1 **Whole Genome Sequencing of *Plasmodium vivax* Isolates Reveals**

2 **Frequent Sequence and Structural Polymorphisms in Erythrocyte**

3 **Binding Genes**

4

5 Anthony Ford[1,2*], Daniel Kepple[2*], Beka Raya Abagero[3], Jordan Connors[1], Richard

6 Pearson[4], Sarah Auburn[5], Sisay Getachew[6, 7], Colby Ford[1], Karthigayan Gunalan[8],

7 Louis H. Miller[8], Daniel A. Janies[1], Julian C. Rayner[9], Guiyun Yan[10], Delenasaw

8 Yewhalaw[3], Eugenia Lo[2]

9 * Co-first authors with equal contribution

10

11 [1] Department of Bioinformatics and Genomics, University of North Carolina at Charlotte,

12 USA

13 [2] Department of Biological Sciences, University of North Carolina at Charlotte, USA

14 [3] Tropical Infectious Disease Research Center, Jimma University, Ethiopia

15 [4] Malaria Programme, Wellcome Trust Sanger Institute, Hinxton, UK

16 [5] Global and Tropical Health Division, Menzies School of Health Research and Charles

17 Darwin University, Darwin, Northern Territory, Australia

18 [6] College of Natural Sciences, Addis Ababa University, Ethiopia

19 [7] Armauer Hansen Research Institute, Addis Ababa, Ethiopia

20 [8] Laboratory of Malaria and Vector Research, NIAID/NIH, Bethesda, USA

21 [9] Department of Clinical Biochemistry, Cambridge Institute for Medical Research,

22 University of Cambridge, Cambridge CB2 OXY, UK

23    [10] Program in Public Health, University of California at Irvine, USA

24

25    **Correspondence:** Anthony Ford, Department of Bioinformatics and Genomics,

26    University of North Carolina at Charlotte, USA; Guiyun Yan, Program in Public Health,

27    University of California at Irvine, USA; Eugenia Lo, Department of Biological Sciences,

28    University of North Carolina at Charlotte

29

30    **Running Title:** Genomic characteristics of *Plasmodium vivax* in Ethiopia

31

32

33

34

35

36

37

38

39

40

41

42

43

## Abstract

*Plasmodium vivax* malaria is much less common in Africa than the rest of the world because the parasite relies primarily on the Duffy antigen/chemokine receptor (*DARC*) to invade human erythrocytes, and the majority of Africans are Duffy negative. Recently, there has been a dramatic increase in the reporting of *P. vivax* cases in Africa, with a high number of them being in Duffy negative individuals, potentially indicating *P. vivax* has evolved an alternative invasion mechanism that can overcome Duffy negativity. Here, we analyzed single nucleotide polymorphism (SNP) and copy number variation (CNV) in Whole Genome Sequence (WGS) data from 44 *P. vivax* samples isolated from symptomatic malaria patients in southwestern Ethiopia, where both Duffy positive and Duffy negative individuals are found. A total of 236,351 SNPs were detected, of which 21.9% was nonsynonymous and 78.1% was synonymous mutations. The largest number of SNPs were detected on chromosomes 9 (33,478 SNPs; 14% of total) and 10 (28,133 SNPs; 11.9%). There were particularly high levels of polymorphism in erythrocyte binding gene candidates including reticulocyte binding protein 2c (*RBP*2c), merozoite surface protein 1 (*MSP*1), and merozoite surface protein 3 (*MSP*3.5, *MSP*3.85 and *MSP*3.9). Thirteen genes related to immunogenicity and erythrocyte binding function were detected with significant signals of positive selection. Variation in gene copy number was also concentrated in genes involved in host-parasite interactions, including the expansion of the Duffy binding protein gene (*PvDBP*) on chromosome 6 and several *PIR* genes. Based on the phylogeny constructed from the whole genome sequences, the expansion of these genes was an independent process among the *P. vivax* lineages in Ethiopia. We further inferred transmission patterns of *P.*

67    *vivax* infections among study sites and showed various levels of gene flow at a small

68    geographical scale. The genomic features of *P. vivax* provided baseline data for future

69    comparison with those in Duffy-negative individuals, and allowed us to develop a panel

70    of informative Single Nucleotide Polymorphic markers diagnostic at a micro-

71    geographical scale.

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

## Introduction

Vivax malaria is the most geographically widespread human malaria, causing over 130 million clinical cases per year worldwide [1]. *Plasmodium vivax* can produce dormant liver-stage hypnozoites within infected hosts, giving rise to relapse infections from months to years. This unique feature of *P. vivax* contributes to an increase in transmission potential and increases the challenge of elimination [2]. Understanding *P. vivax* genome variation will advance our knowledge of parasite biology and host-parasite interactions, as well as identify potential drug resistance mechanisms [3, 4]. Such data will also help identify molecular targets for vaccine development [5-7], and provide new means to track the transmission and spread of drug resistant parasites [8-9].

Compared to *P. falciparum*, *P. vivax* isolates from Southeast Asia (e.g., Thailand and Myanmar), Pacific Oceania (Papua New Guinea), and South America (Mexico, Peru, and Colombia) have significantly higher nucleotide diversity at the genome level [2]. This could be the result of frequent gene flow via human movement, intense transmission, and/or variation in host susceptibility [10-14]. *P. vivax* infections are also much more likely to contain multiple parasite strains in areas where transmission is intense and/or relapse is common [10, 15-18]. In Papua New Guinea, for example, *P. vivax* infections had an approximately 3.5-fold higher rate of polyclonality and nearly double the multiplicity of infection (MOI) than the *P. falciparum* infections [16]. Similar rates of polyclonality and MOI have also been reported in *P. vivax* in Cambodia [6]. It is possible intense transmission has sustained a large and stable parasite population in these regions [17,18]. By contrast, geographical differentiation and selection pressure

112    over generations can lead to fixation of parasite genotypes in local populations. In the

113    Asia-Pacific region, *P. vivax* showed a high level of genetic relatedness through

114    inbreeding among the dominant clones, in addition to strong selection imposed in a

115    number of antimalarial drug resistance genes [19]. In Ethiopia, the chloroquine

116    resistance transporter gene (*Pvcrt*) of *P. vivax* on chromosome 14 had been shown with

117    significant selection in a region upstream of the promotor, highlighting the ability of *P.*

118    *vivax* to rapidly evolve in response to control measures [20]. Apart from mutations, high

119    copy number observed in *Pvcrt* and multidrug resistant gene (*Pvmdr*1) has also been

120    shown to be associated with increased antimalaria drug resistance [21,22].

121        Recent genomic studies have indicated that some highly polymorphic genes in

122    the *P. vivax* genome are associated with red blood cell invasion and immune evasion

123    [10, 12, 19, 23]. They include the merozoite surface protein genes *MSP*1

124    (PVP01_0728900) and *MSP*7 (PVX_082665), Pv-fam-b (PVX_002525), Pv-fam-e

125    (PVX_089875), the reticulocyte binding protein gene *RBP*2c (PVP01_0534300), serine-

126    repeat antigen 3 (*SERA*; PVX_003840), as well as virulent genes (*VIR*) such as *VIR*22

127    (PVX_097530) and *VIR*12 (PVX_083590) [23-29]. Polymorphisms in genes associated

128    with immune evasion and reticulocyte invasion have important implications for the

129    invasion efficiency and severity of *P. vivax* infections. Members of the erythrocyte

130    binding gene family, including reticulocyte binding proteins (*RBP*s), Duffy-binding

131    proteins (*DBP*s), and merozoite surface proteins (*MSP*3 and *MSP*7) have been

132    previously shown to exhibit high sequence variation in *P. vivax* [20, 30]. The

133    polymorphisms in *RBP*1 and *RBP*2 genes may relate to an increased capability of

134    erythrocyte invasion by *P. vivax* [31-33]. It has been suggested that Pv*RBP*2b-TfR1

135    interaction is vital for the initial recognition and invasion of host reticulocytes [34], prior

136    to the engagement of *PvDBP1* and Duffy antigen chemokine receptor (*DARC*) and

137    formation of a tight junction between parasite and erythrocyte [35]. Apart from Pv*RBP*,

138    Reticulocyte Binding Surface Antigen (Pv*RBSA*) [36], an antigenic adhesin, may also

139    play a key role in *P. vivax* parasites binding to target cells, possessing the capability of

140    binding to a population of reticulocytes with a different Duffy phenotype [37, 38].

141    Another erythrocyte binding protein gene (Pv*EBP*), a paralog of *PvDBP1*, which harbors

142    all the hallmarks of a *Plasmodium* red blood cell invasion protein, including conserved

143    Duffy-binding like and C-terminal cysteine-rich domains [39], has been recently shown

144    to be variable in copy number in the Malagasy *P.* vivax [39]. Functional analyses

145    indicated that region II of this gene bound to both Duffy-positive and Duffy-negative

146    reticulocytes, although at a lower frequency compared to *PvDBP*, suggestive of its role

147    in erythrocyte invasion [40]. Both Pv*EBP*1 and Pv*EBP*2 genes exhibit high genetic

148    diversity and are common antibody binding targets associated with clinical protection

149    [41, 42]. Other proteins such as tryptophan-rich antigen gene (*TRAg*), anchored

150    micronemal antigen (*GAMA*), and Rhoptry neck protein (*RON*) have also been

151    suggested to play a role in red cell invasion, especially in low-density infections [43-47].

152    Information of the polymorphisms in these genes will have important implications on the

153    dynamics of host-parasite interactions.

154        Compared to Southeast Asia and South America where *P. vivax* is highly

155    endemic, data on polymorphisms in erythrocyte binding gene candidates of *P. vivax*

156    from Africa is limited. Filling the gap is critical for identifying functional genes in

157    erythrocyte invasion, biomarkers for tracking the African *P. vivax* isolates, as well as

158   potential gene targets for vaccine development. It was previously thought that most

159   African populations were immune to *P.* vivax infections due to the absence of *DARC*

160   gene expression required for erythrocyte invasion. However, several recent reports

161   have indicated the emergence and potential spread of *P. vivax* across Africa [32, 48-

162   50]. The objective of this study was to describe genomic variation of *P. vivax* from

163   Ethiopia. Specifically, we examined the level of genetic polymorphisms in a panel of 64

164   potential erythrocyte binding protein genes that have been suggested to play a role in

165   the parasite-host invasion process. In addition, we inferred transmission patterns of *P.*

166   *vivax* infections from different study sites based on the genetic variants. A recent study

167   by Auburn *et al*. [20] has compared the genetic variants of *P. vivax* from Ethiopia with

168   other geographical isolates. In the present study, we focus on the genomic

169   characteristics of *P. vivax* among different study sites in Ethiopia with the goals to

170   establish a baseline for genome comparison with the Duffy-negative *P. vivax* in our

171   ongoing investigation, as well as to develop a panel of informative Single Nucleotide

172   Polymorphic (SNP) markers diagnostic at a micro-geographical scale.

173

## Materials and Methods

### Ethics statement

176   Scientific and ethical clearance was obtained from the Institutional Scientific and Ethical

177   Review Boards of Jimma and Addis Ababa Universities in Ethiopia, and The University

178   of North Carolina, Charlotte, USA. Written informed consent/assent for study

179   participation was obtained from all consenting heads of households, parents/guardians

180    (for minors under age of 18), and each individual who was willing to participate in the

181    study.

182

**Study area and sample collection**

184    Genomic DNA was extracted from 22 clinical samples collected in Jimma, southwestern

185    Ethiopia during peak transmission season (September – November, 2016; Figure 1).

186    Finger-pricked blood samples were collected from malaria symptomatic (who has fever

187    with axillary body temperature > 37.5°C and with confirmed asexual stages of malaria

188    parasite based on microscopy) or febrile patients visiting the health centers or hospitals

189    at each of the study sites. Thick and thin blood smears were prepared for microscopic

190    examination, and 4-6 ml of venous blood were collected from each *P. vivax*-confirmed

191    patient in K2 EDTA blood collection tubes. For the whole blood samples, we used the

192    Lymphoprep/Plasmodpur-based protocol to deplete the white blood cells and enrich the

193    red blood cell pellets [51]. DNA was then extracted from approximately 1 ml of the red

194    blood cell pellets using Zymo Bead Genomic DNA kit (Zymo Research) following the

195    manufacturer's procedures. The extracted DNA were first assessed by nested and

196    quantitative PCR methods to confirm and quantify *P. vivax* of the infected samples [52].

197    From a larger set of samples, we then performed microsatellite analyses using seven

198    different loci [53].  Only monoclonal samples were selected and proceeded for

199    sequencing. Whole genome sequencing was conducted on the Illumina HiSeq 3000

200    Sequencing Platform at the Wellcome Sanger Institute (European Nucleotide Archive

201    [ENA] accession number of each sample in Table 1). The generated sequence reads

202    were mapped individually to the publicly available reference genome PvP01 from Gene

203   DB using Bowtie version 2 [54]. The original 22 samples were processed to remove

204   reads other than *P. vivax*. The percentage coverage of the *P. vivax* reads in our

205   samples were high enough to not affect the results. An additional 24 sample sequence

206   data were obtained as FASTQ files from the ENA. These samples were collected from

207   Arbaminch, Badowacho, Halaba, and Hawassa in southwestern Ethiopia (Figure 1), the

208   Duffy status of each of these 24 samples is unknown. They were then aligned to the

209   PVP01 reference genome using BWA-MEMv.2 with default settings [55, 56]. The overall

210   quality of each resulting BAM was assessed using FASTQC.  Similarly, we concluded

211   that the percentage of the *P. vivax* reads covered in the additional 24 samples were

212   high enough to reflect the dominant signal of the variants and negate polyclonal

213   influences. Two of our samples displayed a significant decline in average quality in read

214   mapping and were therefore removed from further SNP variant and copy number

215   variation analyses.

216

217   **SNP discovery, annotation, and filtering**

218   Potential SNPs were identified by SAM tools v.1.6 mpileup procedure [57] in conjunction

219   with BCF tools v.1.6 [57] across all 44 sample BAM files using the PVP01 reference

220   genome.  Compared to the Salvador-I, the PVP01 reference genome consists of 14

221   major chromosomal sequences, and provides a greater level of gene annotation power

222   and improved assembly of the subtelomeres [56]. We analyzed only sequence reads

223   that were mapped to these 14 major chromosomal sequences. The hypervariable and

224   subtelomeric regions in our samples were retained during the variant calling procedure

225   and each sample BAM file had duplicates marked using SAMtools 1.6 markdup

226    procedure. For the mpileup procedure, the maximum depth threshold, which determines

227    the number of maximum reads per file at a position, was set to 3,000 million to ensure

228    that the maximum amount of reads for each position was not reached.  Samples were

229    pooled together using a multisampling variant calling approach. The SNPs were then

230    annotated with SnpEff v.4.3T [58] based on the annotated gene information in GeneDB.

231    Filtering was done using the following standard metrics, including Read Position Bias,

232    Mapping Quality vs Strand Bias, Raw read depth, Mapping Quality Bias, Base Quality

233    Bias, and Variant Distant Bias produced by SAM tools and BCF tools during the variant

234    calling procedure. In Snp Sift, data was filtered by choosing SNPs that had a Phred

235    Quality score $\geq$ 40, a raw read depth (DP) $\geq$ 30, and a base quality bias >0.1 [59].  We

236    then calculated the allele frequency for each SNP position for all 44 samples using the

237    frequency procedure in VCF tools v.0.1.15 [60]. The total number of SNPs across all

238    samples, as well as the number of nonsynonymous and synonymous mutations were

239    recorded. Mutations were compared among the 14 chromosomes in addition to a panel

240    of 64 erythrocyte binding genes.

241

242    **Copy number variation analyses**

243    Copy number variation of gene regions was assessed with CNVnator [61]. CNVnator

244    uses mean-shift theory, a partitioning procedure based on an image processing

245    technique and additional refinements including multiple bandwidth partitioning and GC

246    correction [61]. We first calculated the read depth for each bin and correct GC-bias. This

247    was followed by mean-shift based segment partition and signal merging, which

248    employed an image processing technique.  We then performed CNV calling, of which

249 segments with a mean RD signal deviating by at least a quarter from genomic average

250 read depth signal were selected and regions with a *P*-value less than 0.05 were called.

251 A one-sided test was then performed to call additional copy number variants. SAM tools

252 v.1.6 was utilized in our data preprocessing step to mark potential duplicates in the BAM

253 files and followed the CNV detection pipeline [62]. We extracted the read mappings

254 from each of BAM files for all chromosomes. Once the root file was constructed using

255 the extracted reads, we generated histograms of the read depths using a bin size of

256 100. The statistical significance for the windows that showed unusual read depth was

257 calculated and the chromosomes were partitioned into long regions that have similar

258 read depth.

259      To validate the results from CNVnator, we used the GATK4 copy number

260 detection pipeline to further examine gene copy number [63-65]. The read coverage

261 counts were first obtained from pre-processed genomic intervals of a 1000-bp window

262 length based on the PvP01 reference genome. The read fragment counts were then

263 standardized using the Denoise Read Counts that involved two transformations. The

264 first transformation was based on median counts, including the $\log_2$ transformation, and

265 the counts were normalized to center around one. In the second transformation, the tool

266 denoises was used to standardized copy ratios using principal component analysis.

267

268 **Test for positive selection**

269 Regions of positive selection were examined among the 44 Ethiopian *P. vivax* isolates

270 using the integrated haplotype score approach, specifically the SciKit-Allel for python, a

271 package used for analysis of large scale genetic variation data [66]. Before the samples

272   were run through Scikit-Allel, genotypes for each of the samples were phased using

273   BEAGLE [67]. Genes that were detected with signals of positive selection by SciKit-

274   Allel, as well as a panel of 64 potential erythrocyte binding genes were further evaluated

275   using the PAML package (Phylogenetic Analysis by Maximum Likelihood) [68]. Using

276   the codeml procedure in PAML, DNA sequences were analyzed with the maximum

277   likelihood approach in a phylogenetic framework. The synonymous and

278   nonsynonymous mutation rates between protein-coding DNA sequences were then

279   estimated in order to identify potential regions of positive selection. We created two

280   models, the neutral model M1 and the selection model M2. The average $d_N/d_S$ values

281   were estimated across all branches in both M1 and M2 models and the average $d_N/d_S$

282   values across all sites in the M2 model. The $d_N/d_S$ values were compared between the

283   two models using a likelihood ratio test for significant positive selection.

284

285   **Comparison of nucleotide diversity among EBP gene regions**

286   Based on the literature [23-33], we identified 64 gene regions that are potentially related

287   to erythrocyte binding in *P. vivax* (Supplementary Table 1). These included the *DBP*

288   (duffy binding protein), *EBP* (erythrocyte binding protein), *MSP* (merozoite surface

289   protein), and *RBP* (reticulocyte binding protein) multigene families, the tryptophan rich

290   antigen gene family (*TRAg*), GPI-anchored microanemal antigen (*GAMA)*, microneme

291   associated antigen (*MA)*, rhoptry associated adhesin (*RA*), high molecular weight

292   rhoptry protein 3 (*RHOP*3), and rhoptry neck protein (*RON)* genes. Previous study has

293   shown that the transcriptome profiles of the *TRAg* genes were differentially transcribed

294   at the erythrocytic stages, indicating that these genes may play specific roles in blood-

295 stage development [43].  The reticulocyte binding protein multigene family encodes

296 genes that each have a receptor on the surface that is essential for the host-invasion

297 stage of *P. vivax* [69]. The *MSP* multigene family, currently assumed to be a candidate

298 for vaccine generation, also plays a role in the invasion stage of *P. vivax* and is also

299 immunogenic [26]. The nucleotide diversity of 64 potential erythrocyte binding genes

300 were compared among the 44 *P. vivax* sample consensus sequences using DnaSP

301 [70]. The Pairwise-Deletion method where gaps were ignored in each pairwise

302 comparison was used for this calculation.

303

**Genetic relatedness and transmission network analyses**

305 Phylogenetic analyses were performed to infer the genetic relatedness among the 44

306 Ethiopian isolates. Sequence alignment was first conducted using a multiple sequence

307 alignment program in MAFFT v. 7 [71]. The alignment was then trimmed to remove

308 gaps using trimal (the *gappyout* option) that trimmed the alignments based on the gap

309 percentage count over the whole alignment. After sequence editing, we concatenated

310 all alignment files using FASconCAT-G [72], a perl program that allows for

311 concatenation and translation (nucleotide to amino acid states) of multiple alignment

312 files for phylogenetic analysis. We used the maximum likelihood method implemented in

313 the Randomized Accelerated Maximum Likelihood (RAxML) v8 to construct

314 phylogenetic trees [73]. The GTRGAMMA model was used for the best-scoring

315 maximum likelihood tree. The GTR model incorporates the optimization of substitution

316 rates and the GAMMA model accounts for rate heterogeneity.  A total of 100 rapid

317 bootstrap runs were conducted to evaluate the confidence of genetic relationships. In

318    addition, we performed principal component analyses using the glPCA function in R, a

319    subset of the adegenet package [74], to determine the genetic relatedness of the

320    samples among the different study sites in Ethiopia. A transmission network was

321    created using StrainHub, a tool for generating transmission networks using phylogenetic

322    information along with isolate metadata [75]. The transmission network was generated

323    using the locations of the samples as the nodes and calculating the source hub ratio

324    between each location. The source hub ratio was calculated by the number of

325    transitions originating from a node over the total number of transitions related to that

326    node. A node with a ratio close to 1 indicates a source, a ratio close to 0.5 indicates a

327    hub, and a ratio close to 0 indicates a sink for the *P. vivax* infections.

328

## Results

330    **Distribution of SNPs among the chromosomes and EBP genes**

331    A total of 252,973 SNPs were detected among the 44 Ethiopian *P. vivax* samples

332    (Figure 2), with 21.5% (54,336 out of 252,973) nonsynonymous and 78.5% (198,637 out

333    of 252,973) synonymous mutations (Figure 3A). The highest number of SNPs were

334    observed on chromosomes 7 (28,856 SNPs; 11.4%), 9 (28,308 SNPs; 11.2%), and 12

335    (28,190 SNPs; 11.1%); whereas the lowest number of SNPs were observed on

336    chromosomes 3 (6,803 SNPs; 2.7%), 6 (5,044 SNPs; ~2%), and 13 (8,809 SNPs; 3.5%;

337    Figure 3A; Supplementary Table 2).

338

339    The 64 erythrocyte binding genes accounted for 3,607 of the total SNPs, with 1685

340    (46.7%) identified as nonsynonymous and 1922 (53.3%) as synonymous mutations

341    (Figure 3B).  Among these genes, the highest number of SNPs were observed in

342    reticulocyte binding protein gene (*RBP*2c) on chromosome 5, followed by the *MSP*3

343    multigene family (*MSP*3.5, *MSP*3.9 and *MSP*3.8) on chromosome 10.  Nucleotide

344    diversity also showed to be highest in the *RBP* and *MSP*3 multigene families, with an

345    average nucleotide diversity of 1.3% and 2.8%, respectively, among our samples

346    (Figure 3B).  By contrast, the lowest number of SNPs were observed in the Duffy

347    binding protein gene (*DBP*1) on chromosome 6 with a total of 13 SNPs, of which 12

348    were identified as nonsynonymous and one as synonymous mutations (Figure 3B).

349    Likewise, another erythrocyte binding protein (*EBP*2), located also on chromosome 6,

350    was one of the least variable genes with only one nonsynonymous mutation.  The *TRAg*

351    gene family also showed a low level of nucleotide diversity when compared to the other

352    *EBP* gene families with an average nucleotide diversity of 0.2% (Figure 3B).

353

**Gene regions under positive selection**

355    Based on the integrated haplotype scores, positive selection was detected in 13 gene

356    regions (Figure 4). These included the sub-telomeric protein 1 (*STP*1) on chromosome

357    5, the membrane associated erythrocyte binding-like protein (*MAEBL*) on chromosome

358    9, *MSP*3.8 on chromosome 10, as well as various plasmodium interspersed repeats

359    (*PIR)* protein genes on chromosomes 3, 5, 7, 10, 11, and 12 (Figure 4).  Based on

360    PAML, 25 out of the 64 erythrocyte binding genes showed evidence of positive selection

361    (Table 2; Supplementary Table 3). The majority of these genes belong to the *TRAg*

362    multigene family. The *TRAg* genes had an average $d_N/d_S$ ratio of 2.75 across all

363    branches and an average of 5.75 across all sites for the M2 model tested for selection

364 (Table 2). Compared to the other *TRAg* genes, *TRAg*15 had more sites detected under

365 positive selection, with 50 of the sites showing a posterior probability greater than 50%

366 and 43 showing a posterior probability greater than 95% (Table 2). While the *TRAg*4

367 gene had the highest $d_N/d_S$ ratio across all sites among other *TRAg* genes, only six sites

368 were shown under positive selection with a posterior probability greater than 50% and

369 one with a posterior probability greater than 95%.

370 　　　All *RBP* genes, except for *RBP*2c, showed regions with significant signals of

371 positive selection (average $d_N/d_S$ ratio across all sites: 5.11; Table 2). Among them,

372 *RBP*2p1 had the largest number of sites with posterior probabilities greater than 95%

373 (Table 2). Among all the *MSP* genes, only *MSP*5, *MSP*9, and *MSP*10 indicated regions

374 under positive selection. The *MSP*5 and *MSP*9 genes had an average $d_N/d_S$ ratio of

375 3.85 across all sites and 1.11 across branches (Table 2). While *MSP*10 had an average

376 $d_N/d_S$ ratio of 1.16 across all branches and less than 1 across all sites, only seven sites

377 were indicated with posterior probabilities greater than 50% and 95% (Table 2).

378 Although *MSP*3.8 showed potential positive selection based on the integrated haplotype

379 scores (Figure 4), PAML did not show significant evidence of positive selection. For the

380 *DBP* gene family, *DBP*9 showed the highest $d_N/d_S$ ratio across all sites and branches

381 (10.39 and 3.88, respectively; Table 2).

382

383 **Copy number variation and evolution of high-order copy variants**

384 According to CNVnator, 19 gene regions showed copy number variation among our

385 samples (Figure 5; Supplementary Table 4). Among them, 11 gene regions were

386 detected with up to 2-3 copies and 8 gene regions with 4 copies or higher. We observed

387  copy number variation in several *PIR* genes distributed across chromosomes 1, 2, 4, 5,

388  7, 10 and 12 (Figure 5; Supplementary Table 4). Specifically, for the *PIR* genes located

389  on chromosome 2 (including PVP01_0220700, PVP01_0200200, PVP01_0200300, and

390  PVP01_0200100; Figure 5), more than 20% of the samples had 2-3 copies and

391  approximately 2-4% of the samples had 4 copies or higher. Among the 64 erythrocyte

392  binding genes, duplications were observed in *DBP*1 on chromosome 6 and *MSP*3 on

393  chromosome 10.  *DBP*1 ranged from one to as high as five copies, and *MSP*3 ranged

394  from one to as high as three copies among our samples (Figure 5), consistent with

395  previous findings [19, 20, 76]. The remaining erythrocyte binding genes were detected

396  with a single copy across our samples.

397      A maximum likelihood tree constructed based on the whole genome sequences

398  showed an admixture of *P. vivax* isolates with single and multiple *PvDBP* copy number

399  (Figure 6A). The Ethiopian *P. vivax* isolates were divided into six subclades. Subclade I

400  contained *P. vivax* samples mostly from Arbaminch and Badowacho with both one and

401  two *PvDBP* copies. Subclade II contained samples from Jimma and Hawassa with two

402  *PvDBP* copies. Subclade III contained a mixture of *P. vivax* samples from Arbaminch,

403  Halaba, Hawassa, and Jimma with single and high-order *PvDBP* copies. This clade was

404  sister to subclade IV that contained *P. vivax* samples mostly from Jimma (Figure 6A). In

405  subclade IV, no distinct clusters were detected between isolates with single and multiple

406  *PvDBP*.  Subclade V contained samples from Jimma and subclade VI contained

407  samples from Arbaminch, Badowacho, Hawassa, and Halaba. Each of the subclades

408  had samples with both one and two *PvDBP* copies.  Similar patterns were observed in

409  the *MSP*3 and *PIR* genes where *P. vivax* isolates with single and multiple copies were

410   clustered together in separate subclades (Figures 6B-D), suggesting that these gene

411   regions could have expanded multiply among samples at different locations.

412

413   **Gene flow and transmission network of the Ethiopian *P. vivax***

414   The principal component analysis based on the SNP variants showed samples from

415   Arbaminch, Badowacho, Hawassa, and Halaba were genetically closely related but

416   differentiated from Jimma (Figure 7A). The transmission network indicated that

417   Arbaminch was the major source or hub of infections where the infections in Jimma,

418   Hawassa, Badowacho, and Halaba were originated from (Table 3; Figure 7B). On the

419   other hand, no transmission was originated from Halaba, making this location the

420   largest sink of transmissions. The greatest extent of gene flow was observed between

421   Arbaminch and Badowacho (Figure 7B). Hawassa and Jimma showed a source hub

422   ratio of 0.5, indicating that there are equally as many egress transmissions as ingress

423   transmissions (Table 3). Although Jimma and Badowacho/Halaba are in close

424   geographical proximity, no apparent gene flow was observed between these sites.

425

426   **Discussion**

427   Across the genome, the total number of SNPs observed among 44 *P. vivax* isolates in

428   Ethiopia were comparable to those previously reported in South American [77] and

429   Southeast Asian countries [19]. For instance, 303,616 high-quality SNPs were detected

430   in 228 *P. vivax* isolates from Southeast Asia and Oceania in a previous study, of which

431   Sal-I was used as the reference sequence and subtelomeric regions were discarded

432   [19]. Auburn *et al.* [20] found that the average nucleotide diversity in Ethiopia was lower

433    than in Thailand and Indonesia, but higher than in Malaysia. Chromosomes 3, 4, and 5

434    have been previously shown to contain the lowest proportion of synonymous SNPs than

435    the other parts of the genome [12]. In the present study, chromosomes 3 and 6 were

436    found to have the lowest number of both synonymous and nonsynonymous SNPs. This

437    follows observations made in other studies done with nucleotide diversity ranging from

438    0.8 SNPs per kb in North Korea to 0.59 SNPs per kb in Peru [78]. Among the 64

439    erythrocyte binding gene candidates, the MSP and RBP multigene families showed the

440    highest level of genetic variation. This agrees with previous studies that reported a

441    remarkably high diversity in *RBP*2 than in *RBP*1 and its homolog group in *P. falciparum*

442    [31]. In the Greater Mekong Subregion, the *MSP*3 and *PIR* gene families also indicated

443    high levels of genetic diversity with 1.96% and 1.92% SNPs per base respectively,

444    confirming that members of multigene families are highly variable genetically [30, 79].

445    Such diversity suggested that the binding domains of these genes could be under

446    differential selection pressure. This pattern has been observed in previous studies and

447    is likely due to their critical role in reticulocyte invasion, immunogenic properties, and

448    human migration [26, 80-82].

449        Both CNVnator and GATK4 showed high order copies in several *PIR* gene

450    regions.  In addition, the *PIR* and *STP*1 genes were also indicated with significant

451    selection based on the iHS calculations. The *PIR* gene family, which includes *STP*1, are

452    located on the subtelomere regions and is a highly variable multigene family ranging

453    from 1,200 genes in the reference strain PvP01 to 346 genes in monkey-adapted strain

454    Salvador-I [56, 83]. Our analyses included only SNP variants that had a quality score of

455    40 or higher. Also, we used the PVP01 reference genome to map and annotate the

456    subtelomeric regions, with the goal to reflect variability and features across the entire

457    chromosome; whereas previous studies used the Sal-I reference genome with

458    hypervariable and subtelomeric regions removed to minimize mapping errors [19, 84].

459    A recent study in *P. chabaudi* suggested that polymorphisms in *PIR* genes could affect

460    the virulence of the parasites following passage from the mosquitoes [85]. Such a

461    variation in copy number of the *PIR* gene family has also been reported in *P. cynomolgi*

462    and *P. vivax* [86], suggesting that gene duplication could have been occurred

463    repeatedly in the ancestral lineages [86]. The *PIR* multigene family is one of the largest

464    gene families identified so far in *P. vivax* with several different potential functions.

465    Some *PIR* genes encode proteins on the surface of infected red blood cells, which could

466    confer to immune evasion; others encode proteins involved in signaling, trafficking and

467    adhesion functions [83].  Positive selection detected in the *PIR* genes among the

468    Ethiopian *P. vivax* isolates may have important implications on the susceptibility of the

469    mosquito hosts [87].

470         For the *P. vivax* isolates in Southeast Asia, copy number variation was observed

471    in nine gene regions including *DBP*1, *MDR*1, and PVX_101445 (on chromosome 14)

472    with copy number ranging from 3 to 4 [19]. *DBP*1 and *MSP*3 showed higher order

473    copies when compared to other genomic regions. In this study, the highest and most

474    variable copy number variations were detected in the *DBP*1, with copy numbers ranging

475    from one to as high as five.  Likewise, for the *MSP*3, copy numbers ranging from one to

476    as high as four. Based on the phylogeny, *DBP*1 and *MSP*3 expansion had occurred

477    multiple times as tandem copies.  These findings were consistent with earlier studies

478    [19, 76] and suggested that gene expansion may play a key role in host cell invasion

479    [88]. For all other putative erythrocyte binding genes, only a single copy was detected

480    among all samples. A larger sample in future investigations would verify this

481    observation.

482    In the present study, we identified a panel of 64 putative erythrocyte binding gene

483    candidates based on the information from the literature and analyzed their

484    polymorphisms. However, we did validate the function for each of these genes.  Among

485    these 64 putative erythrocyte binding gene candidates, *MAEBL* was shown to be highly

486    conserved in *Plasmodium* [89], had the highest signal for positive selection among the

487    *P. vivax* samples in Ethiopia. In *P. berghei*, *MAEBL* is a sporozoite attachment protein

488    that plays a role in binding and infecting the mosquito salivary gland [89]. In *P.*

489    *falciparum*, *MAEBL* is located in the rhoptries and on the surface of mature merozoites,

490    and expresses at the beginning of schizogony [89]. In *P. vivax*, *MAEBL* is a conserved

491    antigen expressed in blood stages, as well as in the mosquito midgut and salivary gland

492    sporozoites [89, 90]. The *MAEBL* antigen contains at least 25 predicted B-cell epitopes

493    that are likely to elicit antibody-dependent immune responses [91]. Positive selection

494    observed in this gene region among the Ethiopian *P. vivax* isolates could be associated

495    with the immunity-mediated selection pressure against blood-stage antigens. Though

496    *DBP*1 had the highest and most diverse copy number variation, no significant signal of

497    positive selection was detected.

498    It is noteworthy that the calculation of integrated haplotype scores and the

499    accuracy of phasing genotypes using BEAGLE were dependent on the levels of linkage

500    disequilibrium of the whole genomes. The higher the levels of linkage disequilibrium, the

501    more accurate are the phased genotypes and thus the iHS score. Pearson *et al*. [19]

502    found that *P. vivax* experienced drops in linkage disequilibrium after correcting for

503    population structure and other confounders.  Linkage disequilibrium of *P. vivax*

504    genomes has been previously shown to be associated with the rate of genetic

505    recombination and transmission intensity [92-94]. In high transmission sites of Papua

506    New Guinea and the Solomon Islands, no identical haplotypes and no significant

507    multilocus LD were observed, indicating limited inbreeding and random associations

508    between alleles in the parasite populations [95, 96]. However, when transmission

509    intensity declined, similar haplotypes and significant LD were observed possibly due to

510    self-fertilization, inbreeding and/or recombination of similar parasite strains

511    [92]. Multilocus LD is significantly associated with the genetic relatedness of the

512    parasite strains [97], but inversely associated with the proportion of polyclonal infections

513    [98]. In Southwestern Ethiopia, malaria transmission ranged from low to moderate, and

514    LD levels varied markedly among the study sites [53, 99]. To address this limitation in

515    BEAGLE, all genes that were detected with positive selection in BEAGLE were further

516    analyzed with PAML for verification. Future study should include broad samples to

517    thoroughly investigate selection pressure at the population level and the function

518    significance of polymorphisms in the *MAEB*L and *PIR* genes.

519         Previous studies have shown high levels of genetic diversity among *P. vivax*

520    isolates in endemic countries [16, 100, 101]. Such a diversity was directly related to high

521    transmission intensity and/or frequent gene exchange between parasite populations via

522    human movement [4, 12 , 13, 53]. For example, previous studies using microsatellites

523    have demonstrated a consistently high level of intra-population diversity ($H_E$ = 0.83) but

524    low between-population differentiation ($F_{ST}$ ranged from 0.001-0.1] in broader regions of

525    Ethiopia [53, 99]. High heterozygosity was also observed in *P. vivax* populations from

526    Qatar, India, and Sudan (average $H_{E} = 0.78$; 62), with only slight differentiation from *P.*

527    *vivax* in Ethiopia ($F_{ST}$ = 0.19) [102]. Frequent inbreeding among dominant clones [92,

528    95] and strong selective pressures especially in relapse infections [19, 20, 102, 103]

529    may also contribute to close genetic relatedness between and within populations. Thus,

530    in this study, it is not surprising to detect a high level of parasite gene flow among the

531    study sites at a small geographical scale, despite the limited number of samples. In the

532    present study, we successfully employed a transmission network model to identify

533    transmission paths, as well as the source and sink of infections in the region, beyond

534    simply indicating genetic relationships.

535         To conclude, this study elaborated on the genomic features of *P. vivax* in

536    Ethiopia, particularly focusing polymorphisms in erythrocyte binding genes that

537    potentially play a key role in local parasite invasion, a critical question given the mixed

538    Duffy positive and negative populations of Ethiopia. The findings provided baseline

539    information on the genomic variability of *P. vivax* infections in Ethiopia and allowed us to

540    compare the genomic variants of *P. vivax* between Duffy-positive and Duffy-negative

541    individuals as the next step of our ongoing investigation. Further, we are in progress of

542    developing a panel of informative SNP markers to track transmission at a micro-

543    geographical scale.

544    

## Data Availability

546    Additional information is provided as supplementary data accompanies this paper.

547    Sequence data of this study are deposited in the European Nucleotide Archive (ENA)

548    and the accession number of each sample is listed in Table 1.

549

## Acknowledgements

551    We are greatly indebted to the staffs and technicians from Jimma University for field

552    sample collection, the communities and hospitals for their support and willingness to

553    participate in this research.

554

## Funding

556    This research was funded by National Institutes of Health (NIH R15 AI138002 to EL;

557    NIH U19 AI129326 to GY; NIH R01 AI050243 to GY; D43 TW001505 to GY) and The

558    Wellcome Trust 206194/Z/17/Z to JR. The funders had no role in study design, data

559    collection and analysis, decision to publish, or preparation of the manuscript.

560

## Competing interests

562    The authors have declaredthat no competing interests exist.

563

## References

565    1. World Health Organization. World Malaria Report 2018. WHO, Geneva.

566    2. White MT, Shirreff G, Karl S, Ghani AC, Mueller I. Variation in relapse frequency and

567    the transmission potential of *Plasmodium vivax* malaria. Proc Biol Sci. 2016;283:

568    20160048.

569    3. Hemingway J, Shretta R, Wells TNC, Bell D, Djimdé AA, Achee N, et al. Tools and

570    strategies for malaria control and elimination: what do we need to achieve a grand

571    convergence in malaria? PLOS Biol. 2016;14: e1002380.

572    4. Hupalo DN, Luo Z, Melnikov A, Sutton PL, Rogov P, Escalante A, et al. Population

573        genomics studies identify signatures of global dispersal and drug resistance in

574        *Plasmodium vivax*. Nat Genet. 2016;48: 953–8.

575    5. Vallejo AF, Martinez NL, Tobon A, Alger J, Lacerda MV, Kajava AV, et al. Global

576        genetic diversity of the *Plasmodium vivax* transmission-blocking vaccine candidate

577        Pvs48/45. Malar J. 2016;15: 202.

578    6. Tham WH, Beeson JG, Rayer JC. *Plasmodium vivax* vaccine research - we've only

579        just begun. Int J Parasitol. 2017;47: 111-118.

580    7. Kale S, Yadav CP, Rao PN, Shalini S, Eapen A, Srivasatava HC, Sharma SK, Pande

581        V, Carlton JM, Singh OP, Mallick PK. Antibody responses within two leading

582        *Plasmodium vivax* vaccine candidate antigens in three geographically diverse

583        malaria-endemic regions of India. Malar J. 2019;18: 425.

584    8. Baniecki ML, Faust AL, Schaffner SF, Park DJ, Galinsky K, Daniels RF, et al.

585        Development of a single nucleotide polymorphism barcode to genotype *Plasmodium*

586        *vivax* infections. PLoS Negl Trop Dis. 2015;9: e0003539.

587    9. Diez Benavente E, Campos M, Phelan J, Nolder D, Dombrowski JG, Marinho CRF, et

588        al. A molecular barcode to inform the geographical origin and transmission dynamics

589        of *Plasmodium vivax* malaria. PLoS Genet. 2020;16: e1008576.

590    10. Parobek CM, Lin JT, Saunders DL, Barnett EJ, Lon C, Lanteri CA, et al. Selective

591        sweep suggests transcriptional regulation may underlie *Plasmodium vivax* resilience

592        to malaria control measures in Cambodia. Proc. Natl. Acad. Sci. U.S.A. 2016;113: 50.

593    11. Auburn S, Benavente ED, Miotto O, Pearson RD, Amato R, Grigg MJ, et al.

594        Genomic analysis of a pre-elimination Malaysian *Plasmodium vivax* population

595        reveals selective pressures and changing transmission dynamics. Nat Commun.

596        2018;9: 2585.

597    12. Benavente ED, Ward Z, Chan W, Mohareb FR, Sutherland CJ, Roper C, et al.

598        Genomic variation in *Plasmodium vivax* malaria reveals regions under selective

599        pressure. PLoS One. 2017;12: 5.

600    13. Lima-Junior JDC, Pratt-Riccio LR. Major histocompatibility complex and malaria:

601        focus on *Plasmodium vivax* Infection. Front. Immunol. 2016;7: 13.

602    14. Kano FS, Souza AMD, Torres LDM, Costa MA, Souza-Silva FA, Sanchez BAM, et

603       al. Susceptibility to *Plasmodium vivax* malaria associated with DARC (Duffy antigen)

604       polymorphisms is influenced by the time of exposure to malaria. Sci. Rep. 2018;8:

605       13851.

606    15. Ventocilla JA, Nuñez J, Tapia LL, Lucas CM, Manock SR, Lescano AG, et al.

607       Genetic variability of *Plasmodium vivax* in the north coast of Peru and the Ecuadorian

608       amazon basin. Am. J. Trop. Med. Hyg. 2018;99: 27–32.

609    16. Fola AA, Harrison GLA, Hazairin MH, Barnadas C, Hetzel MW, Iga J, et al. Higher

610       complexity of infection and genetic diversity of *Plasmodium vivax* than *Plasmodium*

611       *falciparum* across all malaria transmission zones of Papua New Guinea. Am. J. Trop.

612       Med. Hyg. 2017;16–0716.

613    17. Barry AE, Waltmann A, Koepfli C, Barnadas C, Mueller I. Uncovering the

614       transmission dynamics of *Plasmodium vivax* using population genetics. Pathog. Glob.

615       Health. 2015;109: 142–52.

616    18. Koepfli C, Rodrigues PT, Antao T, Orjuela-Sánchez P, Eede PVD, Gamboa D, et al.

617       *Plasmodium vivax* diversity and population structure across four continents. PLoS

618       Negl. Trop. Dis. 2015;9: e0003872.

619    19. Pearson RD, Amato R, Auburn S, Miotto O, Almagro-Garcia J, Amaratunga C, et al.

620       Genomic analysis of local variation and recent evolution in *Plasmodium vivax*. Nat

621       Genet. 2016;48: 959–64.

622    20. Auburn S, Getachew S, Pearson RD, Amato R, Miotto O, Trimarsanto H, et al.

623       Genomic analysis of *Plasmodium vivax* in southern Ethiopia reveals selective

624       pressures in multiple parasite mechanisms. J. Infect. Dis. 2019;220: 1738–49.

625    21. Costa GL, Amaral LC, Fontes CJF, Carvalho LH, Brito CFAD, Sousa TND.

626       Assessment of copy number variation in genes related to drug resistance in

627       *Plasmodium vivax* and *Plasmodium falciparum* isolates from the Brazilian Amazon

628       and a systematic review of the literature. Malar. J. 2017;16: 152.

629    22. Lin JT, Muth S, Rogers WO, Ubalee R, Kharabora O, Juliano JJ, et al. *Plasmodium*

630       *vivax* isolates from Cambodia and Thailand show high genetic complexity and distinct

631       patterns of *P. vivax* multidrug resistance gene 1 (pvmdr1) polymorphisms. Am. J.

632       Trop. Med. Hyg. 2013;88: 1116–23.

633   23. Cornejo OE, Fisher D, Escalante AA. Genome-wide patterns of genetic
634       polymorphism and signatures of selection in *Plasmodium vivax*. Genome Biol. Evol.
635       2014;7: 106–19.

636   24. Chen E, Salinas ND, Huang Y, Ntumngia F, Plasencia MD, Gross ML, et al. Broadly
637       neutralizing epitopes in the *Plasmodium vivax* vaccine candidate duffy binding
638       protein. Proc. Natl. Acad. Sci. U.S.A. 2016;113: 6277–82.

639   25. Singh V, Gupta P, Pande V. Revisiting the multigene families: *Plasmodium* var and
640       vir genes. J Vector Borne Dis. 2014;51: 75–81.

641   26. Rice BL, Acosta MM, Pacheco MA, Carlton JM, Barnwell JW, Escalante AA. The
642       origin and diversification of the merozoite surface protein 3 (msp3) multi-gene family
643       in *Plasmodium vivax* and related parasites. Mol. Phylogenetics Evol. 2014;78: 172–
644       84.

645   27. Lu F, Li J, Wang B, Cheng Y, Kong D-H, Cui L, et al. Profiling the humoral immune
646       responses to *Plasmodium vivax* infection and identification of candidate immunogenic
647       rhoptry-associated membrane antigen (RAMA). J. Proteom. 2014;102: 66–82.

648   28. Rahul C, Krishna KS, Pawar AP, Bai M, Kumar V, Phadke S, et al. Genetic and
649       structural characterization of PvSERA4: potential implication as therapeutic target for
650       *Plasmodium vivax* malaria. J. Biomol. Struct. Dyn. 2013;32: 580–90.

651   29. Rahul C, Krishna KS, Meera M, Phadke S, Rajesh V. *Plasmodium vivax*: N-terminal
652       diversity in the blood stage SERA genes from Indian isolates. Blood Cells Mol. Dis.
653       2015;55: 30–5.

654   30. Chen S-B, Wang Y, Kassegne K, Xu B, Shen H-M, Chen J-H. Whole-genome
655       sequencing of a *Plasmodium vivax* clinical isolate exhibits geographical
656       characteristics and high genetic variation in China-Myanmar border area. BMC
657       Genom. 2017;18: 131.

658   31. Rayner JC, Corredor V, Tran TM, Barnwell JW, Huber CS, Galinski MR. Dramatic
659       difference in diversity between *Plasmodium falciparum* and *Plasmodium Vivax*
660       reticulocyte binding-like genes. Am. J. Trop. Med. Hyg. 2005;72: 666–74.

661   32. Gunalan K, Niangaly A, Thera MA, Doumbo OK, Miller LH. *Plasmodium vivax*
662       infections of duffy-negative erythrocytes: historically undetected or a recent
663       adaptation? Trends Parasitol. 2018;34: 420–9.

664    33. Luo Z, Sullivan SA, Carlton JM. The biology of *Plasmodium vivax* explored through
665        genomics. Ann. N. Y. Acad. Sci. 2015;1342: 53–61.
666    34. Gruszczyk J, Huang RK, Chan L-J, Menant S, Hong C, Murphy JM, et al. Cryo-EM
667        structure of an essential *Plasmodium vivax* invasion complex. Nature. 2018;559:
668        135–9.
669    35. Chan LJ, Dietrich MH, Nguitragool W, Tham WH. *Plasmodium vivax* reticulocyte
670        binding proteins for invasion into reticulocytes. Cell. Microbiol. 2019; e13110.
671    36. Moreno-Pérez DA, Baquero LA, Chitiva-Ardila DM, Patarroyo MA. Characterising
672        PvRBSA: an exclusive protein from *Plasmodium* species infecting reticulocytes.
673        Parasites Vectors. 2017;10: 243.
674    37. Camargo-Ayala PA, Garzón-Ospina D, Moreno-Pérez DA, Ricaurte-Contreras LA,
675        Noya O, Patarroyo MA. On the evolution and function of *Plasmodium vivax*
676        reticulocyte binding surface antigen (pvrbsa). Front. Genet. 2018;9: 372.
677    38. Roesch C, Popovici J, Bin S, Run V, Kim S, Ramboarina S, et al. Genetic diversity
678        in two *Plasmodium vivax* protein ligands for reticulocyte invasion. PLOS Negl. Trop.
679        Dis. 2018;12: e0006555.
680    39. Hester J, Chan ER, Menard D, Mercereau-Puijalon O, Barnwell J, Zimmerman PA,
681        et al. *De-novo* assembly of a field isolate genome reveals novel *Plasmodium vivax*
682        erythrocyte invasion genes. PLOS Negl. Trop. Dis. 2013;7: e2569.
683    40. Ntumngia FB, Thomson-Luque R, Torres LDM, Gunalan K, Carvalho LH, Adams
684        JH. A novel erythrocyte binding protein of *Plasmodium vivax* suggests an alternate
685        invasion pathway into duffy-positive reticulocytes. mBio. 2016;7: e01261-16.
686    41. Carias LL, Dechavanne S, Nicolete VC, Sreng S, Suon S, Amaratunga C, et al.
687        Identification and characterization of functional human monoclonal antibodies to
688        *Plasmodium vivax* duffy-binding protein. J. Immunol. 2019;202: 2648–60.
689    42. He WQ, Shakri AR, Bhardwaj R, Franca CT, Stanisic DI, Healer J, et al. Antibody
690        responses to *Plasmodium vivax* duffy binding and erythrocyte binding proteins predict
691        risk of infection and are associated with protection from clinical Malaria. PLOS Negl.
692        Trop. Dis. 2019;13: e0006987.

693  43. Wang B, Lu F, Cheng Y, Chen J-H, Jeon H-Y, Ha K-S, et al. Immunoprofiling of the
694      tryptophan-rich antigen family in *Plasmodium vivax*. Infect. Immun. 2015;83: 3083–
695      95.

696  44. Baquero LA, Moreno-Pérez DA, Garzón-Ospina D, Forero-Rodríguez J, Ortiz-
697      Suárez HD, Patarroyo MA. PvGAMA reticulocyte binding activity: predicting
698      conserved functional regions by natural selection analysis. Parasites Vectors.
699      2017;10: 251.

700  45. Arévalo-Pinzón G, Bermúdez M, Curtidor H, Patarroyo MA. The *Plasmodium vivax*
701      rhoptry neck protein 5 is expressed in the apical pole of *Plasmodium vivax* VCG-1
702      strain schizonts and binds to human reticulocytes. Malar. J. 2015;14: 106.

703  46. Tyagi K, Hossain ME, Thakur V, Aggarwal P, Malhotra P, Mohmmed A, et al.
704      *Plasmodium vivax* tryptophan rich antigen PvTRAg36.6 interacts with PvETRAMP
705      and PvTRAg56.6 interacts with PvMSP7 during erythrocytic stages of the parasite.
706      Plos One. 2016;11: e0151065.

707  47. Gunalan K, Sá JM, Barros RRM, Anzick SL, Caleon RL, Mershon JP, et al.
708      Transcriptome profiling of *Plasmodium vivax* in Samira monkeys identifies potential
709      ligands for invasion. Proc. Natl. Acad. Sci. U.S.A. 2019;116: 7053–61.

710  48. Zimmerman PA. *Plasmodium vivax* infection in duffy-negative people in Africa. Am.
711      J. Trop. Med. Hyg. 2017;97: 636–8.

712  49. Battle KE, Lucas TCD, Nguyen M, Howes RE, Nandi AK, Twohig KA, et al. Mapping
713      the global endemicity and clinical burden of *Plasmodium vivax*, 2000–17: a spatial
714      and temporal modelling study. Lancet. 2019;394: 332–43.

715  50. Twohig KA, Pfeffer DA, Baird JK, Price RN, Zimmerman PA, Hay SI, et al. Growing
716      evidence of *Plasmodium vivax* across malaria-endemic Africa. PLOS Negl. Trop. Dis.
717      2019;13: e0007140.

718  51. Auburn S, Campino S, Clark TG, Djimde AA, Zongo I, Pinches R, et al. An Effective
719      Method to Purify Plasmodium falciparum DNA Directly from Clinical Blood Samples
720      for Whole Genome High-Throughput Sequencing. Plos One. 2011;6: e22213.

721  52. Lo E, Yewhalaw D, Zhong D, Zemene E, Degefa T, Tushune K, et al. Molecular
722      epidemiology of *Plasmodium vivax* and *Plasmodium falciparum* malaria among duffy-
723      positive and duffy-negative populations in Ethiopia. Malar. J. 2015;14: 84.

724   53. Lo E, Hemming-Schroeder E, Yewhalaw D, Nguyen J, Kebede E, Zemene E, et al.

725       Transmission dynamics of co-endemic *Plasmodium vivax* and *P. falciparum* in

726       Ethiopia and prevalence of antimalarial resistant genotypes. PLOS Negl. Trop. Dis.

727       2017;11: e0005806.

728   54. Langmead B, Salzberg SL. Fast gapped-read alignment with bowtie 2. Nat.

729       Methods. 2012;9: 357–9.

730   55. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler

731       transform. Bioinformatics. 2009;25:1754–60.

732   56. Auburn S, Böhme U, Steinbiss S, Trimarsanto H, Hostetler J, Sanders M, et al. A

733       new *Plasmodium vivax* reference sequence with improved assembly of the

734       subtelomeres reveals an abundance of pir genes. Wellcome Open Res. 2016;1: 4.

735   57. Li H. A statistical framework for SNP calling, mutation discovery, association

736       mapping and population genetical parameter estimation from sequencing data.

737       Bioinformatics. 2011;27: 2987–93.

738   58. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for

739       annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. Fly.

740       2012;6: 80–92.

741   59. Cingolani P, Patel VM, Coon M, Nguyen T, Land SJ, Ruden DM, et al. Using

742       *Drosophila melanogaster* as a model for genotoxic chemical mutational studies with a

743       new program, SnpSift. Front. Genet. 2012;3: 35.

744   60. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, Depristo MA, et al. The

745       variant call format and VCFtools. Bioinformatics. 2011;27: 2156–8.

746   61. Abyzov A, Urban AE, Snyder M, Gerstein M. CNVnator: An approach to discover,

747       genotype, and characterize typical and atypical CNVs from family and population

748       genome sequencing. Genome Res. 2011;21: 974–84.

749   62. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,

750       Durbin R, 1000 Genome Project Data Processing Subgroup. The sequence

751       alignment/map format and SAMtools. Bioinformatics. 2009;25: 2078–2079.

752   63. Mckenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The

753       genome analysis toolkit: a mapreduce framework for analyzing next-generation DNA

754       sequencing data. Genome Res. 2010;20: 1297–303.

755    64. Depristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A

756        framework for variation discovery and genotyping using next-generation DNA

757        sequencing data. Nat Genet. 2011;43: 491–8.

758    65. Auwera GAVD, Carneiro MO, Hartl C, Poplin R, Angel GD, Levy-Moonshine A, et al.

759        From fastQ data to high-confidence variant calls: the genome analysis toolkit best

760        practices pipeline. Curr Protoc Bioinformatics. 2013;43: 1-11.

761    66. Miles A, Harding N. cggh/scikit-allel: v1.1.8 [Internet]. Zenodo. 2017. Available from:

762        https://zenodo.org/record/822784#.XKIe6yhKiUk

763    67. Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing

764        data inference for whole genome association studies by use of localized haplotype

765        clustering. Am J Hum Genet. 2007;81: 1084-97.

766    68. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. Mol. Biol. Evol.

767        2007;24: 1586–91.

768    69. Prajapati SK, Singh OP. Insights into the invasion biology of *Plasmodium vivax*.

769        Front. Cell. Infect. Microbiol. 2013;3: 8.

770    70. Rozas J. DNA sequence polymorphism analysis using DnaSP. Methods Mol. Biol.

771        Bioinformatics DNA Seq. Anal. 2009;537: 337–50.

772    71. Katoh K, Standley DM. MAFFT Multiple sequence alignment software version 7:

773        improvements in performance and usability. Mol. Biol. Evol. 2013;30: 772–80.

774    72. Kück P, Longo GC. FASconCAT-G: extensive functions for multiple sequence

775        alignment preparations concerning phylogenetic studies. Front. Zool. 2014;11: 81.

776    73. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses

777        with thousands of taxa and mixed models. Bioinformatics 2006;22: 2688–90.

778    74. Jombart T. Adegenet: a R package for the multivariate analysis of genetic markers.

779        Bioinformatics. 2008;24: 1403–5.

780    75. Schneider ADB, Ford CT, Hostager R, Williams J, Cioce M, Çatalyürek ÜV, et al.

781        StrainHub: a phylogenetic tool to construct pathogen transmission networks.

782        Bioinformatics. 2019; btz646.

783    76. Lo E, Hostetler JB, Yewhalaw D, Pearson RD, Hamid MMA, Gunalan K, et al.

784        Frequent expansion of *Plasmodium vivax* duffy binding protein in Ethiopia and its

785        epidemiological significance. PLOS Negl. Trop. Dis. 2019;13: e0007222.

77. Oliveira TCD, Rodrigues PT, Menezes MJ, Gonçalves-Lopes RM, Bastos MS, Lima NF, et al. Genome-wide diversity and differentiation in new world populations of the human malaria parasite *Plasmodium vivax*. PLOS Negl. Trop. Dis. 2017;11: e0005824.

78. Neafsey DE, Galinsky K, Jiang RHY, Young L, Sykes SM, Saif S, et al. The malaria parasite *Plasmodium vivax* exhibits greater genetic diversity than *Plasmodium falciparum*. Nat Genet. 2012;44: 1046–50.

79. Dharia NV, Bright AT, Westenberger SJ, Barnes SW, Batalov S, Kuhen K, et al. Whole-genome sequencing and microarray analysis of *ex vivo Plasmodium vivax* reveal selective pressure on putative drug resistance genes. Proc. Natl. Acad. Sci. USA 2010;107: 20045–50.

80. Chen J-H, Chen S-B, Wang Y, Ju C, Zhang T, Xu B, et al. An immunomics approach for the analysis of natural antibody responses to *Plasmodium vivax* infection. Mol Biosyst. 2015;11: 2354–63.

81. Shen H-M, Chen S-B, Wang Y, Xu B, Abe EM, Chen J-H. Genome-wide scans for the identification of *Plasmodium vivax* genes under positive selection. Malar. J. 2017;16: 238.

82. Mascorro CN, Zhao K, Khuntirat B, Sattabongkot J, Yan G, Escalante AA, et al. Molecular evolution and intragenic recombination of the merozoite surface protein MSP-3α from the malaria parasite *Plasmodium vivax* in Thailand. Parasitology. 2005;131: 25–35.

83. Cunning D, Lawton J, Jarra W, Preiser P, Langhorne J. The pir multigene family of *Plasmodium*: antigenic variation and beyond. Mol. Biochem. Parasitol. 2010;170: 65-73.

84. Loy DE, Plenderleithc LJ, Sundararamana SA, Liua W, Gruszczyke J, ChenG YJ, Trimbolia S, Learna GH, MacLeanc OA, Morgan ALK, Lia Y, Avittoa AN, Gilesa J, Calvignac-Spencerg S, Sachseg A, Leendertzg FH, Speedeh S, Ayoubai A, Peetersi M, Rayner JC, Tham WH, Sharp PM, Hahna BH. Evolutionary history of human *Plasmodium vivax* revealed by genome-wide analyses of related ape parasites. Proc Natl Acad Sci USA 2018;115: e8450-8459.

816   85. Spence PJ, Jarra W, Lévy P, Reid AJ, Chappell L, Brugat T, et al. Vector

817       transmission regulates immune control of *Plasmodium* virulence. Nature. 2013;498:

818       228–31.

819   86. Tachibana S-I, Sullivan SA, Kawai S, Nakamura S, Kim HR, Goto N, et al.

820       *Plasmodium cynomolgi* genome sequences provide insight into *Plasmodium vivax*

821       and the monkey malaria clade. Nat Genet. 2012;44: 1051–5.

822   87. Carlton JM, Adams JH, Silva JC, Bidwell SL, Lorenzi H, Caler E, et al. Comparative

823       genomics of the neglected human malaria parasite *Plasmodium vivax*. Nature.

824       2008;455: 757–63.

825   88. Urusova D, Carias L, Huang Y, Nicolete VC, Popovici J, Roesch C, et al. Structural

826       basis for neutralization of *Plasmodium vivax* by naturally acquired human antibodies

827       that target DBP. Nat. Microbiol. 2019;4: 1486–96.

828   89. Kariu T, Yuda M, Yano K, Chinzei Y. MAEBL is essential for malarial sporozoite

829       infection of the mosquito salivary gland. J. Exp. Med. 2002;195: 1317–23.

830   90. Preiser P, Renia L, Singh N, Balu B, Jarra W, Voza T, et al. Antibodies against

831       MAEBL ligand domains M1 and M2 inhibit sporozoite development *in-vitro*. Infect.

832       Immun. 2004;72: 3604–8.

833   91. Doolan DL, Dobano C, Baird JK. Acquired immunity to malaria. Clin. Microbiol. Rev.

834       2009;22: 13–36.

835   92. Waltmann A, Koepfli C, Tessier N, Karl S, Fola A, Darcy AW, et al. Increasingly

836       inbred and fragmented populations of *Plasmodium vivax* associated with the

837       eastward decline in malaria transmission across the Southwest Pacific. PLOS Negl.

838       Trop. Dis. 2018;12: e0006146.

839   93.  Eede PVD, D'Alessandro U, Erhart A, Thang ND, Anné J, Overmeir CV, et al. High

840       complexity of *Plasmodium vivax* infections in symptomatic patients from a rural

841       community in central Vietnam detected by microsatellite genotyping. Am. J. Trop.

842       Med. Hyg. 2010;82: 223–7.

843   94.  Liu Y, Auburn S, Cao J, Trimarsanto H, Zhou H, Gray K-A, et al. Genetic diversity

844       and population structure of *Plasmodium vivax* in central China. Malar. J. 2014;13:

845       262.

846  95. Jennison C, Arnott A, Tessier N, Tavul L, Koepfli C, Felger I, et al. *Plasmodium*

847  *vivax* populations are more genetically diverse and less structured than sympatric

848  *Plasmodium falciparum* populations. PLOS Negl. Trop. Dis. 2015;9: e0003634.

849  96. Koepfli C, Timinao L, Antao T, Barry AE, Siba P, Mueller I, et al. A large

850  *Plasmodium vivax* reservoir and little population structure in the South Pacific. PLoS

851  ONE. 2013;8: e66041.

852  97. Lin E, Kiniboro B, Gray L, Dobbie S, Robinson L, Laumaea A, et al. Differential

853  patterns of infection and disease with *P. falciparum* and *P. vivax* in young Papua New

854  Guinean children. PLoS ONE. 2010;5: e9047.

855  98. Koepfli C, Ross A, Kiniboro B, Smith TA, Zimmerman PA, Siba P, et al. Multiplicity

856  and diversity of *Plasmodium vivax* infections in a highly endemic region in Papua

857  New Guinea. PLOS Negl. Trop. Dis. 2011;5: e1424.

858  99. Getachew S, To S, Trimarsanto H, Thriemer K, Clark TG, Petros B, et al. Variation

859  in complexity of infection and transmission stability between neighboring populations

860  of *Plasmodium vivax* in southern Ethiopia. Plos One. 2015;10: e0140780.

861  100. Friedrich LR, Popovici J, Kim S, Dysoley L, Zimmerman PA, Menard D, et al.

862  Complexity of infection and genetic diversity in Cambodian *Plasmodium vivax*. PLOS

863  Negl. Trop. Dis. 2016;10: e0004526.

864  101. Lo E, Lam N, Hemming-Schroeder E, Nguyen J, Zhou G, Lee M-C, et al. Frequent

865  spread of *Plasmodium vivax* malaria maintains high genetic diversity at the Myanmar-

866  China border, without distance and landscape barriers. J. Infect. Dis. 2017;216:

867  1254–63.

868  102. Abdelraheem MH, Bansal D, Idris MA, Mukhtar MM, Hamid MMA, Imam ZS, et al.

869  Genetic diversity and transmissibility of imported *Plasmodium vivax* in Qatar and

870  three countries of origin. Sci. Rep. 2018;8: 8870.

871  103. Popovici J, Friedrich LR, Kim S, Bin S, Run V, Lek D, et al. Genomic analyses

872  reveal the common occurrence and complexity of *Plasmodium vivax* relapses in

873  Cambodia. mBio. 2018;9: e01888-17.

874

875

## Tables

**Table 1**. Information of whole genome sequences of 44 *Plasmodium vivax* isolates from Ethiopia. The European Nucleotide Archive (ENA) accession number for all files.

**Table 2.** A shortlist of 25 erythrocyte binding genes that showed signals of positive selection based on the Likelihood Ratio Test of the M1 (neutral model) and M2 models (selection model) in PAML.

**Table 3.** Transmission network metrics among study sites calculated by StrainHub.

## Figures

**Figure 1.** An overview of the *P vivax* sample collection locations including Arbaminch, Badowacho, Hawassa, Halaba, and Jimma in southwestern Ethiopia.

**Figure 2.** A summary representation of the *P. vivax* genome, with the outer ring as an ideogram representing the 14 nuclear chromosomes and sizes of each. The second track represented the average coverage for each chromosome among the 44 Ethiopian samples. The third track containing the gray vertical dashes represented the distribution of genes across the 14 chromosomes. The forth track that contained the red vertical lines represented the 64 erythrocyte binding gene candidates. The fifth inner track with the light blue background represented the $d_N/d_S$ ratio calculated by partitioning the chromosomes into genomic regions and $d_N/d_S$ directly. The three outliers (yellow dots) represented three unknown plasmodium protein genes that were

899    detected with significant positive selection.  The sixth track indicated the overall copy

900    number variation calculated using CNVnator.  Red dots represented genes with copy

901    number variation among the Ethiopian genomes.

902

903    **Figure 3.** (A) A distribution of the nonsynonymous and synonymous mutations of each

904    chromosome. A higher proportion of synonymous mutations was observed compared to

905    nonsynonymous mutations. Chromosomes 7, 9, and 12 have the most mutations

906    overall, with chromosomes 6 and 3 having the fewest number of mutations. (B) Number

907    of mutation sites and the nucleotide diversity of 64 erythrocyte binding genes.  The

908    *PvRBP* and *PvMSP* multigene families have the highest number of polymorphic sites

909    when compared to the others, with *PvRBP*2c the highest number of nonsynonymous

910    and synonymous mutations, followed by *PvMSP*3 and *PvMSP*1.  Approximately 40% of

911    the mutations were nonsynonymous.  These genes were also indicated with the highest

912    nucleotide diversity.

913

914    **Figure 4.** Signal of positive selection across the 14 chromosomes among all *P. vivax*

915    samples.  Genes that showed significant signal of positive selection included *STP*1,

916    *MAEBL*, *MSP*3.8, and *PIR* gene regions.  *PvMSP*3.8 gene may play a role in the

917    erythrocyte invasion. *MAEBL* is a membrane associated erythrocyte binding like protein

918    that may have a function associated with erythrocyte invasion.

919

920    **Figure 5.** A total of 28 gene regions that were detected with copy number variation.

921    Annotation of these genes can be found in Supplementary Table 4. Among them,

922    *PvDBP*1 (PVP01_0623800) and *PvMSP*3 (PVP01_1030900) were associated with

923    erythrocyte invasion.  Other genes that were found to have high-order copy number

924    were *PIR* protein genes or unknown exported plasmodium proteins.

925

926    **Figure 6**. An unrooted whole genome phylogenetic tree of the 44 Ethiopian samples

927    showing the evolution of (A) *PvDBP*; (B) *PvMSP*3; (C) *PIR* gene on chromosome 2; and

928    (D) *PIR* gene on chromosome 11. The Ethiopian isolates were divided into three

929    subclades. Subclade I contained samples mostly from the Arbaminch and Badowacho.

930    Subclade II contained a mixture of isolates from Arbaminch, Halaba, Hawassa and

931    Jimma. Subclade III contained samples from Jimma. No distinct clusters were observed

932    between isolates with single and multiple *PvDBP*, *PvMSP*3, and *PIR* genes. These

933    patterns suggest that these gene regions could have expanded multiply among samples

934    at different locations.

935

936    **Figure 7.** (A) Principal component analysis plot based on the SNP information from our

937    variant analysis. Samples obtained from Jimma were clustered together, whereas

938    samples from Arbaminch, Badowacho, Hawassa, and Halaba were mixed together with

939    the exception of two samples from Hawassa. This clustering pattern suggested that

940    there was considerable genetic variation among study sites even at a small

941    geographical scale. (B) The transmission network, created using the StrainHub

942    program, indicated that Arbaminch was the major source of infection in Jimma, Halaba,

943    Badowacho and Hawassa.  The greatest extent of gene flow (indicated by the boldest

944    arrow) was observed between Arbaminch and Badowacho. Even though Jimma,

945　Badowacho and Halaba are geographically in close proximity, gene flow was not
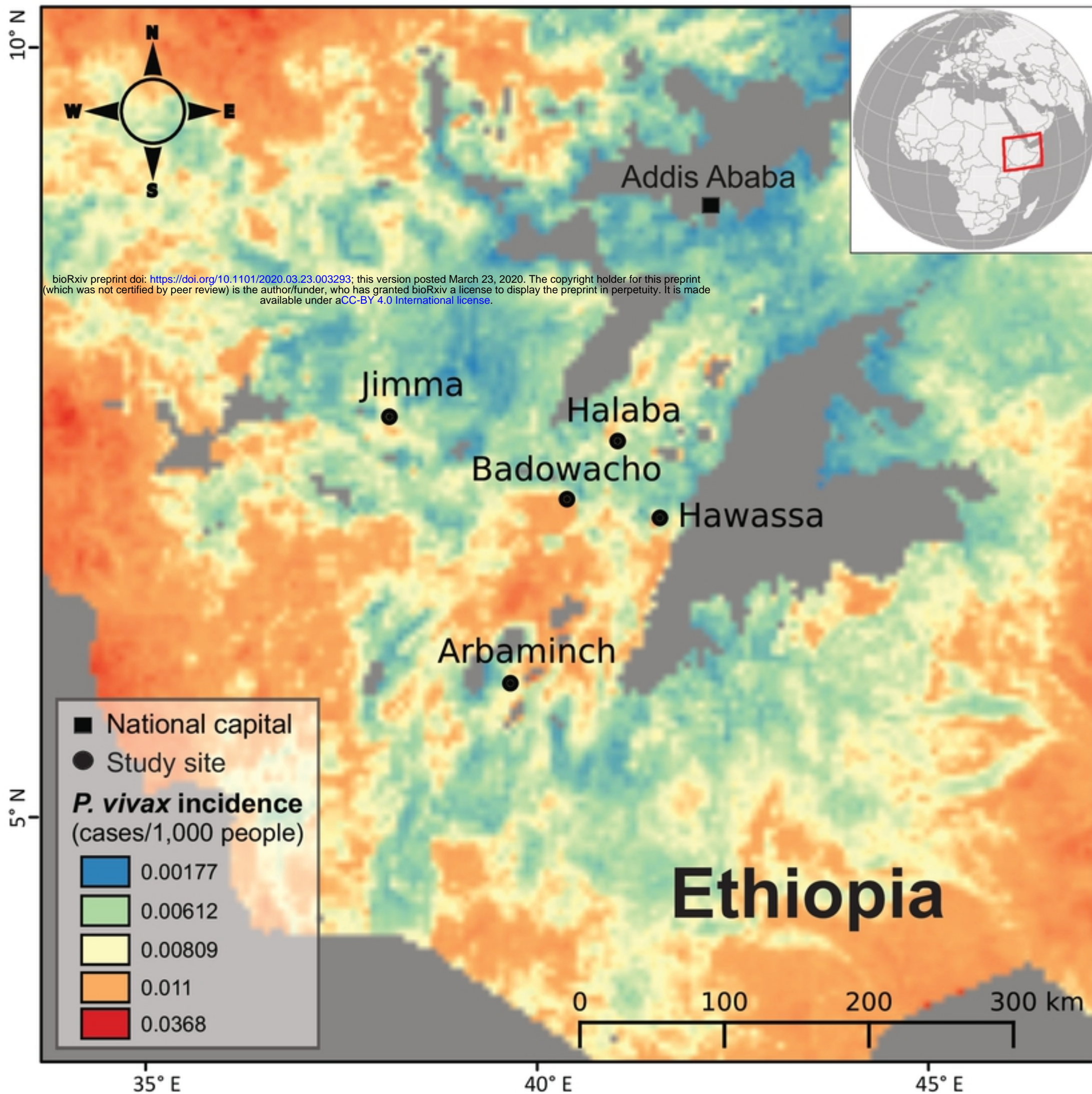
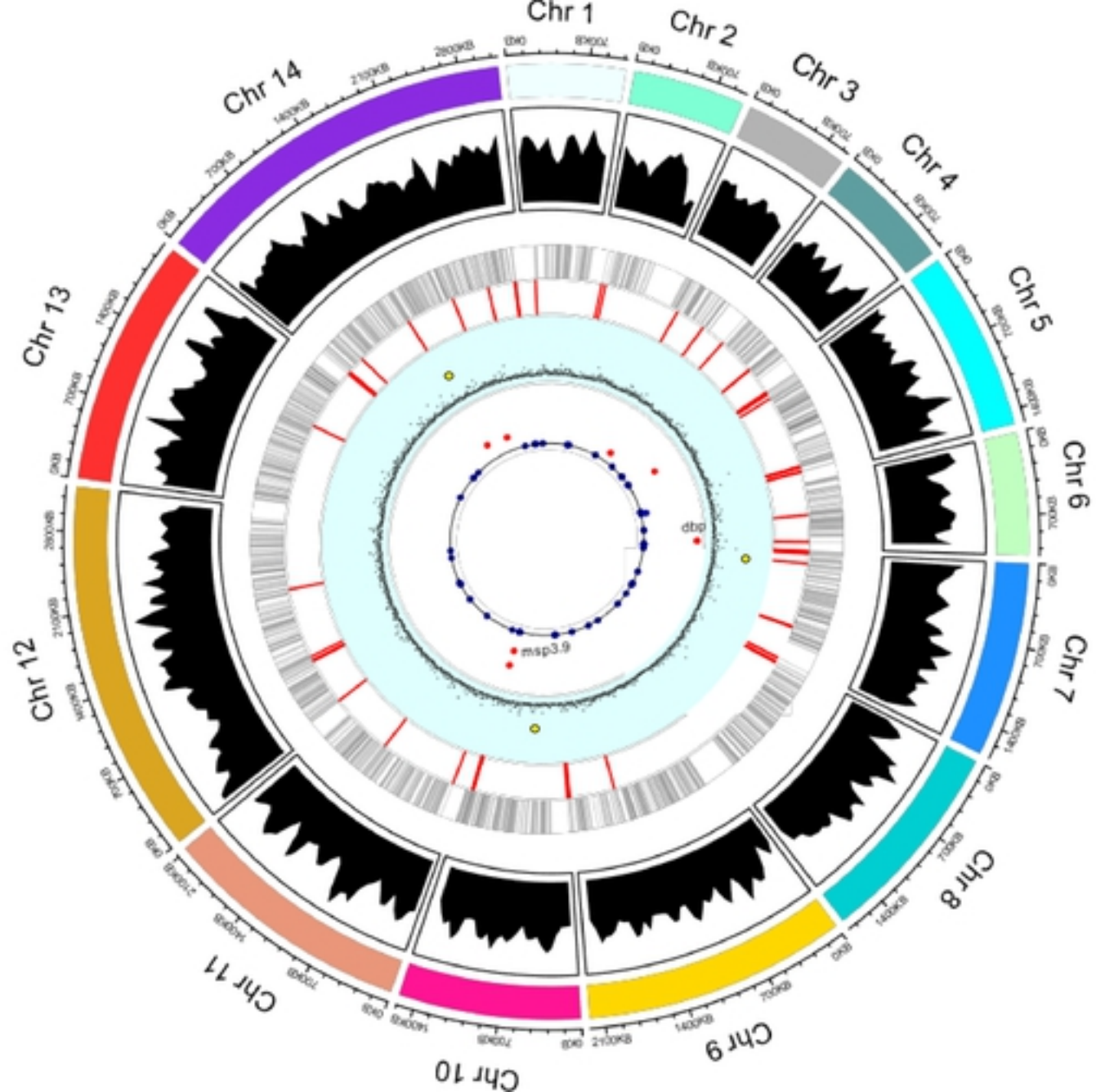946　intense among these sites.

947

948　**Supplementary files**

949　**Supplementary Table 1.** Distribution of SNP variants in the 64 *P. vivax* erythrocyte

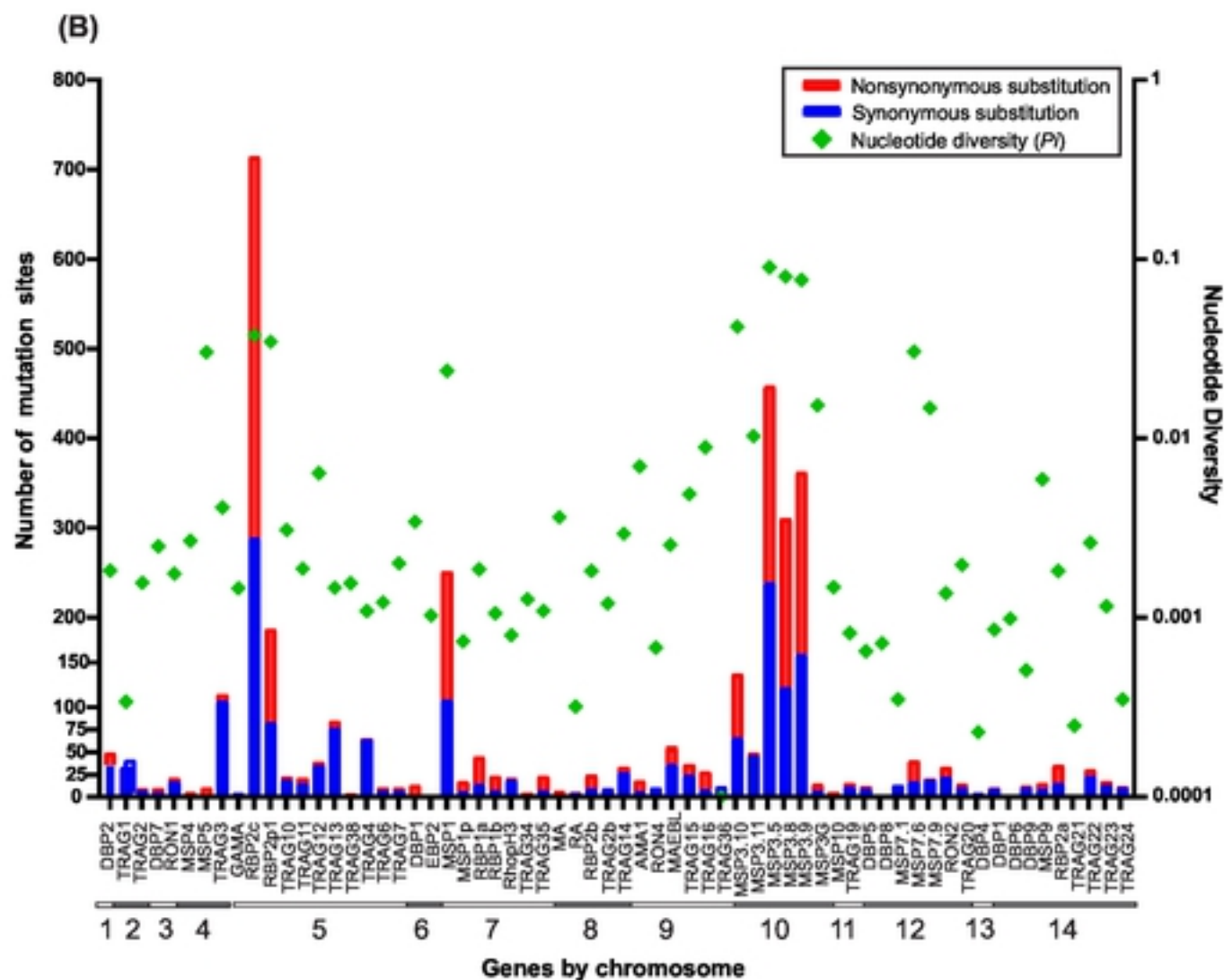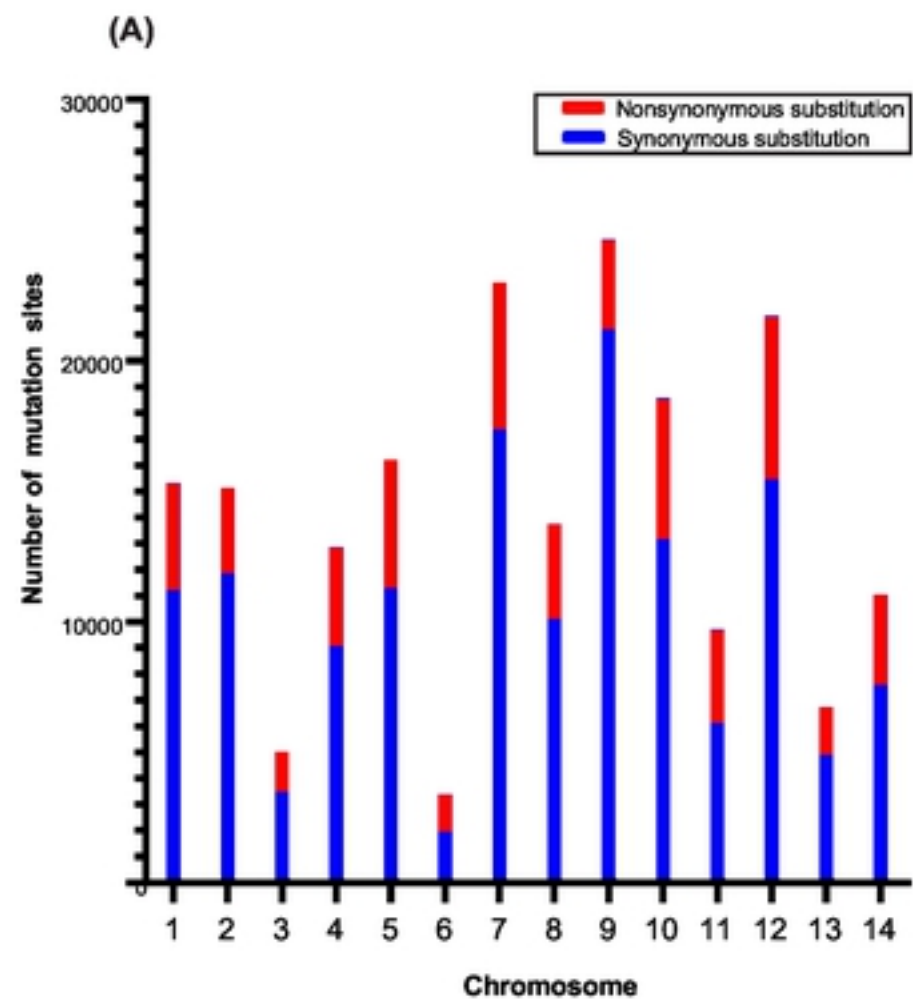950　binding gene candidates among the 44 Ethiopian genomes.

951

952　**Supplementary Table 2.** Distribution of single nucleotide polymorphism (SNP) variants

953　across *P. vivax* chromosomes of the 44 Ethiopian genomes.
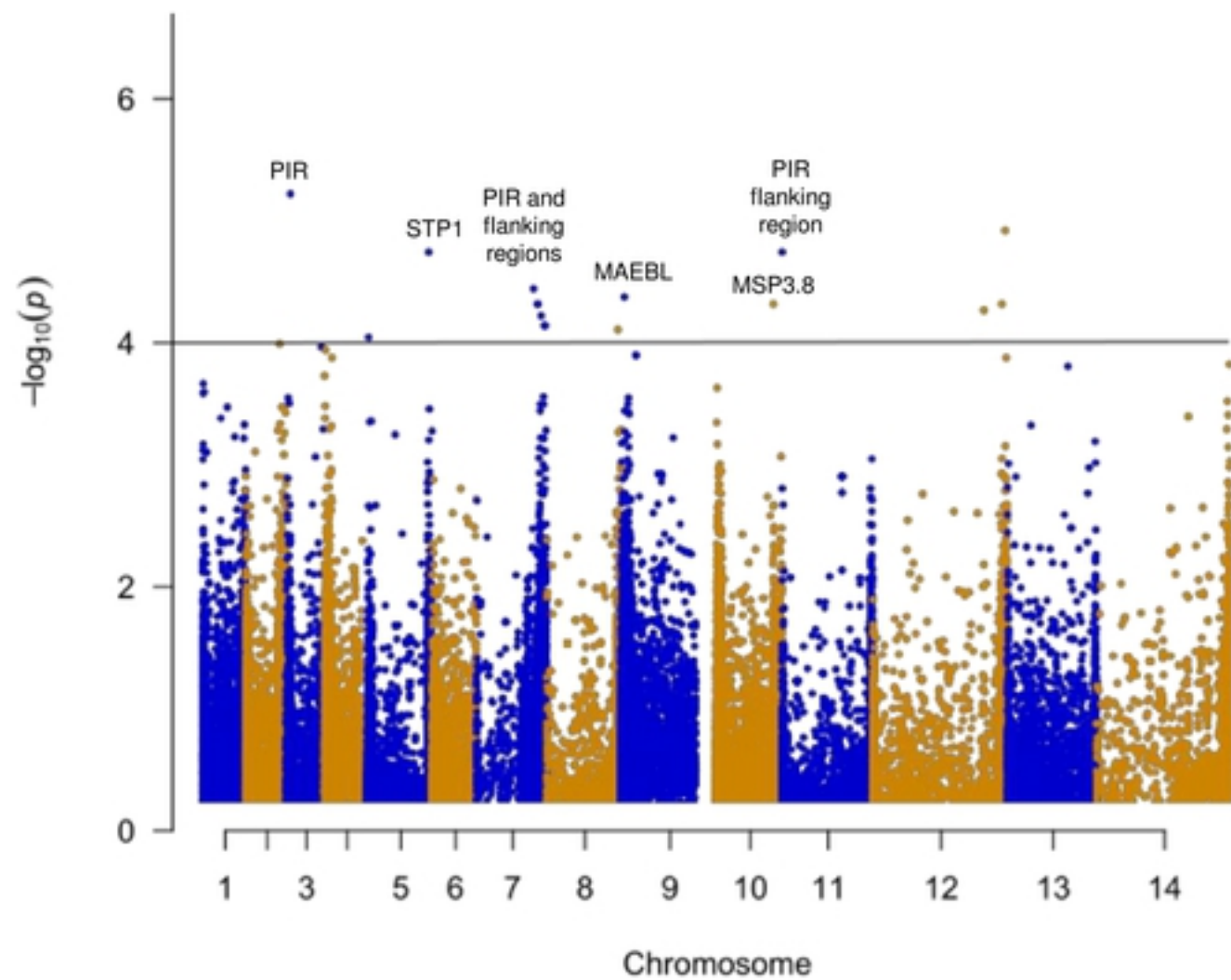
954

955　**Supplementary Table 3.** Likelihood Ratio Test results of the M1 (neutral model) and

956　M2 models (selection model) in PAML of all the 64 erythrocyte binding gene candidates.
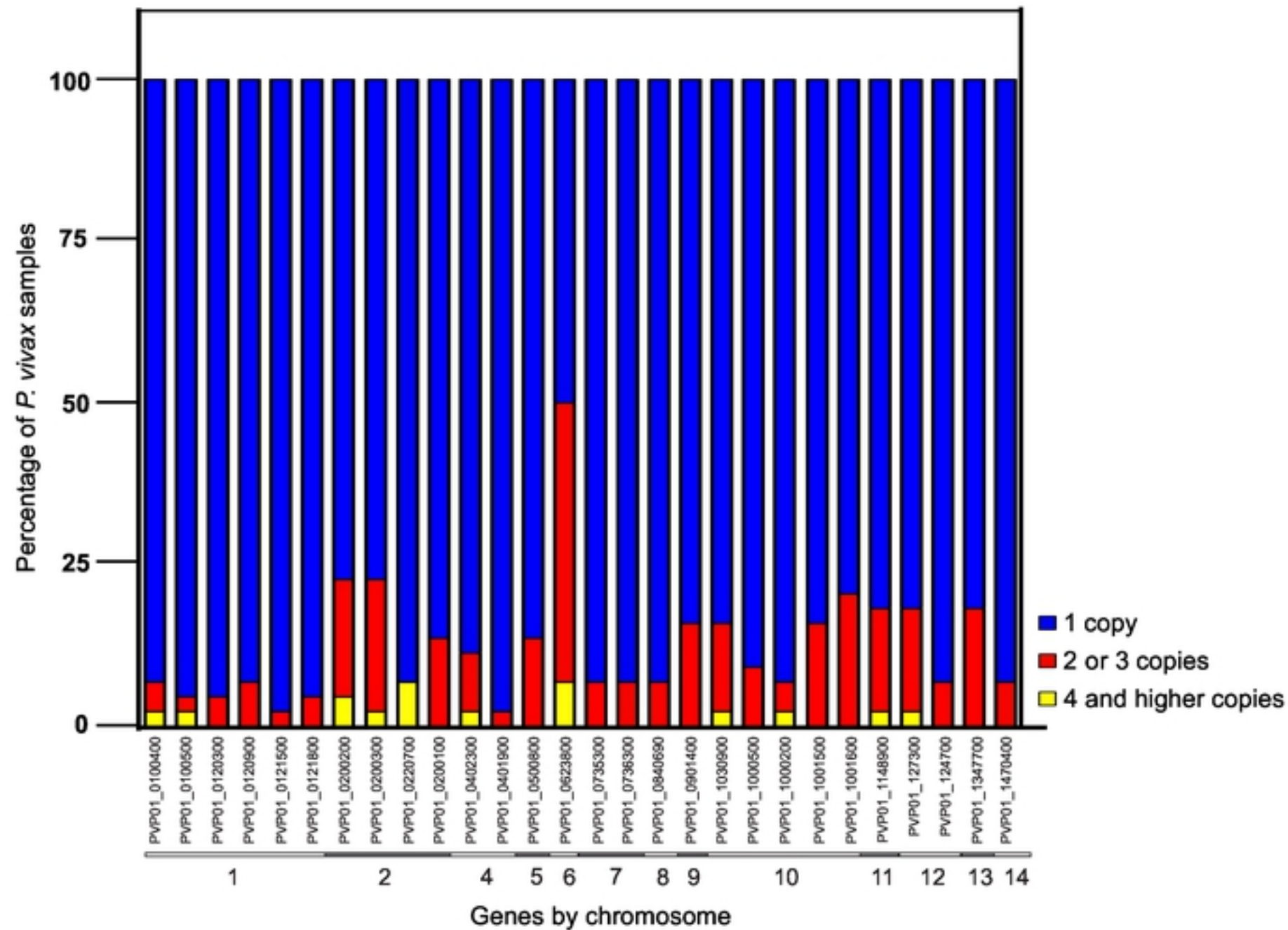
957

958　**Supplementary Table 4.** Gene regions that were detected with copy number variation

959　among the 44 Ethiopian *P. vivax* isolates based on CNVnator. Among them, only two

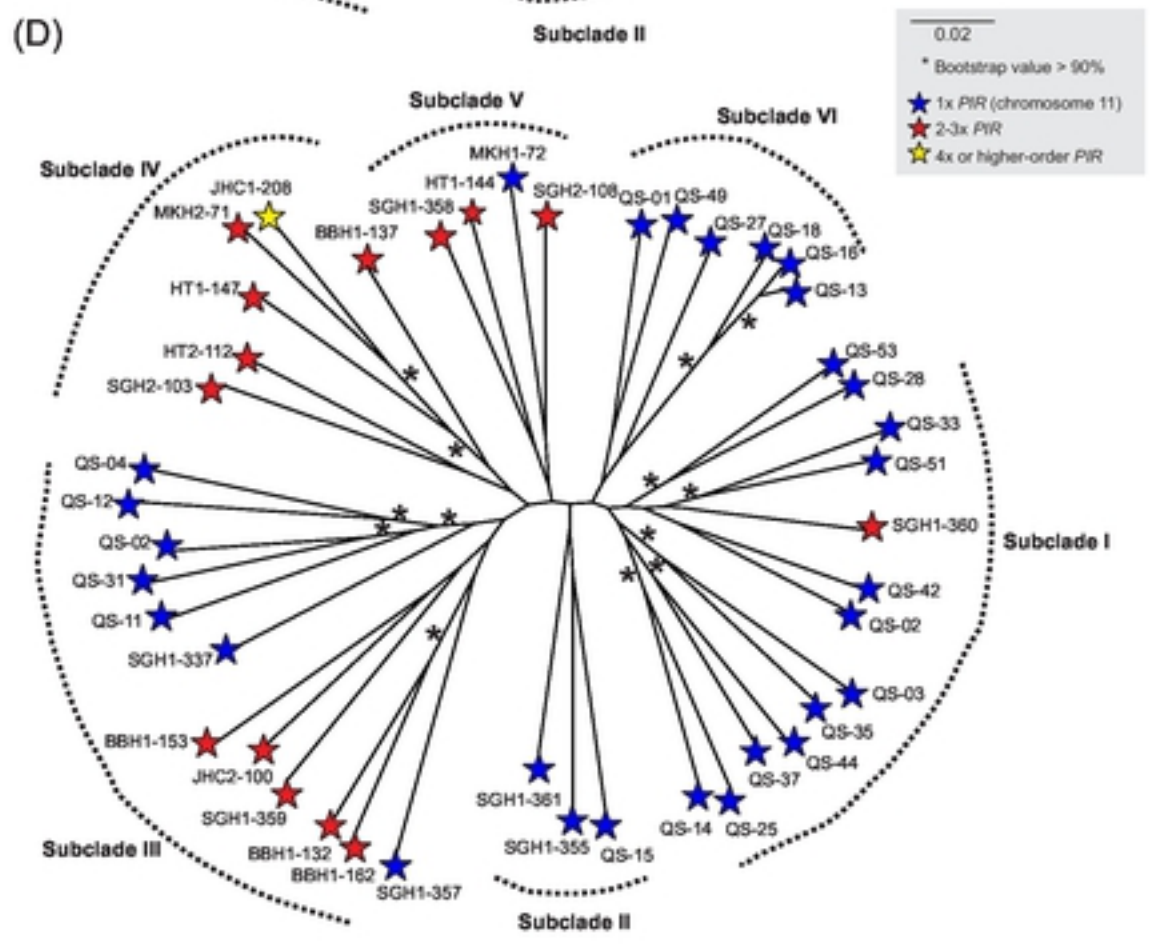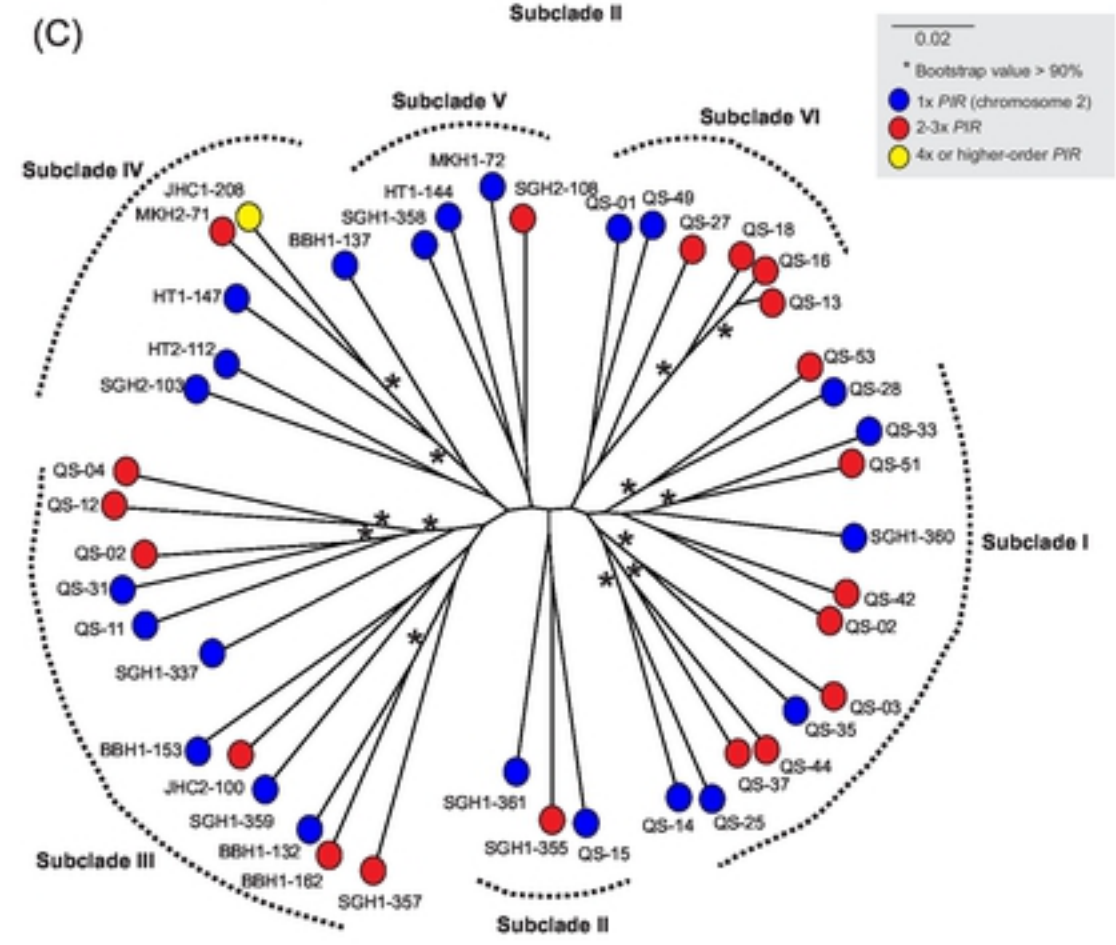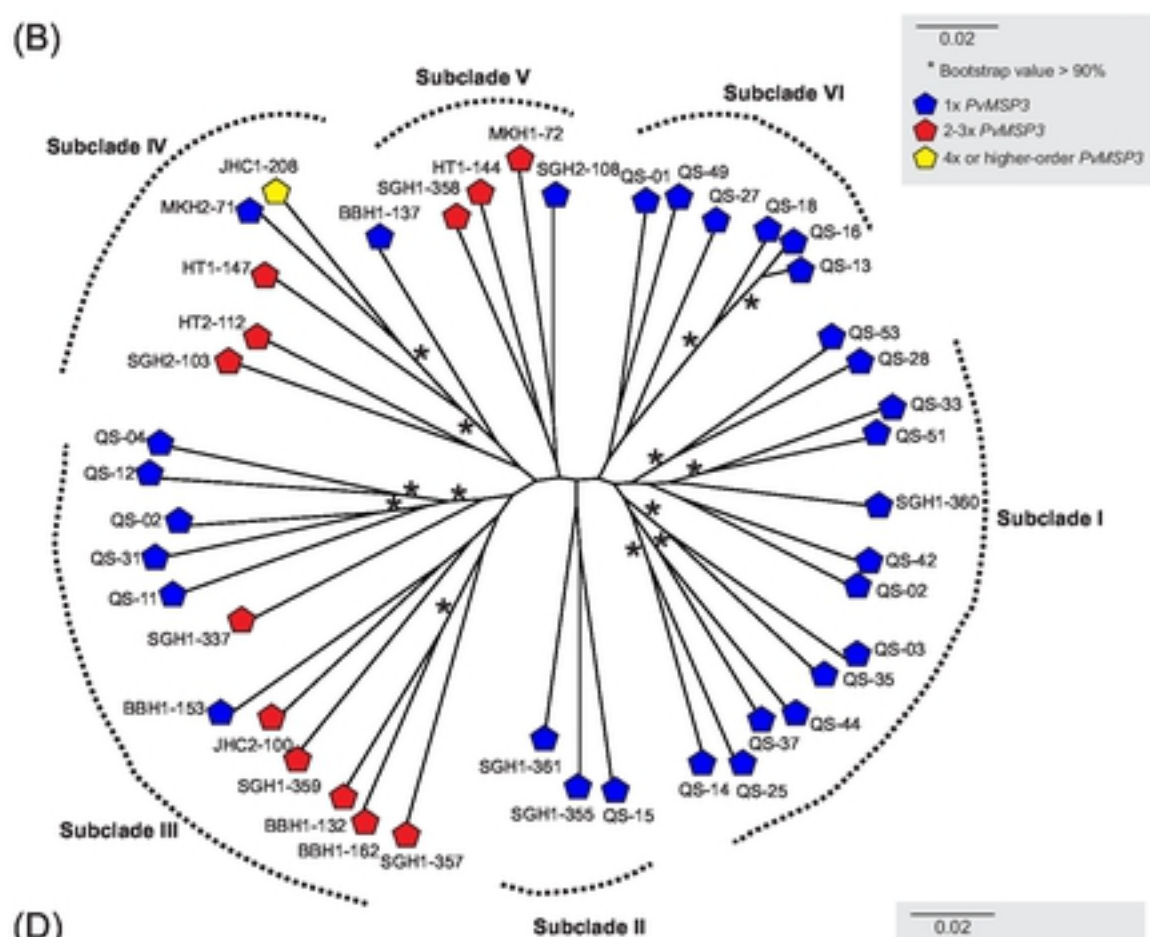960　erythrocyte binding gene candidates *PvDBP*1 and *PvMSP*3 were detected with high-
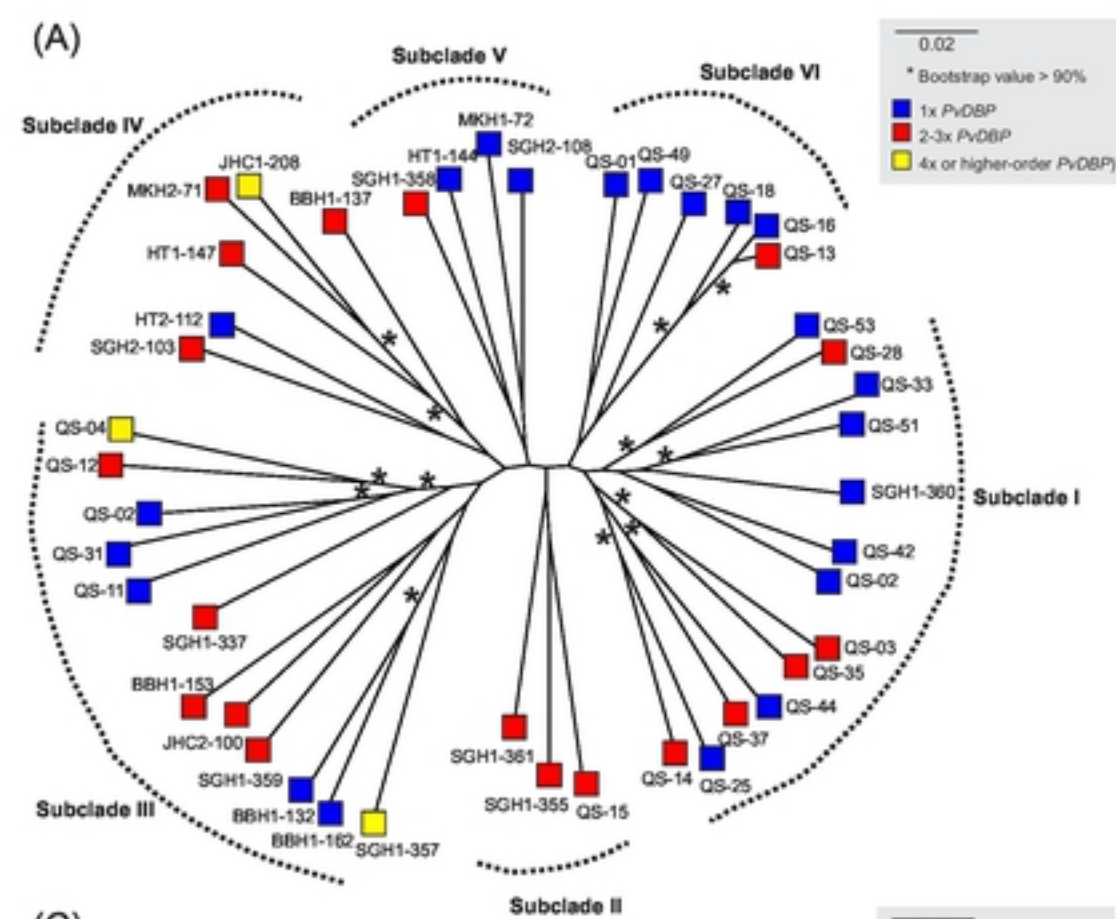
961　order copies.

![Bar chart showing Percentage of *P. vivax* samples (y-axis, 0 to 100) versus Genes by chromosome (x-axis). Legend: blue = 1 copy, red = 2 or 3 copies, yellow = 4 and higher copies.]

Y-axis: Percentage of *P. vivax* samples

X-axis labels: PVP01_0100400, PVP01_0100500, PVP01_0120300, PVP01_0120900, PVP01_0121500, PVP01_0121800, PVP01_0200200, PVP01_0200300, PVP01_0220700, PVP01_0200100, PVP01_0402300, PVP01_0401900, PVP01_0500800, PVP01_0623800, PVP01_0735300, PVP01_0736300, PVP01_0840690, PVP01_0901400, PVP01_1030900, PVP01_1000500, PVP01_1000200, PVP01_1001500, PVP01_1001600, PVP01_1148900, PVP01_1273000, PVP01_1244700, PVP01_1347700, PVP01_1470400

Chromosome groups: 1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14

X-axis title: Genes by chromosome

Legend:
- 1 copy
- 2 or 3 copies
- 4 and higher copies

(A)

Arbaminch
Badowacho
Halaba
Hawassa
Jimma

PC2

PC3

PC1

(B)

Hawassa

Badowacho

Jimma

Arbaminch

Halaba