and phylogenomics. Mol Biol

469     Evol [Internet]. 2017; Available from: http://dx.doi.org/10.1093/molbev/msx319

470     42. Chin C-S, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, et al. Nonhybrid,

471     finished microbial genome assemblies from long-read SMRT sequencing data. Nat Methods.

472     2013;10:563–9.

473     43. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an

474     integrated tool for comprehensive microbial variant detection and genome assembly

475     improvement. PLoS One. 2014;9:e112963.

476     44. Ewing B, Hillier L, Wendl MC, Green P. Base-calling of automated sequencer traces using

477     phred. I. Accuracy assessment. Genome Res. 1998;8:175–85.

478     45. Waterhouse RM, Aganezov S, Anselmetti Y, Lee J, Ruzzante L, Reijnders MJ, et al.

479     Leveraging evolutionary relationships to improve Anopheles genome assemblies [Internet].

480     2018. Available from: http://dx.doi.org/10.1101/434670

481     46. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and

482     accurate long-read assembly via adaptivek-mer weighting and repeat separation. Genome Res.

483     2017;27:722–36.

484     47. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.

485     2012;9:357–9.

486     48. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the gap: upgrading

487     genomes with Pacific Biosciences RS long-read sequencing technology. PLoS One.

488     2012;7:e47768.

489     49. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and

490     open software for comparing large genomes. Genome Biol. 2004;5:R12.

491    50. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, et al. Circos: an

492    information aesthetic for comparative genomics. Genome Res. 2009;19:1639–45.

493    51. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, et al. De novo assembly

494    of haplotype-resolved genomes with trio binning. Nat Biotechnol [Internet]. 2018; Available from:

495    http://dx.doi.org/10.1038/nbt.4277

496    52. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform.

497    Bioinformatics. 2010;26:589–95.

498    53. Bishara A, Liu Y, Weng Z, Kashef-Haghighi D, Newburger DE, West R, et al. Read clouds

499    uncover variation in complex regions of the human genome. Genome Res. 2015;25:1570–80.

500    54. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact

501    alignments. Genome Biol. 2014;15:R46.

502    55. Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A, et al.

503    Accurate detection of complex structural variations using single-molecule sequencing. Nat

504    Methods. 2018;15:461–8.

505    56. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory

506    requirements. Nat Methods. 2015;12:357–60.

507    57. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length

508    transcriptome assembly from RNA-Seq data without a reference genome. Nat Biotechnol.

509    2011;29:644–52.

510    58. Wu TD, Watanabe CK. GMAP: a genomic mapping and alignment program for mRNA and

511    EST sequences. Bioinformatics. 2005;21:1859–75.

512    59. Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, Lander ES, et al. Juicebox

513    Provides a Visualization System for Hi-C Contact Maps with Unlimited Zoom. Cell Syst.

514    2016;3:99–101.

515    60. Cabanettes F, Klopp C. D-GENIES: dot plot large genomes in an interactive, efficient and

516    simple way. PeerJ. 2018;6:e4958.

517