1 # *Plasmodium vinckei* genomes provide insights into the

2 # pan-genome and evolution of rodent malaria parasites

3

4 Abhinay Ramaprasad [1,2,6], Severina Klaus [2,3], Olga Douvropoulou [1], Richard

5 Culleton [2,4*] and Arnab Pain [1,5*]

6 [1] Pathogen Genomics Group, BESE Division, King Abdullah University of Science

7 and Technology (KAUST), Thuwal 23955-6900, Kingdom of Saudi Arabia

8 [2] Malaria Unit, Department of Pathology, Institute of Tropical Medicine (NEKKEN),

9 Nagasaki University, 1-12-4 Sakamoto, Nagasaki 852-8523, Japan

10 [3] Biomedical Sciences, University of Heidelberg, Heidelberg, Germany

11 [4] Division of Molecular Parasitology, Proteo-Science Center, Ehime University, 454

12 Shitsukawa, Toon, Ehime 791-0295, Japan

13 [5] Center for Zoonosis Control, Global Institution for Collaborative Research and

14 Education (GI-CoRE), Hokkaido University, N20 W10 Kita-ku, Sapporo 001-0020,

15 Japan

16 [6] Present address: Malaria Biochemistry Laboratory, Francis Crick Institute, London

17 NW1 1AT, UK

18 *arnab.pain@kaust.edu.sa, culleton.richard.oe@ehime-u.ac.jp

19

20

21

22

23

24

25

# Abstract

**Background**

Rodent malaria parasites (RMPs) serve as tractable tools to study malaria parasite biology and host-parasite-vector interactions. *Plasmodium vinckei* is the most geographically widespread of the four RMP species collected in sub-Saharan Central Africa. Several *P. vinckei* isolates are available but relatively less characterized than other RMPs, thus hindering their use in experimental studies. We have generated a comprehensive resource for *P. vinckei* comprising of high-quality reference genomes, genotypes, gene expression profiles and growth phenotypes for ten *P. vinckei* isolates.


**Results**

The *P. vinckei* subspecies have diverged widely from their common ancestor and have undergone genomic structural variations. The subspecies from Katanga, *P. v. vinckei*, has a uniquely smaller genome, a reduced multigene family repertoire and is also amenable to genetic manipulation making it an ideal parasite for reverse genetics. Comparing *P. vinckei* genotypes reveals region-specific selection pressures particularly on genes involved in mosquito transmission. The erythrocyte membrane antigen 1 and *fam-c* families have expanded considerably among the lowland forest-dwelling *P. vinckei* parasites. Genetic crosses can be established in *P. vinckei* but are limited at present by low transmission success under the experimental conditions tested in this study.


**Conclusions**

50  *Plasmodium vinckei* isolates display a large degree of phenotypic and genotypic

51  diversity and could serve as a resource to study parasite virulence and

52  immunogenicity. Inclusion of *P. vinckei* genomes provide new insights into the

53  evolution of RMPs and their multigene families. Amenability to genetic crossing and

54  genetic manipulation make them also suitable for classical and functional genetics to

55  study *Plasmodium* biology.

56

## Keywords

57

58  *Plasmodium vinckei*, Malaria, Rodent malaria parasites, Genomics, Transcriptomics,

59  Genetics, Parasite evolution, Multigene families

60

## Background

61

62  Rodent malaria parasites (RMPs) serve as tractable models for experimental

63  genetics and as valuable tools to study malaria parasite biology and host-parasite-

64  vector interactions [1-4]. Between 1948 and 1974, several rodent malaria parasites

65  were isolated from wild thicket rats (shining thicket rat - *Grammomys poensis* or

66  previously known as *Thamnomys rutilans*, and woodland thicket rat - *Grammomys*

67  *surdaster*) and infected mosquitoes in sub-Saharan Africa and were adapted to

68  laboratory-bred mice and mosquitoes. The isolates were classified into four species,

69  namely, *Plasmodium berghei, Plasmodium yoelii, Plasmodium chabaudi and*

70  *Plasmodium vinckei. Plasmodium berghei* and *P. yoelii* are sister species forming the

71  classical *berghei* group, whereas *P. chabaudi* and *P. vinckei* form the classical

72  *vinckei* group of RMPs [5-7]. *Plasmodium chabaudi* has been used for studying drug

73  resistance, host immunity and immunopathology in malaria [8-12]. *Plasmodium yoelii*

74  and *P. berghei* are extensively used as tractable models to study liver and mosquito

75    stages of the parasite [13, 14]. Efficient transfection techniques [15-18] have been

76    established in all three RMPs and they are widely used as *in vivo* model systems for

77    large scale functional studies [19-22]. Reference genomes for these three RMP

78    species are available [23, 24]. Recently, the quality of these genomes has been

79    significantly improved using next-generation sequencing [11, 25, 26].

80

81    *P. vinckei* is the most geographically widespread RMP species, with isolates

82    collected from many locations in sub-Saharan Africa (Figure 1A). Subspecies

83    classifications were made for 18 *P. vinckei* isolates in total based on parasite

84    characteristics and geographical origin, giving rise to five subspecies; *P. v. vinckei*

85    (Democratic Republic of Congo), *P. v. petteri* (Central African Republic), *P. v. lentum*

86    (Congo Brazzaville), *P. v. brucechwatti* (Nigeria) and *P. v. subsp.* (Cameroon) [27-

87    32]. Blood, exo-erythrocytic and sporogonic stages of a limited number of isolates of

88    the five subspecies have been characterized; *P. v. vinckei* line 67 or CY [33], *P. v.*

89    *petteri* line CE [28], *P. v. lentum* line ZZ [29, 31], *P. v. brucechwatti* line 1/69 or DA

90    [30, 34] and several parasite lines of *P. v. subsp.* [35]. Enzyme variation studies [5,

91    36] and multi-locus sequencing data [6, 7] have indicated that there is significant

92    phenotypic and genotypic variation among *P. vinckei* isolates.

93    The rodent malaria parasites isolated from Cameroon in 1974 by J. M. Bafort are

94    currently without subspecies names, being designated as *P. yoelii subsp.*, *P. vinckei*

95    *subsp.* and *P. chabaudi subsp.*. We now present the full genome sequence data of

96    isolates from these subspecies and show they form distinct clades within their parent

97    species. Therefore, we propose the following subspecies names; *Plasmodium yoelii*

98    *cameronensis*, from the country of origin; *Plasmodium vinckei baforti*, after J. M.

99    Bafort, the original collector of this subspecies; and *Plasmodium chabaudi*

4

100  *esekanensis*, from Eséka, Cameroon, the town from the outskirts of which it was

101  originally collected.

102

103  Very few studies have employed *P. vinckei* compared to the other RMP species

104  despite    the    public    availability    of    several    *P.    vinckei*    isolates

105  (http://www.malariaresearch.eu/content/rodent-malaria-parasites). *Plasmodium    v.*

106  *vinckei* v52 and *P. v. petteri* CR have been used to study parasite recrudescence

107  [37], chronobiology [38] and artemisinin resistance [39]. They are also the only

108  isolates for which draft genome assemblies with annotation are available as part of

109  the    Broad    Institute    *Plasmodium*    100    Genomes    initiative

110  (https://www.ncbi.nlm.nih.gov/bioproject/163123).

111

112  A high-quality reference genome for *P. vinckei* and detailed phenotypic and

113  genotypic data are lacking for the majority of *P. vinckei* isolates hindering wide-scale

114  adoption of this RMP species in experimental malaria studies.

115

116  We now present a comprehensive genome resource for *P. vinckei* comprising of

117  high-quality reference genomes for five *P. vinckei* isolates (one from each

118  subspecies) and describe the genotypic diversity within the *P. vinckei* clade through

119  the sequencing of five additional *P. vinckei* isolates (see Figure 1A inset). With the

120  aid of high-quality annotated genome assemblies and gene expression data, we

121  evaluate the evolutionary patterns of multigene families across all RMPs and within

122  the subspecies of *P. vinckei*.

123

124    We also describe the growth and virulence phenotypes of these isolates and show

125    that *P. vinckei* is amenable to genetic manipulation and can be used to generate

126    experimental genetic crosses.

127

128    Furthermore, we sequenced the whole genomes of seven isolates of the subspecies

129    of *P. chabaudi* (*P. c. esekanensis*) and *P. yoelii* (*P. y. yoelii*, *P. y. nigeriensis*, *P. y.*

130    *killicki* and *P. y. cameronensis*) in order to resolve evolutionary relationships among

131    RMP isolates.

132

133    The data presented here enable the use of the *P. vinckei* clade of parasites for

134    laboratory-based experiments driven by high-throughput genomics technologies and

135    will significantly expand the number of RMPs available as experimental models to

136    understand the biology of malaria parasites.

137

138    **Results**

139    ***Plasmodium vinckei* isolates display extensive diversity in virulence**

140    We followed the infection profiles of ten *P. vinckei* isolates in CBA/J mice (five

141    biological replicates per group) to study their virulence traits. Some of these isolates

142    were available as uncloned lines and so were first cloned by limiting dilution

143    (Additional File 1). As reported previously [40], *P. vinckei* parasites are

144    morphologically indistinguishable from each other, prefer to invade mature

145    erythrocytes, are largely synchronous during blood stage growth and display a

146    characteristically rich abundance of haemozoin crystals in their trophozoites and

147    gametocytes (Figure 1B).

148

149    Parasitaemia was determined daily to measure the growth rate of each isolate and

150    host RBC density and weight were measured as indications of "virulence" (harm

151    to the host) (Figure 1C, Additional file 2 and 3).

152

153    The *P. v. vinckei* isolate *Pvv*CY, was highly virulent and reached a parasitaemia of

154    89.4% $\pm$ 1.4 (standard error of mean; SEM) on day 6 post inoculation of 1 x $10^6$

155    blood stage parasites intravenously, causing host mortality on that day. Both strains

156    of *P. v. brucechwatti*, *Pvb*DA and *Pvb*DB, were virulent and killed the host on day 7

157    or 8 post infection (peak parasitaemia of around 70%). The *P. v. lentum* parasites

158    *Pvl*DS and *Pvl*DE, were not lethal and were eventually cleared by the host immune

159    system, with *Pvl*DS's clearance more prolonged than that of *Pvl*DE (parasitaemia

160    clearance rates; *Pvl*DS = 10.35 %day$^{-1}$; SE = 1.105; p-value of linear fit =0.0025;

161    *Pvl*DE = 16.46 %day$^{-1}$; SE = 3.873; p-value =0.023). The *P. v. petteri* isolates

162    *Pvp*CR and *Pvp*BS reached peak parasitaemia along similar timelines (6-7 dpi), but

163    *Pvp*CR was virulent (peak parasitaemia = 60.35 % $\pm$ 2.38 on day 6) and could

164    sometimes kill the host while *Pvp*BS maintained a mild infection.

165

166    Of the three isolates of *P. vinckei baforti*, *Pvs*EL and *Pvs*EE were similar in their

167    growth profiles and their perceived effect on the host, while in contrast, *Pvs*EH was

168    highly virulent, causing host mortality at day 5, the earliest among all *P. vinckei*

169    parasites.

170

171    RBC densities reduced during the course of infection proportionally to the rise in

172    parasitaemia in all the *P. vinckei* infection profiles studied. There were differences,

173    however, in the patterns of host weight loss. Mild infections by *P. v. lentum* isolates

174    (maximum weight loss in *Pvl*DE = 0.43 mg ± 0.41 and *Pvl*DS = 1.77 mg ± 0.38), *P. v.*

175    *petteri* BS (0.58 mg ± 0.22) and *P. v. baforti* EE (1.66 mg ± 0.31, 0.09 mg ± 0.56) did

176    not cause any significant weight loss in mice, whereas the virulent strains, *P. v.*

177    *petteri* CR (4.04 mg ± 0.18), *P. v. brucechwatti* isolates (*pvb*DA = 3.5 mg ± 0.39 and

178    *pvb*DB = 2.05 mg ± 1.68) caused around a 20% decrease in weight. Virulent strains

179    *Pvv*CY (1.74 mg ± 0.15) and *Pvs*EH (0.52 mg ± 0.13) did not cause any significant

180    weight loss during their infection before host death occurred.

181

182    **_Plasmodium vinckei_ reference genome assembly and annotation**

183    High-quality reference genomes for five *P. vinckei* isolates, one from each

184    subspecies; *P. v. vinckei* CY (*Pvv*CY), *P. v. brucechwatti* DA (*Pvb*DA), *P. v. lentum*

185    DE (*Pvl*DE), *P. v. petteri* CR (*Pvp*CR) and *P. v. baforti* EL (*Pvs*EL) were assembled

186    from single-molecule real-time (SMRT) sequencing. PacBio long reads of 10-20

187    kilobases (kb) and with a high median coverage of >155X across the genome

188    (Additional File 1) enabled *de novo* assembly of each of the 14 chromosomes as

189    single unitigs (high confidence contig) (see Table 1). PacBio assembly base call

190    errors were corrected using high-quality 350bp and 550bp insert PCRfree Illumina

191    reads. A small number of gaps remain in the assemblies, but these are mainly

192    confined to the apicoplast genomes and to the *Pvs*EL and *Pvl*DE genomes that were

193    assembled from 10kB-long PacBio reads instead of 20kB. The *Pvp*CR and *Pvv*CY

194    assemblies, with each chromosome in one piece, are a significant improvement over

195    their existing fragmented genome assemblies (available through PlasmoDB v.30).

196

197    *Plasmodium vinckei* genome sizes range from 19.2 to 19.5 Mb except for *Pvv*CY

198    which has a smaller genome size of 18.3 Mb, similar to that of *P. berghei* (both

8

199   isolates are from the same Katanga region). While we were not able to resolve the

200   telomeric repeats at the ends of some of the chromosomes, all the resolved

201   telomeric repeats had the RMP-specific sub-telomeric repeat sequences

202   CCCTA(G)AA. The mitochondrial and apicoplast genomes were ~6Kb and ~30 kb

203   long respectively, except for the apicoplast genomes of *Pvp*CR and *Pvs*EL for which

204   we were able to resolve only partial assemblies due to low read coverage (see

205   Additional File 4).

206

207   Gene models were predicted by combining multiple lines of evidence to improve the

208   quality of those predictions. These include publicly available *P. chabaudi* gene

209   models, *de novo* predicted gene models and transcript models from strand-specific

210   RNA-seq data of different blood life cycle stages. Consensus gene models were then

211   manually corrected through comparative genomics and visualization of mapped

212   RNAseq reads. As a result, we annotated 5,073 to 5,319 protein-coding genes, 57-

213   67 tRNA genes and 40-48 rRNA genes in each *P. vinckei* genome. Functional and

214   orthology analyses with the predicted *P. vinckei* proteins showed that the core

215   genome content in *P. vinckei* parasites is highly conserved among the species and

216   are comparable to other rodent and primate malaria species.

217

218   ***Plasmodium vinckei* genome assemblies reveal novel structural variations**

219   Comparative analysis of *P. vinckei* and other RMP genomes shows that *P. vinckei*

220   genomes exhibit the same high level of synteny seen within RMP genomes, but with

221   a number of chromosomal rearrangements. These events can be identified by

222   breaks in synteny (synteny breakpoints- SBPs) observed upon aligning and

223   comparing genome sequences.

9

224

225    We aligned *P. vinckei* and other RMP genomes to identify synteny blocks between

226    their chromosomes. Similar to previous findings in RMP genomes [26, 41] (Additional

227    file 12A), we observed large scale exchange of material between non-homologous

228    chromosomes, namely three reciprocal translocation events and one inversion

229    (Figure 2A, Additional File 5 and 12). A pan-*vinckei* reciprocal translocation of

230    ~0.6Mb (with 134 genes) and ~0.4 Mb (with 99 genes) long regions between

231    chromosomes VIII and X was observed between *P. vinckei* and *P. berghei* (whose

232    genome closely resembles that of the putative RMP ancestor [41]). Within the *P.*

233    *vinckei* subspecies, two reciprocal translocations separate *P. v. petteri* and *P. v.*

234    *baforti* from the other three subspecies. One pair of exchanges (~1 Mb and ~0.55

235    Mb) was observed between chromosomes V and XIII, and another smaller pair

236    (~150Kb and ~70Kb) between chromosomes V and VI. These events have left the

237    Chromosomes V of *Pvv*CY-*Pvb*DA-*Pvl*DE and *Pvp*CR-*Pvs*EL groups with only a

238    ~0.15 Mb region of synteny between them, consisting of 48 genes while the

239    remaining 304 genes have been rearranged with chromosome VI and XII.

240

241    There also exists a small, *Pvv*CY-specific inversion of a ~100 kb region in

242    chromosome XIV. All the synteny breakage points (SBPs) were verified manually

243    and were supported by PacBio read coverage ruling out the possibility of a

244    misassembly at the breakpoint junctions. The SBPs in chromosomes V and VI were

245    near rRNA units, loci previously described as hotspots for such rearrangement

246    events [42, 43].

247

248 **A pan-RMP phylogeny reveals high genotypic diversity within the *P. vinckei***

249 **clade**

250 In order to re-evaluate the evolutionary relationships among RMPs, we first inferred

251 a well-resolved species-level phylogeny that takes advantage of the manually

252 curated gene models in eight available high-quality RMP genomes representing all

253 RMPs. A maximum-likelihood phylogeny tree was inferred through partitioned

254 analysis using RAxML, of a concatenated protein alignment (2,281,420 amino acids

255 long) from 3,920 single-copy, conserved core genes in eleven taxa (eight RMPs, *P.*

256 *falciparum*, *P. knowlesi* and *P. vivax*; see Figure 2B and Additional File 6).

257

258 In order to assess the genetic diversity within RMP isolates, we sequenced

259 additional isolates for four *P. vinckei* subspecies (*Pvb*DB, *Pvl*DS, *Pvp*BS, *Pvs*EH and

260 *Pvs*EE), *P. yoelii yoelii* (*Pyy*33X, *Pyy*CN and *Pyy*AR), *P. yoelii nigeriensis* (*Pyn*D), *P.*

261 *yoelii killicki* (*Pyk*DG), *P. yoelii subsp.* (*Pys*EL) and *P. chabaudi subsp.* (*Pcs*EF)

262 (Additional File 1). This, along with existing sequencing data for 13 RMP isolates

263 (from [26, 44]), were used to infer an isolate-level, pan-RMP maximum likelihood

264 phylogeny based on 1,010,956 high-quality SNPs in non-subtelomeric genes that

265 were called by mapping all reads onto the *Pvv*CY reference genome (Figure 2C and

266 Additional File 7). Both phylogenies were well-resolved with robust 100% bootstrap

267 support obtained for the amino-acid based phylogeny and 78% or higher bootstrap

268 support for the SNPs-based phylogeny (majority-rule consensus tree criterion was

269 satisfied at 50 bootstraps for both the phylogenies).

270

271 Both protein alignment-based and SNP-based phylogenies show significant

272 divergence among the *P. vinckei* subspecies compared to the other RMPs. All *P.*

11

273     *vinckei* subspecies have begun to diverge from their common ancestor well before

274     sub-speciation events within *P. yoelii* and *P. chabaudi*.

275

276     A total of 521,934 polymorphic positions were found within the *P. vinckei* core coding

277     regions consisting of 4,644 non-subtelomeric genes across *P. vinckei* isolates. The

278     Katangan isolate, *P. v. vinckei*, has undergone significant divergence from the

279     common *vinckei* ancestor and is the most diverged of any RMP subspecies

280     sequenced to date. Number of SNPs ranged from 292,240 to 318,344 in pair-wise

281     comparisons of isolates with *Pvv*CY, with 4,237 to 4,263 genes (out of the 4,644

282     core genes) having at least 5 SNPs. *Plasmodium v. brucechwatti* has also diverged

283     significantly, while the divergence of *P. v. lentum* is comparable to that of *P. y.*

284     *nigeriensis*, *P. y. kilicki* and *P. c. subsp.* from their respective putative ancestors.

285     Genetic diversity within *P. v. petteri* and *P. v. baforti* isolates are similar to that

286     observed within *P. yoelii* and *P. chabaudi* isolates while *P. v. lentum* and *P. v.*

287     *brucechwatti* isolates have exceptionally high and low divergences respectively.

288     Number of SNPs in pair-wise comparisons between the closest subspecies, *P. v.*

289     *petteri* and *P. v. baforti*, ranged from 53,554 (*P. v. baforti* EE) to 69,454 (*P. v. baforti*

290     EL) with 2,358 genes having at least 5 SNPs.

291

292     Our robust phylogeny based on a comprehensive set of genome-wide sequence

293     variations confirms previous estimates of RMP evolution based on isoenzyme

294     variation [5] and gene sequences of multiple housekeeping loci [6, 7], except for the

295     placement of *P. y. nigeriensis* D which we show to be diverged earlier than *P. y.*

296     *kilicki* DG (supported by a bootstrap value of 100).

297

298 **Molecular evolution within *P. vinckei* isolates**

299 Using SNP data (Additional file 7), we then assessed the differences in selection

300 pressure on the geographically diverse *P. vinckei* isolates by calculating the gene-

301 wise Ka/Ks ratio as a measure of enrichment of non-synonymous mutations in a

302 gene (signifying positive selection). We first compared the Katangan isolate (*Pvv*CY)

303 from the highland forests in the DRC with the non-Katangan isolates from the

304 lowland forests elsewhere.

305

306 We made pairwise comparisons of the four non-Katangan *P. vinckei* subspecies with

307 *P. v. vinckei* which revealed several genes under significant positive selection

308 (Figure 3C). Notably, we identified three genes involved in mosquito transmission,

309 namely, a gamete-release protein (GAMER), a secreted ookinete protein (PIMMS43,

310 previously known as PSOP25) and a thrombospondin-related anonymous protein

311 (TRAP), featuring in all Katangan/non-Katangan subspecies comparisons.

312

313 GAMER (PVVCY_1202630 being the representative ortholog in *P. v. vinckei*;

314 PBANKA_1225400 in *P. berghei*) had high Ka/Ks values in all comparisons (except

315 for *P. v. vinckei*- *P. v. brucechwatti*) and is essential for gamete egress [45].

316 PIMMS43 (PVVCY_1102000; PBANKA_1119200) and TRAP (PVVCY_1305250;

317 PBANKA_1349800) showed high Ka/Ks values in all comparisons and are essential

318 for ookinete evasion from mosquito immune system [46] and sporozoite infectivity of

319 mosquito salivary glands and host hepatocytes [47] respectively.

320

321 Several exported proteins and surface antigens were also identified to have

322 undergone positive selection. PVVCY_0100120 (PCHAS_0100651 being the gene

323 ortholog in *P. chabaudi*) has a circumsporozoite-related antigen PFAM domain

324 (PF06589) and is a conserved protein found in all RMPs except *P. berghei*.

325 PVVCY_1200100 (PBANKA_1002600) is a merozoite surface antigen, p41 [48] that

326 is secreted following invasion [49].

327

328 To assess presence of geographic location-specific selection pressures among the

329 lowland forest isolates, *P. v. brucechwatti*, *P. v. lentum* and *P. v. baforti* were

330 compared with *P. v. petteri* CR from the CAR. To see if similar selection pressures

331 have acted on other RMP species, we also analysed the *P. yoelii* and *P. chabaudi*

332 isolates from these regions that we had sequenced in this study.

333

334 Several exported and rhoptry-associated proteins were identified as been under

335 positive selection in each comparison but in contrast to comparisons with *P. v.*

336 *vinckei*, there was no overlap of positively selected genes among the non-Katangan

337 isolates. However, we identified a conserved rodent malaria protein of unknown

338 function (PVVCY_0501990; PBANKA_051950) that seems to be under significant

339 positive selection with high Ka/Ks values (ranging from 2.14 to 4.39) in all *P. vinckei*

340 comparisons except *P. v. petteri - P. v. baforti*. The *P. yoelii* ortholog of this protein

341 was also positively selected among *P. y. yoelii, P. y. nigeriensis* and *P. y. killicki* but

342 was not under selection within the *P. y. yoelii* isolates, signifying region-specific

343 selection pressures.

344

345 A 28kDa ookinete surface protein (P28; PVVCY_0501540; PBANKA_0514900)

346 seem to be under positive selection in the Nigerian *P. v. brucechwatti* as it features

347 in both *P. v. vinckei - P. v. brucechwatti* and *P. v. brucechwatti - P. v. petteri*

348   comparisons. The protein is also seen positively selected among corresponding

349   Nigerian and Central African Republic *P. yoelii* isolates (*P. y. nigeriensis – P. y.*

350   *yoelii*). A protein phosphatase (PPM8; PVVCY_0903370; PBANKA_0913400) has

351   also undergone positive selection in all three RMP species between CAR and Congo

352   isolates (*P. c. chabaudi- P. c. adami* comparison from [26]).

353

354   **Evolutionary patterns within the RMP multigene families**

355   We were able to accurately annotate members of the ten RMP multigene families in

356   the *P. vinckei* genomes owing to the well-resolved sub-telomeric regions in the

357   Pacbio assemblies and manually curated gene models (see Table 1 and Additional

358   file 8). *P. v. vinckei*, similar to its sympatric species, *P. berghei*, had a lower

359   multigene family repertoire with copy numbers strikingly less than other *P. vinckei*

360   subspecies (exceptions were the *pir*, *etramp* and lysophospholipase families).

361   Multigene family sizes in the four non-Katangan *P. vinckei* subspecies were similar

362   to *P. chabaudi* except for expansion in the *ema1* and *fam-c* multigene families

363   (Figure 3A).

364

365   Next, we inferred maximum likelihood-based phylogenies for the ten multigene

366   families, in order to identify structural differences amongst their members and to

367   determine family evolutionary patterns across RMP species and *P. vinckei*

368   subspecies (Figure 3B, Additional file 3 and 4). Overall, we identified robust clades

369   (with bootstrap value >70) that fell into the following categories, i) pan-RMP, with

370   orthologous genes from the four RMP species (dark grey), ii) *berghei* group, with

371   genes from *P. berghei* and *P. yoelii* alone, iii) *vinckei* group, with genes from *P.*

372   *chabaudi* and any or all *P. vinckei* subspecies, iv) *P. vinckei*, with genes only from *P.*

15

373  *vinckei* subspecies  and v) non-Katangan, with genes from all *P. vinckei* subspecies

374  except *P. v. vinckei*.

375

376  In general, a high level of orthology was observed between *P. chabaudi* and *P.*

377  *vinckei* genes forming several *vinckei* group clades (marked in orange in Figure 3) in

378  contrast to more species-specific clades of paralogous genes being formed in *P.*

379  *berghei* and *P. yoelii*. Thus, family expansions in *P. chabaudi* and *P. vinckei* seem to

380  have occurred in the common *vinckei* group ancestor prior to speciation.

381  We rebuilt the phylogenetic trees for *pir*, *fam-a* and *fam-b* families in order to see if

382  the previously defined clades [25, 26] were also maintained in *P. vinckei*. Overall, we

383  were able to reproduce the tree structures for the three families using the ML method

384  with *P. vinckei* gene family members now added to them.

385

386  For *pirs*, we obtained four long-form and eight short-form clades as in [26] (Additional

387  file 9, tree 10) albeit with lower bootstrap support, possibly due to our overly stringent

388  automated trimming of the sequence alignment (see Methods). With a few

389  exceptions, *P. vinckei pir* genes majorly populated two clades - L1 and S7 and a

390  subclade S1g. These clades, previously shown to be *P. chabaudi*-dominant, hold

391  equal or near equal proportions of *P. vinckei* pirs too. The only other *P. chabaudi*-

392  dominant clade, L4, remains as a completely *P. chabaudi*-specific gene expansion.

393  No *P. vinckei* species- or subspecies-specific clades are evident except for two

394  subclades that could be inferred as *Pvv*CY-specific expansions within L1 and S7

395  (marked i and ii in Additional file 9, tree 10). Speciation of *P. vinckei* subspecies from

396  their common ancestor seems to have been accompanied by gene gain in L1 and S7

397  clades and gene loss in S1g subclade. There is an almost linear increase of around

16

398    20 genes in *Pvv*CY, *Pvb*DA and *Pvl*DE in clade L1 pirs and a near doubling of clade

399    S7 pirs in *Pvp*CR.

400

401    The *fam-a* and *fam-b* phylogenies (Additional file 9, tree 3 and 4 respectively) show

402    that previously identified ancestral lineages [25] are maintained in *P. vinckei* too. The

403    addition of *P. vinckei* genes resolved the ancestral clade of internal *fam-a* genes in

404    chromosome 13 further into several well-supported *vinckei* group clades and a

405    *berghei* group clade (marked as A in Additional file 9, tree 3). The 19 other *fam-a*

406    clades and five *fam-b* clades consisting of positionally conserved orthologous genes

407    are also conserved in *P. vinckei* (pan-RMP clades marked with * in Additional file 9,

408    tree 3 and 4). The *fam-a* family has expanded in the non-Katangan *P. vinckei*

409    subspecies through independent events of gene duplication in their common

410    ancestor giving rise to several non-Katangan clades (marked as B in Additional file 9,

411    tree 3). There is only a moderate *P. vinckei*-specific expansion in *fam-b* giving rise to

412    three clades (marked as A in Additional file 9, tree 4) that includes *Pvv*CY genes too,

413    pointing to gene duplications in the *P. vinckei* common ancestor. In both the

414    phylogenies, species and subspecies-specific gene duplication events within the

415    *vinckei* group are rare but do occur (marked as i-iv in Additional file 9, tree 3 and 4).

416

417    The *fam-d* multigene family is present as a single ancestral copy in *P. berghei*

418    internally on chromosome IX but is expanded into a gene cluster in the same loci in

419    *P. yoelii* (5 genes) and *P. chabaudi* (21 genes). Similar expansions have occurred in

420    *P. vinckei* subspecies and phylogenetic analysis shows the presence of six robust

421    clades within this family (Additional file 9, tree 6). Clade I is clearly the ancestral

422    clade from which all other *fam-d* genes have been derived as it consists of the single

17

423    *P. berghei* gene and its orthologs in other RMPs, positionally conserved to be the

424    outermost gene of the *fam-d* cluster in each RMP.

425

426    While the *fam-d* family in *P. yoelii* is completely a product of paralogous expansion

427    within Clade I, the *fam-d* families in the *vinckei* group seem to have expanded *via*

428    five ancestral lineages forming clades II-VI. A subset of orthologs in Clade II (marked

429    with * in Additional file 9, tree 6) are positionally conserved among *vinckei* group

430    parasites, located immediately after the *fam-d* ancestral copy and could therefore

431    represent the Clade II ancestral gene in the *vinckei* group common ancestor. *Pvv*CY

432    has a smaller *fam-d* repertoire of 6 genes derived from only three of the five *vinckei*

433    group lineages (Clade II, IV and VI), apart from the conserved ancestral copy.

434

435    ML-based trees for haloacid dehalogenase-like hydrolase (*hdh*), putative reticulocyte

436    binding proteins (*p235*) and lysophospholipases (*lpl*) have generally well-resolved

437    topologies with robust bootstrap support for their nodes and some clades contain

438    syntenic orthologous genes (clades marked with * in Additional file 9, trees 7, 8 and

439    9) to member genes in *P. falciparum,* for example, *Pf*HAD2, *Pf*HAD3, *Pf*HAD4 and

440    *Pf*RH6. Poor bootstrap support was obtained for the *etramp* tree (Additional file 9,

441    tree 2), however clades were identified for some members including *uis3*, *uis4* and

442    *etramp10.2*.

443

444    In order to assess the general level of expression of multi-gene family members in

445    blood stage parasites, we superimposed blood-stage RNAseq data onto the

446    phylogenetic trees. Life stage specific expression of multigene family members in

447    five RMPs – *P. berghei* (rings, trophozoites, schizonts and gametocytes) [26], *P.*

448  *chabaudi* and *P. v. vinckei* (rings, trophozoites and schizonts) [50], *P. v. petteri* and

449  *P. v. lentum* (rings, trophozoites and gametocytes) and mixed blood stage

450  expression levels in *P. yoelii* [44], *P. v. brucechwatti* and *P. v. baforti* were assessed

451  (Additional File 10).

452

453  The level of gene transcription was designated low for genes with normalized FPKM

454  (Fragments Per Kilobase of transcript per Million mapped reads) less than 36,

455  medium if between 36 and 256 and high for genes with FPKM above 256. Both the

456  levels and life stage specificity of gene expression within the various clades were

457  generally conserved across the RMPs signifying that orthologs in structurally distinct

458  clades might have conserved functions across the different RMPs. In general, the

459  proportion of transcribed genes in all multigene families in *P. vinckei* was similar to

460  that observed in *P. chabaudi* and slightly higher than *P. berghei* and *P. yoelii*

461  (excluding the families with only one or two *P. berghei* or *P. yoelii* members).

462

463  **The erythrocyte membrane antigen 1 and *fam-c* sub-telomeric multigene**

464  **families are expanded in non-Katangan *P. vinckei* parasites.**

465  The erythrocyte membrane antigen 1 (EMA1) was first identified and described in *P.*

466  *chabaudi* and is associated with the host RBC membrane [51]. These genes encode

467  for a ~800 aa long protein and consist of two exons; a first short exon carrying a

468  signal peptide followed by a longer exon carrying a PcEMA1 protein family domain

469  (Pfam ID- PF07418). The gene encoding EMA1 is present only as a single copy in *P.*

470  *yoelii* or as two copies in *P. berghei* but has expanded to 14 genes in *P. chabaudi*.

471  We see similar gene expansions of 15 to 21 members in the four non-Katangan *P.*

472  *vinckei* parasites (*Pvb*DA, *Pvl*DE, *Pvp*CR and *Pvs*EL; Figure 3A and Additional file

473     3). However, almost half of these genes are pseudogenes with a conserved SNP

474     (C>A) at base position 14 that introduces a TAA stop codon (S5X) within the signal

475     peptide region, followed by a few more stop codons in the rest of the gene (see

476     Additional File 12B). Apart from one or two cases, the S5X mutation is found in all

477     pseudogenes belonging to the *ema1* family and is *vinckei*-specific (it is not present in

478     the single *P. chabaudi* pseudogene).

479

480     A ML-based phylogeny inferred for the 99 *ema1* genes was in general well-resolved

481     with robust branch support for most nodes (see Figure 3B). Four distinct *vinckei*

482     group-specific clades (Clade I to IV), two *vinckei*-specific clades (Clade IV and V)

483     and a non-Katangan *P. vinckei* - specific clade (VI) with good basal support

484     (bootstrap value of 75-100) were identified. Clade I-IV each consist of *ema1* genes

485     positionally conserved across *P. chabaudi* and all five *P. vinckei* subspecies (in

486     chromosomes I, VII, IX and X respectively) and are actively transcribed during blood

487     stages.

488

489     Of the two *P. berghei ema1* genes, one forms a distal clade with the single *P. yoelii*

490     gene while the other is paraphyletic within Clade IV, pointing to presence of two

491     *ema1* loci in the common RMP ancestor, one of which was possibly lost during

492     speciation of *P. yoelii*. All seven *Pvv*CY *ema1* genes are found within clades I to V

493     with gene duplication events in Clade III and V.

494

495     Family expansion in *P. chabaudi* is mainly driven by gene duplication giving rise to *P.*

496     *chabaudi*-specific clades. In contrast, family expansion within non-Katangan *P.*

497     *vinckei* parasites is mainly driven by expansion of pseudogenized *ema1* genes (41

498  genes). Except for some *P. v. brucechwatti*-specific gene expansions, the

499  pseudogenes do not form subspecies-specific clades suggesting that the expansion

500  must have occurred in their non-Katangan *P. vinckei* common ancestor.

501

502  Gene expression data shows that the members of the *P. vinckei*-specific Clade IV

503  are heavily transcribed during blood stages but most *ema1* pseudogenes that share

504  ancestral lineage with Clade IV are very weakly transcribed. Taken together, a core

505  repertoire of conserved *ema1* genes arising from 4-5 independent ancestral lineages

506  are actively transcribed during blood stages of *P. vinckei* and *P. chabaudi*. The *ema1*

507  multigene family expansion in *P. vinckei* is largely due to duplications of *ema1*

508  pseudogenes, all carrying a S5X mutation and lacking transcription.

509

510  The *fam-c* proteins are exported proteins characterized by *pyst-c1* and *pyst-c2*

511  domains, first identified in *P. yoelii* [42]. There is a considerable expansion of this

512  family in the non-Katangan *P. vinckei* strains resulting in 59 to 65 members, twice

513  that of *Pvv*CY and other RMPs. The *fam-c* genes are exclusively found in the sub

514  telomeric regions and are composed of two exons and an intron, of which the first

515  exon is uniformly 80 bps long (with a few exceptions). *fam-c* proteins are

516  approximately 100-200 amino acids long, and more than one third of the proteins in

517  *P. vinckei* contain a transmembrane domain (75.5%) and a signal peptide (88.9%)

518  but most of them lack a PEXEL-motif (motif was detected in only 4% of the genes

519  compared to 24% in other RMPs).

520

521  An ML-based tree of all *fam-c* genes in the eight RMP species shows the presence

522  of four distinctly distal clades (marked as A in Figure 3B) with robust basal support

21

523     (96-100). Two of them are pan-RMP and two are *vinckei* group-specific, each

524     consisting of *fam-c* genes positionally conserved across the member subspecies

525     (taking into account the genome rearrangement between chromosome V and VI

526     within the *vinckei* clade). Most members of these clades show medium to high gene

527     expression during asexual blood stages.

528

529     The remainder of the tree's topology does not have good branch support (<70) with

530     the exception of some terminal nodes, but it does demonstrate the significant

531     expansion of this gene family within non-Katangan *P. vinckei* parasites (clades

532     shaded in blue).

533

534     There is evidence of significant species- and subspecies-specific expansions with

535     striking examples in *P. yoelii*, *P. chabaudi* and in *P. v. brucechwatti* (marked i, ii and

536     iii in Figure 3B respectively), though they do not form well-supported clades. Most

537     *fam-c* genes in *P. yoelii* seem to have originated from such independent *P. yoelii* -

538     specific expansion events. *P. chabaudi* and *P. v. vinckei fam-c* genes are found

539     more widely dispersed throughout the tree suggesting divergence of this family in the

540     *vinckei* group common ancestor. On the other hand, subspecies-wise distinctions

541     among the non-Katangan *P. vinckei fam-c* genes are less resolved as they form both

542     paralogous and orthologous groups between the four subspecies with several

543     ortholog pairs strongly supported by bootstrap values.

544

545     Thus, *fam-c* gene family expansion in the non-Katangan *P. vinckei* subspecies

546     seems to have been driven by both gene duplications in their common ancestor and

547     subspecies-specific gene family expansions subsequent to subspeciation. Around

22

548    half of the *fam-c* genes have detectable transcripts in asexual or sexual blood

549    stages. Most of the transcribed genes have medium (36<FPKM<256) or high-level

550    expression (FPKM > 256) and blood-stage specific expression data for *P. chabaudi*

551    and *P. v. vinckei* show peak transcription among the asexual blood stages at ring

552    and schizont stages.

553

554    **Genetic crossing can be performed between *P. vinckei* isolates**

555    The availability of several isolates within each *P. vinckei* subspecies with varying

556    growth rates and wide genetic diversity makes them well-suited for genetic studies.

557    Therefore, we attempted genetic crossing of the two *P. vinckei baforti* isolates,

558    *Pvs*EH and *Pvs*EL, that displayed differences in their growth rates. Optimal

559    transmission temperature and vector stages were initially characterized for *P. v.*

560    *baforti* EE, EH and EL. Each isolate was inoculated into three CBA mice and on day

561    3 post infection, around 100 female *A. stephensi* mosquitoes were allowed to

562    engorge on each mouse at different temperatures - 21°C, 23°C and 26°C. All three

563    *P. v. baforti* isolates were able establish infections in mosquitoes at 23°C and 26°C,

564    producing at least 50 mature oocysts on day 15 post-feed, but failed to transmit at

565    21°C (Figure 1A (a) and Additional file 6). Four to five oocysts of 12.5-17.5 μm

566    diameter were observed at day 8 post-feed in the mosquito midgut and around a

567    hundred mature oocysts of 50 um diameter could be observed at day 15 post-feed.

568    Some of these mature oocysts had progressed into sporozoites but only a very few

569    appeared upon disruption of the salivary glands.

570

571    To perform a genetic cross between *Pvs*EH and *Pvs*EL, a mixed inoculum containing

572    equal proportions of *Pvs*EH and *Pvs*EL parasites was injected into CBA mice and a

573    mosquito feed was performed on both day 3 and day 4 post-infection to increase the

574    chances of a successful transmission (Additional 7 B). For each feed, around 160

575    female *A. stephensi* mosquitoes were allowed to take a blood meal from two

576    anaesthetized mice at 24°C for 40 minutes without interruption.

577

578    Upon inspection of mosquito midguts for the presence of oocysts on day 9 post-feed,

579    100% infection was observed (all midguts inspected contained oocysts) for both day

580    3 and day 4 feeds. Around 25-100 oocysts were found per midgut in day 3 fed

581    mosquitoes and 5-40 oocysts per midgut in day 4 fed mosquitoes. On day 12 post-

582    feed, mature oocysts and also a high number of sporozoites were found in the

583    midguts, but upon disrupting the salivary glands on day 20 post-feed, only a few

584    sporozoites were found in the suspension.

585

586    Sporozoites from day 3 and 4 fed mosquitoes were injected into ICR mice (D3 and

587    D4 respectively) and five days later, both mice became positive for blood stage

588    parasites. In order to confirm that a genetic cross has taken place, four clones were

589    obtained from D4 by limiting dilution to screen for presence of both *Pvs*EH and

590    *Pvs*EL alleles within the chromosomes. Based on the SNPs identified between

591    *Pvs*EH and *Pvs*EL, we amplified 600 to 1,000 bp regions from polymorphic genes on

592    both ends of the 14 chromosomes that contained isolate-specific SNPs and

593    performed Sanger sequencing of the amplicons (primer sequences in Additional file

594    6).

595

596    Both *Pvs*EL-specific (11) and *Pvs*EH-specific (17) markers were found in the 28

597    markers sequenced (one marker, PVSEL_0600390, could not be amplified). Also,

598     five chromosomes clearly showed evidence of chromosomal cross-over since they

599     contained markers from both isolates (see Figure 4A), thus confirming a successful

600     *P. vinckei* genetic cross. However, all four clones had the same pattern of

601     recombination which suggests that the diversity of recombinants in the cross-

602     progeny was low and a single recombinant parasite might have undergone

603     significant clonal expansion.

604

605     ***P. vinckei* parasites are amenable to genetic manipulation**

606     We asked if *P. vinckei* parasites can be genetically modified by applying existing

607     transfection and genetic modification techniques routinely used in other RMPs.

608     *Plasmodium v. vinckei* CY was chosen to test this because the isolate naturally

609     established a synchronous infection in mice and reaches a high parasitaemia, which

610     results in an abundance of schizonts for transfection. We aimed to produce a *Pvv*CY

611     line that constitutively expresses GFPLuc (green fluorescent protein- firefly

612     luciferase) fusion protein, similar to those produced in *P. berghei* and *P. yoelii* [52,

613     53]. A recombination plasmid, *pPvvCY-Δp230p-gfpLuc*, was constructed to target

614     and replace the dispensable wildtype P230p locus in *P. v. vinckei* CY

615     (PVVCY_0300700) with a gene cassette encoding for GFPLuc and a *hdhfr*

616     selectable marker cassette (Figure 4B).

617

618     Transfection of purified *Pvv*CY schizonts with 20 µg of linearized p*Pvv*CY-Δp230p-

619     gfpLuc plasmid by electroporation, followed by marker selection using

620     pyrimethamine yielded pyrimethamine-resistant transfectant parasites (PvGFP-

621     Luc$_{con}$) on day 6 after drug treatment. Stable transfectants were cloned by limiting

622     dilution and plasmid integration in these clones was confirmed by PCR. Constitutive

623     expression of GFPLuc in PvGFP-Luc$_{con}$ asexual and sexual blood stage parasites

624     was confirmed by fluorescence live cell imaging (Figure 4C). GFPLuc expression in

625     PvGFP-Luc$_{con}$ oocysts was confirmed by fluorescence imaging of mosquito midguts

626     7 days after blood meal.

627

## Discussion

629     Of the four RMP species that have been adapted to laboratory mice, *P. berghei*, *P.*

630     *yoelii* and *P. chabaudi* have been extensively used to investigate malaria parasite

631     biology. Adopting these RMPs as tractable experimental models has been facilitated

632     by continuous efforts in characterizing their phenotypes, sequencing their genomes

633     and establishing protocols for parasite maintenance, genetic crossing and genetic

634     modification. Here, we extend these efforts to *Plasmodium vinckei*.

635

636     We have systematically studied ten *P. vinckei* isolates and produced a

637     comprehensive resource of their reference genomes, transcriptomes, genotypes and

638     phenotypes to help establish *P. vinckei* as a useful additional experimental model for

639     malaria.

640

641     Enzyme variation and molecular phylogeny studies indicate that the five subspecies

642     of *P. vinckei* have diverged significantly from each probably due to the geographical

643     isolation of these parasites in different locations around the African Congo basin.

644     This diversity calls for a reference genome for each subspecies in order to capture

645     large-scale changes in their genomes such as chromosomal structural variations and

646     gene copy number variations that might have played a role in their subspeciation. To

647     accurately capture these events, we used a combination of Pacbio and Illumina

26

648  sequencing that allowed us to produce an end-to-end assembly of *P. vinckei*

649  chromosomes. This, coupled with manual curation of the predicted gene models, led

650  to the creation of five high-quality reference genomes for *P. vinckei* that are a

651  significant improvement to the existing fragmented genomes available for *P. v.*

652  *vinckei* and *P. v. petteri*.

653

654  Comparative synteny analysis between *P. vinckei* and other RMP genomes reveals

655  structural variations at both the species and the subspecies levels. Assuming that

656  the observed variations have occurred only once, a putative pathway of genome

657  rearrangements during RMP evolution can be inferred. No rearrangements have

658  occurred during *P. berghei* and *P. yoelii* speciation and their genomes are likely to be

659  identical to the RMP ancestor [41]. A reciprocal translocation between chromosomes

660  VIII and X has accompanied the speciation of *P. vinckei*, and this is mutually

661  exclusive from the reciprocal translocation between chromosome VII and IX that has

662  occurred during *P. chabaudi* speciation. Following this, there has been a small

663  inversion in chromosome X during the subspeciation of *P. v. vinckei* and

664  translocations between chromosomes V, VI and XIII during the subspeciation of *P. v.*

665  *petteri* (which are then carried over to *P. v. baforti*).

666

667  We generated additional sequencing data for several *P. vinckei* isolates and made

668  available at least two genotypes per *P. vinckei* subspecies (except for *P. v. vinckei*

669  for which only one isolate is available) so as to facilitate future studies that might

670  employ *P. vinckei* parasites to study phenotype-genotype relationships. Similarly, we

671  also supplemented the existing genotype information for other RMPs by sequencing

672  several isolates from additional subspecies of *P. chabaudi* and *P. yoelii*. Our data

673     thus comprises of genotypes from sympatric species from each region of isolation

674     allowing us to re-evaluate the genotypic diversity and evolution among RMP isolates.

675

676     A genome-wide SNP-based phylogeny shows that the divergences between different

677     subspecies are proportional to the level of isolation of the habitat for all RMP

678     species. *Plasmodium vinckei*, *P. yoelii* and *P. chabaudi* isolates from sites in

679     Cameroon have very similar genotypes to their counterparts in the Central African

680     Republic denoting similar evolutionary pressures and perhaps the presence of gene

681     flow across these regions,   while isolates from Brazzaville (Congo) are more

682     diverged probably due to the different environmental conditions in these locations

683     [40].

684

685     Subspecies from West Nigeria and the DRC are highly diverged compared to

686     subspecies from the rest of Africa. The distinctiveness of *P. berghei* and *P. v.*

687     *vinckei,* both from the DRC is most likely due to climactic and host-vector differences

688     in the highland forests of Katanga. Highland forests are an altitude of 1000-7000 m

689     with mean temperature of 21C   whereas the lowland forests lie at an altitude less

690     than 800m with a mean temperature of 25 C. Different host-vector systems are

691     prevalent in the lowland forests (*Grammomys poensis (previously known as*

692     *Thamnomys rutilans)* - specific mosquito species unknown) and the highland

693     Katangan forests (*Grammomys surdaster -Anopheles dureni millecampsi*). The

694     associated selection pressures seem to have mainly influenced their transmission,

695     as reflected by their lower optimal transmission temperatures and the high Ka/Ks

696     ratios observed for three proteins that play critical functions in this process. Recently,

697     several more rodent host and mosquito vector species have been identified in the

698  forests of Gabon [54] implying that a diverse set of host-vector systems could have

699  existed for RMPs. Thus, diversification of RMP species into several subspecies

700  within these isolated ecological niches might have been driven by evolutionary forces

701  resulting from the diverse host, vector and environmental conditions experienced at

702  each locale.

703

704  Malaria parasite genomes contain several highly polymorphic multigene families

705  located in the sub-telomeric chromosomal regions that encode a variety of exported

706  proteins involved in processes such as immune evasion, cytoadherence, nutrient

707  uptake and membrane synthesis. Multigene families are thought to have evolved

708  rapidly under the influence of immune and other evolutionary pressures resulting in

709  copy number variations and rampant sequence reshuffling that ultimately leads to

710  phenotypic plasticity in *Plasmodium*.

711

712  Previously, phylogenetic analyses of *pir, fam-a and fam-b* genes from three RMP

713  species have shown that structurally distinct genes exist within these families

714  forming robust clades with varied levels of orthology/paralogy. Identifying sub-

715  families that have structurally diversified within the multigene families can help to

716  better understand their functions and to this end, we constructed phylogenetic trees

717  for the ten multigene families with genes from all four RMP species. Due to the scale

718  of the analysis, we applied automated trimming to our alignments and limited our

719  tree inference method to maximum likelihood. While this resulted in poor bootstrap

720  values for some clades in the *pir*, *fam-a* and *fam-b* trees compared to previous

721  phylogenetic analyses [25, 26], our method was able to retrieve similar tree

722 topologies to those previously inferred and in general produced trees with good

723 nodal support for the rest of the multigene families.

724

725 Robust pan-RMP clades identified in our study represent ancestral lineages

726 consisting of structural orthologs that perform conserved functions across all RMPs

727 and will be useful for future work with these families. We show that certain ancestral

728 lineages can expand in a particular species or subspecies in response to selective

729 pressures resulting in distinct evolutionary histories for each family. For example, the

730 *pir* family expansion is mostly species-specific and driven by frequent gene

731 conversion after speciation, whereas the expansion of the *fam-a* gene repertoire

732 seems to have occurred initially in the RMP ancestor followed by species-specific

733 expansions.

734

735 Inclusion of *P. vinckei pirs* in the RMP *pir* family phylogeny show that *P. vinckei pir*s

736 do not form independent clades of their own and instead populate three *P. chabaudi*-

737 dominant clades. This suggests that some of the *pir* clades were established earlier

738 on when the classical *vinckei* and *berghei* group of parasites split from their common

739 RMP ancestor resulting in *vinckei* group-specific clades like L1, S7 and S1g.

740

741 Similarly, the addition of *P. vinckei* genes resolved the ancestral clade of internal

742 *fam-a* genes into several well-supported *vinckei* group clades and a *berghei* group

743 clade. We observe similarly high level of orthology between *P. chabaudi* and *P.*

744 *vinckei* genes in other multigene families forming several *vinckei* group clades in

745 contrast to more species-specific clades of paralogous genes in *P. berghei* and *P.*

746 *yoelii*. For example, within the *fam-d* family, five ancestral lineages can be identified

747  in the *vinckei* group as opposed to only one paralogous *P. yoelii*-specific expansion

748  within the *berghei* group. Taken together, it seems that family expansions in *P.*

749  *chabaudi* and *P. vinckei* have occurred in the common *vinckei* group ancestor prior

750  to speciation and that multigene families have evolved quite differently across the

751  *vinckei* and *berghei* groups of RMPs. These might be related to the striking

752  differences in the basic phenotypes of these two groups of parasites.

753

754  We also observed size expansions in the *ema1* and *fam-c* families within the non-

755  Katangan *P. vinckei* parasites, all being isolates from the lowlands around the Congo

756  Basin. *Ema1* family expansions seem to be specific to lowlands dwelling *vinckei*

757  group parasites as they are expanded in both non-Katangan *P. vinckei* and *P.*

758  *chabaudi.* However, unlike *P. chabaudi*, the duplicated gene members in non-

759  Katangan *P. vinckei* are all pseudogenized by a S5X mutation effectively rendering

760  the functional repertoire to be just 6-8 genes, similar to highlands dwelling Katangan

761  *P. v. vinckei*. Thus, it could be speculated that even under similar selective

762  pressures, *ema1* family expansions contribute to parasite fitness in *P. chabaudi* but

763  may not be required for the survival of sympatric *P. vinckei* parasites. The *P. vinckei*

764  *ema1* pseudogenes could still serve as silent donor genes that recombine into

765  functional variants to bring about antigenic variation [55]. In the case of the *fam-c*

766  gene family, the expansion is specific to non-Katangan *P. vinckei* subspecies since

767  *P. chabaudi*, *P. yoelii* and *P. v. vinckei* all have similar repertoire sizes. The

768  expansions seem to be driven by gene duplications initially in their non-Katangan

769  common ancestor and again after subspeciation.

770

771     The effect of the difference in habitats is even more pronounced in the Katangan

772     parasite, *P. v. vinckei*. It has a smaller genome and a compact multigene family

773     repertoire reminiscent of the only other Katangan isolate, *P. berghei* and its genetic

774     distance from other members of the *P. vinckei* clade is in the same order of

775     magnitude as that between separate species within the RMPs. The reduced

776     multigene family repertoire mainly consists of members belonging to pan-RMP or

777     *vinckei*-group specific ancestral lineages making it an ideal *vinckei* group parasite to

778     study the localization and function of variant proteins.

779

780     We tested whether *Pvv*CY was amenable to genetic manipulation using standard

781     transfection protocols already established for other RMPs. We were able to

782     successfully knock-in a GFP-luciferase fusion cassette to *Pvv*CY to produce a GFP-

783     Luc reporter line for *P. v. vinckei,* following a transfection protocol routinely used for

784     modifying *P. yoelii* in our lab [44, 56]. We were able to visualise GFP-positive

785     parasites during different blood stages and in oocysts thus confirming stable GFPLuc

786     expression. We were unable to visualise other life stages (sporozoites and liver

787     stages) due to our failure to produce viable salivary gland sporozoites in this

788     parasite.

789

790     The transfection of *P. chabaudi* has been challenging due to its slow proliferation

791     rate and schizont sequestration resulting in low merozoite yield, thus necessitating

792     optimized transfection protocols. In contrast, *Plasmodium v. vinckei* reaches high

793     parasitaemia without being immediately lethal to the host (90% parasitaemia on day

794     6) and is highly synchronous yielding a large number of schizonts. A predominant

795     population of schizonts appear near midnight in *P. v. vinckei* infections, at which

796  point, they can be Percoll-purified from exsanguinated blood and transfected with

797  DNA.

798

799  *P. vinckei* and *P. chabaudi*, while being distinct species, share several

800  characteristics that are common among *vinckei* group RMPs, such as a predilection

801  for mature erythrocytes, synchronous infections and the sequestration of schizonts

802  from peripheral circulation [33, 37, 57-59]. Thus, *P. v. vinckei* can serve as an ideal

803  experimental model for functional studies targeting these aspects of parasite biology.

804

805  The availability of several RMP isolates with phenotypic differences aids their use in

806  study of parasite fitness and transmission success in mixed infections [60, 61] and

807  for the identification of genes involved in parasite virulence, strain-specific immunity,

808  drug resistance and host-cell preference using genetic crosses [44, 62-64]. With this

809  in mind, we studied the virulence of ten *P. vinckei* isolates to identify differences in

810  their growth rate and their effect on the host.

811

812  Some of these isolates have been previously characterized [40], but we

813  systematically profiled additional representative isolates for each subspecies (where

814  available) under comparative conditions in the same host strain. We identified pairs

815  of isolates with contrasting virulence phenotypes within two *P. vinckei* subspecies –

816  *P. v. petteri* (*Pvp*CR and *Pvp*BS) and *P. v. baforti* (*Pvs*EH and *Pvs*EL or *Pvs*EE).

817  These isolate pairs would be ideal candidates for studies utilising genetic crossing to

818  identify genetic *loci* linked to virulence using Linkage Group Selection [44].

819

820  Since *P. vinckei* subspecies have significantly diverged from each other, isolates

821  within the same subspecies are more likely to recombine than isolates from different

822  subspecies. However, intra-specific hybrids between *P. v. petteri* and *P. v. baforti*

823  may also be possible (as demonstrated earlier in *P. yoelii* [65]) since these two

824  subspecies are closely related (see Figure 2C). However, difficulties in transmitting

825  *P. vinckei* parasites have been reported previously with either the gametocytes

826  failing to produce midgut infections or sporozoites failing to invade the salivary

827  glands or infections resulting in non-infective sporozoites [27, 30, 35]. Repeated

828  attempts to create a cross between two *P. vinckei baforti* isolates failed to produce

829  any detectable recombinants due to low frequency of mosquito transmission [35].

830  Here, we renewed these efforts with different *P. vinckei* isolates to see if we could

831  establish a *P. vinckei* genetic cross. Two attempts were made to create a *pvp*CR X

832  *pvp*BS cross and further two attempts were made to create a *pvs*EL X *pvs*EH cross.

833  However, in all attempts the sporozoites failed to optimally invade the salivary glands

834  and we managed to isolate only a few in the *P. v. subsp* cross, subsequently

835  obtaining a cross progeny in mice. While we were able to demonstrate a successful

836  genetic cross by showing the presence of alleles from both isolates in the cross

837  progeny, the recombinant diversity was quite low probably due to the transmission

838  bottleneck. We are currently further investigating the optimal conditions for

839  transmitting *P. vinckei*.

840

## Conclusions

842  In this study, we have created a comprehensive resource for the rodent malaria

843  parasite *Plasmodium vinckei*, comprising of five high-quality reference genomes, and

844  blood stage-specific transcriptomes, genotypes and phenotypes for ten isolates. We

845    have employed state-of-the-art sequencing technologies to produce largely complete

846    genome assemblies and highly accurate gene models that were manually polished

847    based on strand-specific RNA sequencing data. The unfragmented nature of our

848    genome assemblies allowed us to characterize structural variations within *P. vinckei*

849    subspecies, which, to the best of our knowledge, is the first time that large-scale

850    genome re-arrangements have been found among subspecies of a *Plasmodium*

851    species.

852

853    The biological or phenotypic significance, if any, of such alterations are poorly

854    understood, but it seems likely that they may drive speciation through the promotion

855    of reproductive isolation of species or subspecies. Through our extensive

856    sequencing efforts, we have generated genotype data for seventeen RMP isolates

857    comprising of five *P. vinckei*, four *P. yoelii* and one *P. chabaudi* subspecies, thus

858    making at least one genotype available for all subspecies of the RMP that previously

859    lacked any sequencing data. We also systematically characterised the virulence

860    phenotypes of the ten *P. vinckei* isolates to capture the phenotypic diversity among

861    them. Combined, these efforts will greatly aid genetic linkage studies to resolve

862    genotype-phenotype relationships.

863

864    In order to understand the evolutionary relationships among the RMP isolates, we

865    have carried out a combination of analyses to describe the genotypic diversity

866    molecular evolution of these parasites. While our phylogenies more or less agree

867    with previous biochemical and molecular data-based studies, our reconstruction

868    based on sequence variations on a genome scale provides higher resolution to the

869    divergence estimates. Taking advantage of the high-quality RMP genomes produced

870 from our work and previous studies, we also undertook a comprehensive

871 phylogenetic analysis of multigene families across all RMP species and identify

872 various structurally diversified sub-families with distinct evolutionary histories. This

873 will enable future studies on the critical role of multigene families in parasite

874 adaptation, and to aid this, we have made searchable and interactive versions of the

875 phylogenies publicly available through the iTOL online tool [66].

876

877 While genome rearrangements have occurred during speciation and sub speciation

878 events, diversification of the multigene families seem to have occurred earlier when

879 the RMPs split into *vinckei* and *berghei* groups of parasites. Thus, structural, copy

880 number and nucleotide-level variations among the RMPs have occurred at various

881 points during the evolution of RMPs in response to a variety of evolutionary

882 pressures. The gene expression data from our study, covering specific blood stages

883 for some *P. vinckei* subspecies, show conserved expression of multigene family

884 members across RMPs. While not comprehensive, it complements existing RMP

885 transcriptomes and will aid functional studies in the *P. vinckei* model. Taken

886 together, our study provides a comprehensive view of the phenotypic and genotypic

887 diversity within RMPs and functional diversification of the multigene families in

888 response to selection pressures.

889

890 The synchronicity of *P. v. vinckei* infection and its unique ability to sustain high

891 parasitaemia without killing its host culminating in good schizont yields make this

892 parasite an attractive model for reverse genetics studies, especially those on

893 multigene families owing to its reduced repertoire. We have successfully

894 demonstrated genetic manipulation in *P. v. vinckei* but encountered difficulties in

895    producing large numbers of recombinant parasites through genetic crossing.

896    Attempts to transmit isolates from three different *P. vinckei* subspecies in *A.*

897    *stephensi* mosquitoes failed in our hands as sporozoites repeatedly failed to infect

898    the salivary glands. Careful optimisation of transmission parameters and serial

899    mosquito passages of the *P. vinckei* parasites might help in improving their

900    transmission efficiency and could aid genetic linkage studies with these parasites.

901

902

903

904    # Methods

905    **Parasite lines and experiments using mice and mosquitos**

906    The parasite lines used in this study and their original isolate information are detailed

907    in Supplementary Table 1. Frozen parasite stabilates of cloned or uncloned lines

908    were revived and inoculated intravenously into ICR mice. Five *P. vinckei* isolates

909    (*Pvv*CY, *Pvb*DA, *Pvb*DB, *Pvl*DE and *Pvs*EE) and the *P. yoelii nigeriensis* isolate

910    (*Pyn*D) were uncloned stabilates and were cloned by limiting dilution to obtain clonal

911    parasite lines.

912

913    Laboratory animal experimentation was performed in strict accordance with the

914    Japanese Humane Treatment and Management of Animals Law (Law No. 105 dated

915    19 October 1973 modified on 2 June 2006), and the Regulation on Animal

916    Experimentation at Nagasaki University, Japan. The protocol was approved by the

917    Institutional Animal Research Committee of Nagasaki University (permit:

918    12072610052).

919

920   Six to eight weeks old female ICR or CBA mice were used in all the experiments.

921   The

922   mice were housed at 23°C and maintained on a diet of mouse feed and water. Mice

923   infected with malaria parasites were given 0.05% para-aminobenzoic acid (PABA)-

924   supplemented water to assist parasite growth.

925

926   All mosquito transmission experiments were performed using *Anopheles stephensi*

927   mosquitoes were housed in a temperature and humidity-controlled insectary at 24°C

928   and 70% humidity. Mosquito larvae were fed with mouse feed and yeast mixture and

929   adult mosquitoes were maintained on 10% glucose solution supplemented with

930   0.05% PABA.

931

932   **Parasite growth profiling**

933   For each isolate, an inoculum containing 1 X $10^6$ parasitized RBCs was injected

934   intravenously to five CBA mice. Blood smears, haematocrit readings (Beckman

935   Coulter Counter) and body weight readings were taken daily for 20 days or until host

936   mortality to monitor parasitaemia, anaemia and weight loss. Blood smears were fixed

937   with 100% methanol and stained with Geimsa's solution. The average parasitaemia

938   was calculated from parasite and total RBC counts taken at three independent

939   microscopic fields.

940

941   **Genomic DNA isolation and whole genome sequencing**

942   Parasitized whole blood was collected from the brachial arteries of infected mice and

943   blood sera was removed by centrifugation. RBC pellets were washed once with PBS

944   and leukocyte-depleted using CF11 (Sigma Cat# C6288) cellulose columns. Parasite

945    pellets were obtained by gentle lysis of RBCs with 0.15% saponin solution. Genomic

946    DNA extraction from the parasite pellet was performed using DNAzol reagent

947    (Invitrogen CAT # 10503027) as per manufacturer's instructions.

948

949    Single-molecule sequencing was performed for five *P. vinckei* isolates. 5-10 ug of

950    gDNA was sheared using a Covaris g-TUBE shearing device to obtain target sizes of

951    20kB (for *Pvv*CY, *Pvb*DA and *Pvp*CR) and 10kB (for samples *Pvl*DE and *Pvs*EL).

952    Sheared DNA was concentrated using AMPure magnetic beads and SMRTbell

953    template libraries were generated as per Pacific Biosciences instructions. Libraries

954    were sequenced using P6 polymerase and chemistry version 4 (P6C4) on 3-6 SMRT

955    cells and sequenced on a PacBio RS II. Reads were filtered using SMRT portal v2.2

956    with default parameters. Read yields were 352,693, 356,960, 765,596, 386,746 and

957    675,879 reads for *Pvv*CY, *Pvb*DA, *Pvl*DE, *Pvp*CR and *Pvs*EL respectively totalling

958    around 2.7 to 4.7 Gb per sample. Mean subread lengths ranged from 6.15 to 9.1 kB.

959    N50 of 11.7 kB and 19.2 kB were obtained for 10 and 20 kB libraries respectively.

960

961    PCR-free Illumina sequencing was performed for all RMP isolates. 1-2 ug of DNA

962    was sheared using Covaris E series to obtain fragment sizes of 350 and 550bp.

963    350bp and 550bp PCR-free libraries were prepared using TruSeq PCR-free DNA

964    library preparation kits according to the manufacturer's instructions. Libraries were

965    sequenced on the Illumina HiSeq2000 platform with 2 X 100bp paired-end read

966    chemistry. Read yields ranged from 8-22 million reads for each library (see

967    Additional File 1).

968

969    **Genome assembly and annotation**

970 Genome assembly from long single molecule sequencing reads was performed

971 using FALCON (v0.2.1)[67] with length cutoff for seed reads used for initial mapping

972 set as

973 2,000bp and for pre-assembly set as 12,000bp. The falcon sense options were set

974 as- "–min idt 0.70 –min cov 4 –local match count threshold 2 –max n read 200"

975 and overlap filtering settings were set as "–max diff 240 –max cov 360 –min cov

976 5 –bestn 10". 28-40 unitigs were obtained and smaller unitigs were discarded as

977 they were exact copies of the regions already present in the larger unitigs.

978

979 PCR-free reads were used to correct base call errors in the unitigs using ICORN2

980 [68], run with default settings and for 15 iterations. The unitigs were classified as

981 chromosomes based on their homology with *P. chabaudi* chromosomes (GeneDB

982 version 3). In *Pvl*DE and *Pvs*EL samples, some of the chromosomes were made of

983 two to three unitigs with overlapping ends which were then fused and the gaps were

984 removed manually. Apicoplast and mitochondrial genomes were assembled from

985 PCR-free reads alone using Velvet assembler [69].

986

987 Syntenic regions between genome sequences were identified using MUMmer v3.2

988 [70]. Synteny breakpoints were identified manually and were confirmed not to be

989 misassemblies by verifying that they had continuous read coverage from PacBio and

990 Illumina reads. Artemis Comparison tool [71, 72] and Integrative Genomics Viewer

991 [73] were used for this purpose. The structural variations were illustrated using

992 CIRCOS [74].

993 *De novo* gene predictions were made using AUGUSTUS [75] trained on *P. chabaudi*

994 gene models. RNA sequencing reads were mapped onto the reference genome

995 using TopHat [76] to infer splice junctions. AUGUSTUS predicted gene models,

996 junctions.bed file from TopHat and *P. chabaudi* gene models were fed into MAKER

997 [77] to create consensus gene models that were then manually curated based on

998 RNAseq evidence in Artemis Viewer and Artemis Comparison tool [71, 72].

999 Ribosomal RNA (rRNA) and transfer RNA (tRNA) were annotated using RNAmmer

1000 v1.2 [78]. Gene product calls were assigned to *P. vinckei* gene models based on

1001 above identified orthologous groups using custom scripts. Functional domain

1002 annotations were inferred from InterPro database using InterProScan v5.17 [79].

1003 Transmembrane domains were predicted by TMHMMv2.0 [80], signal peptide

1004 cleavage sites by SignalP v4.0 [81], presence of PEXEL/VTS motif detected using

1005 ExportPredv4.0 [82] (with PEXEL score cutoff of 4.3).

1006

1007 **Transcriptomics**

1008 Total RNA was isolated for four *P. vinckei* isolates (*Pvb*DA, *Pvp*CR, *Pvl*DS and

1009 *Pvs*EL) from mixed blood stages using TRIzol (Invitrogen) following the

1010 manufacturer's protocol. For *Pvp*CR and *Pvl*DS, additionally, total RNA was isolated

1011 from ring, trophozoite and gametocyte enriched fractions obtained using a Nycodenz

1012 gradient.

1013 Strand-specific mRNA sequencing was performed from total RNA using TruSeq

1014 Stranded mRNA Sample Prep Kit LT (Illumina) according to the manufacturer's

1015 instructions. Briefly, polyA+ mRNA was purified from total RNA using oligo-dT

1016 dynabead selection. First strand cDNA was synthesised using randomly primed

1017 oligos followed by second strand synthesis where dUTPs were incorporated to

1018 achieve strand-specificity. The cDNA was adapter-ligated and the libraries amplified

1019    by PCR. Libraries were sequenced in Illumina Hiseq2000 with paired-end 100bp

1020    read chemistry.

1021

1022    Stage-specific RNAseq data for *Pvv*CY's intraerythrocytic growth stages were

1023    obtained from an earlier study [50]. Gene expression was captured every 6 hours

1024    during *Pvv*CY's 24 h IDC with three replicates, of which 6h, 12h and 24h timepoints

1025    were used in this study to denote gene expression at ring, trophozoite and schizont

1026    stages respectively. Similarly, for *P. chabaudi* AS, gene expression was captured

1027    every 3h during its IDC with two replicates in a recent study [56], of which the 5.5h,

1028    11.5h and 23.5 h timepoints on day 2 were chosen to denote ring, trophozoite and

1029    schizont stages respectively. *P. yoelii* and *P. berghei* transcriptome data were

1030    obtained from [26] and [44] respectively.

1031

1032    **SNP calling and molecular evolution analysis**

1033    Illumina paired-end reads for a total of 30 RMP isolates produced in this study or

1034    sourced from previous studies (see Additional File 13) were used for SNP calling. In

1035    the case of isolates sequenced in this study, the 350bp fragment size PCR-free

1036    sequencing data was used. First, to produce a high quality pan-RMP SNP dataset

1037    for phylogeny construction, all quality-trimmed reads were mapped onto the *Pvv*CY

1038    reference genome using BWA tool [83] with default parameters. MAPQ values of the

1039    mapped reads were fixed and duplicated reads removed using CleanSam,

1040    FixMateInformation and MarkDuplicates commands in picardtools

1041    (http://broadinstitute.github.io/picard) and only uniquely mapped reads were retained

1042    using samtools with parameter -q 1 (http://www.htslib.org/). Raw SNPs were called

1043    from the mapped reads using samtools mpileup and bcftools with following

1044  parameters- minimum base quality of 20, minimum mapping quality of 10 and ploidy

1045  of 1. SNPs with quality (QUAL) less than 20, read depth (DP) less than 10, mapping

1046  quality (MQ) less than 2 and allele frequency (AF1) less than 80% were removed.

1047  Further, only SNPs present in protein-coding genes were retained and those present

1048  in low-complexity regions (predicted by DustMasker [84]) and sub-telomeric

1049  multigene family members were excluded.

1050

1051  The filtered SNPs from different samples were merged and SNP positions with

1052  missing calls in more than six samples were removed. This filtered high-quality set of

1053  1,020,956 SNP positions were used to infer maximum likelihood phylogeny (see

1054  Additional File 7).

1055

1056  For inferring Ka/Ks ratios between *P. vinckei* isolates and *Pvv*CY, filtered SNPs

1057  obtained above were merged as before but excluding *Pvp*BS due to its high missing

1058  call rate. Only SNP positions with no missing calls in any sample were retained and

1059  morphed onto *Pvv*CY gene sequences using gatk command

1060  FastaAlternateReferenceMaker [85] to produce isolate-specific gene sequences

1061  which were then used for pairwise sequence comparisons to identify synonymous

1062  and non-synonymous substitutions. Ka/Ks ratios were calculated using KaKs

1063  Calculator [86] and averaged across isolates if more than one was available for a

1064  subspecies.

1065

1066  For comparisons against *P. v. petteri*, *P. yoelii* and *P. chabaudi*, sample reads were

1067  mapped onto *Pvp*CR, *P. yoelii* 17X and *P. chabaudi* AS genomes respectively and

1068  subsequent steps were followed as before. Similar to *Pvp*BS, *Pys*EL was excluded

1069  from Ka/Ks analysis due to high missing rate.

1070

1071  **Phylogenetic analysis**

1072  For constructing species-level phylogenies, orthologous proteins were identified

1073  between the five *P. vinckei* genomes, three RMP genomes, *P. falciparum, P.*

1074  *knowlesi* and *P. vivax* genomes using OrthoMCL v2.0.9 [87] with inflation parameter

1075  as 1.5, BLAST hit evalue cutoff as $1e^{-5}$ and percentage match cutoff as 50%.

1076

1077  One-to-one orthologous proteins from each of the 3,920 ortholog groups that form

1078  the core proteome were aligned using MUSCLE [88]. Alignments were trimmed

1079  using trimAl [89] removing all gaps and concatenated into a partitioned alignment

1080  using catsequence. An initial RAxML [90] run was performed on individual

1081  alignments to identify best amino acid substitution model under the Akaike

1082  Information Criterion (--auto-prot=aic). These models were then used to run a

1083  partitioned RAxML analysis on the concatenated protein alignment using

1084  PROTGAMMA model for rate heterogeneity.

1085

1086  For constructing isolate-level phylogeny, the vcf files containing high-quality SNPs

1087  were first converted to a matrix for phylogenetic analysis using vcf2phylip

1088  (https://github.com/edgardomortiz/vcf2phylip). RAxML tree inference was performed

1089  using GTRGAMMA model for rate heterogeneity along with ascertainment bias

1090  correction (--asc-corr=stamatakis) since we used only variant sites.

1091

1092 Maximum likelihood trees for multigene families were constructed based on

1093 nucleotide sequence alignments of member genes that included intron sequences if

1094 present (except in the case of *pir* family where introns were excluded). Alignments

1095 were performed using MUSCLE with default parameters, frame-shifts edited

1096 manually in AliView [91] followed by automated trimming with trimAl using -gappyout

1097 parameter. In all the phylogenies, bootstrapping was conducted until the majority-

1098 rule consensus tree criterion (-I autoMRE) was satisfied (usually 150-300 replicates).

1099 Phylogenetic trees were visualized and annotated in the iTOL server [66].

1100

1101 **Plasmid construction and transfection in *P. vinckei***

1102 The p*Pvv*CY-p230p-gfpLuc plasmid was constructed using MultiSite Gateway

1103 cloning

1104 system (Invitrogen). attB-flanked 5'and 3'homology arms were obtained by

1105 amplifying

1106 800bp regions upstream and downstream of PVVCY_0300700. These fragments

1107 were subjected to independent BP recombination with pDONRP4-P1R (Invitrogen) to

1108 generate entry plasmids pENT12-5U and pENT41-3U, respectively. Similarly, the

1109 gfpLuc cassette from pL1063 was amplified and subjected to LR reaction to obtain

1110 pENT23-gfpLuc. BP reaction was performed using the BP Clonase II enzyme mix

1111 (Invitrogen) according to the manufacturer's instructions.

1112

1113 *Plasmodium vinckei vinckei* CY schizont-enriched fraction was collected by

1114 differential centrifugation on 50% Nycodenz in incomplete RPMI1640 medium, and

1115 20 ug of ApaI- and StuI-double digested linearized transfection constructs were

1116 electroporated to $1 \times 10^7$ of enriched schizonts using a Nucleofector device (Amaxa)

1117   with human T-cell solution under program U-33. Transfected parasites were

1118   intravenously injected into

1119   7-week-old ICR female mice, which were treated by administering pyrimethamine in

1120   the drinking water (0.07 mg/mL) 24 hours later for a period of 4-7 days. Drug

1121   resistant parasites were cloned by limiting dilution with an inoculum of 0.3

1122   parasites/100 uL injected into 10 female ICR mice. Two clones were obtained, and

1123   integration of the transfection constructs was confirmed by PCR amplification with a

1124   unique set of primers for the modified p230p gene locus. Live imaging of parasites

1125   was performed on thin smears of parasite-infected blood prepared on glass slides

1126   stained with Hoechst 33342. Fluorescent and differential interference contrast (DIC)

1127   images were captured using an AxioCam MRm CCD camera (Carl Zeiss, Germany)

1128   fixed to an Axio imager Z2 fluorescent microscope with a Plan-Apochromat 100 ×/1.4

1129   oil immersion lens (Carl Zeiss) and Axiovision software (Carl Zeiss). GFP-expressing

1130   *P. vinckei* oocysts in mosquito midguts were imaged in SMZ25 microscope (Nikon).

1131

1132   **Mosquito transmission and genetic crossing of *P. vinckei* parasites**

1133   To determine the optimal transmission temperature for *P. vinckei baforti* isolates,

1134   infected CBA mice were anaesthetized on day 3 post-inoculation and ~100 female

1135   *Anopheles stephensi* mosquitoes (7 to 12 days post emergence) were allowed to

1136   take a blood meal for 30 min without interruption after confirming presence of

1137   gametocytes by microscopy. Three batches of ~100 mosquitoes were fed at three

1138   different temperatures - 21°C, 23°C and 26°C. The fed mosquitoes were maintained

1139   at the feed temperatures and at 70% humidity. To check for presence of

1140   oocysts/sporozoites, mosquitoes were dissected, and their midguts or salivary

1141 glands were suspended in a drop of PBS solution atop a glass slide, covered by a

1142 coverslip and studied under a microscope.

1143

1144 For genetic crossing, isolates were harvested from donor mice and mixed to achieve

1145 a 1:1 ratio and 1 x $10^6$ parasites of this mixture was inoculated into four female CBA

1146 mice. Three days after inoculation, after confirming the presence of gametocytes,

1147 two infected CBA mice were anaesthetized and placed on two mosquito cages, each

1148 containing around 80 mosquitoes each. Mosquitoes were allowed to feed on the

1149 mice without interruption for 40 minutes at 24°C. A fresh feed was again performed

1150 on the 4th day post-inoculation with the other two CBA mice and two fresh cages of

1151 mosquitoes. 5-10 female mosquitoes from each cage were dissected on the 9th and

1152 12th day after the blood meal to check for presence of oocysts in the mosquito

1153 midguts. Twenty days after the blood meal, the mosquitoes were dissected and the

1154 salivary glands were removed, placed in 0.5-0.7 ml PBS solution and gently

1155 disrupted

1156 to release sporozoites. The suspensions from day 3 and day 4 feeds were injected

1157 intravenously into an ICR mouse each. When the mice were positive for blood-stage

1158 parasites, they were sub-inoculated into ten ICR mice with an inoculum of 0.6

1159 parasites/100uL to obtain clones from the potential cross progeny by limiting dilution.

1160 Eight days post infection, four mice were positive for parasites and these clones

1161 were screened for the presence of both *Pvs*EH and *Pvs*EL alleles within the

1162 chromosomes.

1163

## Declarations

1164

**Ethics approval and consent to participate**

1165

47

1166 Laboratory animal experimentation was performed in strict accordance with the

1167 Japanese Humane Treatment and Management of Animals Law (Law No. 105 dated

1168 19 October 1973 modified on 2 June 2006), and the Regulation on Animal

1169 Experimentation at Nagasaki University, Japan. The protocol was approved by the

1170 Institutional Animal Research Committee of Nagasaki University (permit:

1171 12072610052).

1172

1173 **Consent for publication**

1174 Not applicable

1175 **Availability of data and materials**

1176 All genome sequences, gene annotations and sequencing data files generated in

1177 this study can be found in ENA Study: PRJEB19355. All the datasets would also be

1178 available via EuPathDB portal at the time of publication. All parasite resources will be

1179 made available to the scientific community via the BEI Resources

1180 (https://www.beiresources.org/). Searchable and interactive versions of the

1181 phylogeny trees produced in this study can be accessed at

1182 https://itol.embl.de/shared/2lCr6w0mdDENs.

1183 **Competing interests**

1184 The authors declare that they have no competing interests.

1185

1190

**Authors' contributions**

RC and AP conceived the study. AR and RC conducted all rodent and mosquito experiments. AR and OG prepared sequencing libraries. AR, SK and RC conducted genetic cross experiments. AR collected data and performed all bioinformatic analyses. AR wrote the manuscript and all authors contributed to it.

**Authors' information (optional)**

Not applicable.

# References

1. Carlton JM, Hayton K, Cravo PV, Walliker D: **Of mice and malaria mutants: unravelling the genetics of drug resistance using rodent malaria models**. *Trends Parasitol* 2001, **17**(5):236-242.
2. Culleton RL, Abkallo HM: **Malaria parasite genetics: doing something useful**. *Parasitol Int* 2015, **64**(3):244-253.
3. Matz JM, Kooij TW: **Towards genome-wide experimental genetics in the in vivo malaria model parasite Plasmodium berghei**. *Pathog Glob Health* 2015, **109**(2):46-60.
4. De Niz M, Heussler VT: **Rodent malaria models: insights into human disease and parasite biology**. *Curr Opin Microbiol* 2018, **46**:93-101.
5. Carter R: **Studies on enzyme variation in the murine malaria parasites Plasmodium berghei, P. yoelii, P. vinckei and P. chabaudi by starch gel electrophoresis**. *Parasitology* 1978, **76**(3):241-267.

1221    6.    Perkins SL, Sarkar IN, Carter R: **The phylogeny of rodent malaria parasites:**
1222           **simultaneous analysis across three genomes**. *Infect Genet Evol* 2007, **7**(1):74-83.
1223    7.    Ramiro RS, Reece SE, Obbard DJ: **Molecular evolution and phylogenetics of rodent**
1224           **malaria parasites**. *BMC Evol Biol* 2012, **12**:219.
1225    8.    Cravo PV, Carlton JM, Hunt P, Bisoni L, Padua RA, Walliker D: **Genetics of mefloquine**
1226           **resistance in the rodent malaria parasite Plasmodium chabaudi**. *Antimicrob Agents*
1227           *Chemother* 2003, **47**(2):709-718.
1228    9.    Langhorne J, Quin SJ, Sanni LA: **Mouse models of blood-stage malaria infections:**
1229           **immune responses and cytokines involved in protection and pathology**. *Chem*
1230           *Immunol* 2002, **80**:204-228.
1231    10.   Spence PJ, Jarra W, Levy P, Reid AJ, Chappell L, Brugat T, Sanders M, Berriman M,
1232           Langhorne J: **Vector transmission regulates immune control of Plasmodium**
1233           **virulence**. *Nature* 2013, **498**(7453):228-231.
1234    11.   Brugat T, Reid AJ, Lin J, Cunningham D, Tumwine I, Kushinga G, McLaughlin S, Spence
1235           P, Bohme U, Sanders M *et al*: **Antibody-independent mechanisms regulate the**
1236           **establishment of chronic Plasmodium infection**. *Nat Microbiol* 2017, **2**:16276.
1237    12.   Stephens R, Culleton RL, Lamb TJ: **The contribution of Plasmodium chabaudi to our**
1238           **understanding of malaria**. *Trends Parasitol* 2012, **28**(2):73-82.
1239    13.   Prudencio M, Mota MM, Mendes AM: **A toolbox to study liver stage malaria**. *Trends*
1240           *Parasitol* 2011, **27**(12):565-574.
1241    14.   Guttery DS, Roques M, Holder AA, Tewari R: **Commit and Transmit: Molecular**
1242           **Players in Plasmodium Sexual Development and Zygote Differentiation**. *Trends*
1243           *Parasitol* 2015, **31**(12):676-685.
1244    15.   Pfander C, Anar B, Schwach F, Otto TD, Brochet M, Volkmann K, Quail MA, Pain A,
1245           Rosen B, Skarnes W *et al*: **A scalable pipeline for highly effective genetic**
1246           **modification of a malaria parasite**. *Nat Methods* 2011, **8**(12):1078-1082.
1247    16.   Marr E, Milne R, Anar B, Girling G, Schwach F, Mooney J, Nahrendorf W, Spence P,
1248           Cunningham D, Baker D *et al*: **An enhanced toolkit for the generation of knockout**
1249           **and marker-free fluorescent Plasmodium chabaudi [version 1; peer review: 2**
1250           **approved]**. *Wellcome Open Research* 2020, **5**(71).
1251    17.   Janse CJ, Ramesar J, Waters AP: **High-efficiency transfection and drug selection of**
1252           **genetically transformed blood stages of the rodent malaria parasite Plasmodium**
1253           **berghei**. *Nat Protoc* 2006, **1**(1):346-356.
1254    18.   Jongco AM, Ting LM, Thathy V, Mota MM, Kim K: **Improved transfection and new**
1255           **selectable markers for the rodent malaria parasite Plasmodium yoelii**. *Mol Biochem*
1256           *Parasitol* 2006, **146**(2):242-250.
1257    19.   Bushell E, Gomes AR, Sanderson T, Anar B, Girling G, Herd C, Metcalf T, Modrzynska
1258           K, Schwach F, Martin RE *et al*: **Functional Profiling of a Plasmodium Genome**
1259           **Reveals an Abundance of Essential Genes**. *Cell* 2017, **170**(2):260-272 e268.
1260    20.   Stanway RR, Bushell E, Chiappino-Pepe A, Roques M, Sanderson T, Franke-Fayard B,
1261           Caldelari R, Golomingi M, Nyonda M, Pandey V *et al*: **Genome-Scale Identification of**
1262           **Essential Metabolic Processes for Targeting the Plasmodium Liver Stage**. *Cell* 2019,
1263           **179**(5):1112-1128 e1126.
1264    21.   Antonova-Koch Y, Meister S, Abraham M, Luth MR, Ottilie S, Lukens AK, Sakata-Kato
1265           T, Vanaerschot M, Owen E, Jado JC *et al*: **Open-source discovery of chemical leads**
1266           **for next-generation chemoprotective antimalarials**. *Science* 2018, **362**(6419).

1267    22.     de Koning-Ward TF, Gilson PR, Crabb BS: **Advances in molecular genetic systems in**
1268         **malaria**. *Nat Rev Microbiol* 2015, **13**(6):373-387.
1269    23.   Carlton JM, Angiuoli SV, Suh BB, Kooij TW, Pertea M, Silva JC, Ermolaeva MD, Allen JE,
1270         Selengut JD, Koo HL *et al*: **Genome sequence and comparative analysis of the model**
1271         **rodent malaria parasite Plasmodium yoelii yoelii**. *Nature* 2002, **419**(6906):512-519.
1272    24.     Hall N, Karras M, Raine JD, Carlton JM, Kooij TW, Berriman M, Florens L, Janssen CS,
1273         Pain A, Christophides GK *et al*: **A comprehensive survey of the Plasmodium life cycle**
1274         **by genomic, transcriptomic, and proteomic analyses**. *Science* 2005, **307**(5706):82-
1275         86.
1276    25.     Fougere A, Jackson AP, Bechtsi DP, Braks JA, Annoura T, Fonager J, Spaccapelo R,
1277         Ramesar J, Chevalley-Maurel S, Klop O *et al*: **Variant Exported Blood-Stage Proteins**
1278         **Encoded by Plasmodium Multigene Families Are Expressed in Liver Stages Where**
1279         **They Are Exported into the Parasitophorous Vacuole**. *PLoS Pathog* 2016,
1280         **12**(11):e1005917.
1281    26.   Otto TD, Bohme U, Jackson AP, Hunt M, Franke-Fayard B, Hoeijmakers WA, Religa AA,
1282         Robertson L, Sanders M, Ogun SA *et al*: **A comprehensive evaluation of rodent**
1283         **malaria parasite genomes and gene expression**. *BMC Biol* 2014, **12**:86.
1284    27.     Bafort JM: **The biology of rodent malaria with particular reference to Plasmodium**
1285         **vinckei vinckei Rodhain 1952**. *Ann Soc Belges Med Trop Parasitol Mycol* 1971,
1286         **51**(1):5-203.
1287    28.     Carter R, Walliker D: **New observations on the malaria parasites of rodents of the**
1288         **Central African Republic - Plasmodium vinckei petteri subsp. nov. and Plasmodium**
1289         **chabaudi Landau, 1965**. *Ann Trop Med Parasitol* 1975, **69**(2):187-196.
1290    29.     Carter R, Walliker D: **Malaria parasites of rodents of the Congo (Brazzaville):**
1291         **Plasmodium chabaudi adami subsp. nov. and Plasmodium vinckei lentum Landau,**
1292         **Michel, Adam and Boulard, 1970**. *Ann Parasitol Hum Comp* 1976, **51**(6):637-646.
1293    30.     Killick-Kendrick R: **Parasitic Protozoa of the blood of rodents. V. Plasmodium**
1294         **vinckei brucechwatti subsp. nov. A malaria parasite of the thicket rat, Thamnomys**
1295         **rutilans, in Nigeria**. *Ann Parasitol Hum Comp* 1975, **50**(3):251-264.
1296    31.     Landau I, Michel JC, Adam JP, Boulard Y: **The life cycle of Plasmodium vinckei**
1297         **lentum subsp. nov. in the laboratory; comments on the nomenclature of the**
1298         **murine malaria parasites**. *Ann Trop Med Parasitol* 1970, **64**(3):315-323.
1299    32.     Bafort J: **New isolations of murine malaria in Africa: Cameroon**. In: *5th International*
1300         *Congress of Protozoology, New York City Abstracts of papers p343: 1977*.
1301    33.     Bafort J: **Etude du cycle biologique du Plasmodium v. vinckei Rodhain 1952**. *Ann*
1302         *Soc Belge Méd Trop* 1969, **49**(6):533-628.
1303    34.     Bafort JM, Molyneux DH: **Liver infections with Plasmodium v. vinckei in hosts**
1304         **refractory to blood infection**. *Trans R Soc Trop Med Hyg* 1971, **65**(1):13.
1305    35.     Lainson FA: **Observations on the morphology and electrophoretic variation of**
1306         **enzymes of the rodent malaria parasites of Cameroon, Plasmodium yoelii, P.**
1307         **chabaudi and P. vinckei**. *Parasitology* 1983, **86 (Pt 2)**:221-229.
1308    36.     Carter R: **Enzyme variation in Plasmodium berghei and Plasmodium vinckei**.
1309         *Parasitology* 1973, **66**(2):297-307.
1310    37.     LaCrue AN, Scheel M, Kennedy K, Kumar N, Kyle DE: **Effects of artesunate on**
1311         **parasite recrudescence and dormancy in the rodent malaria model Plasmodium**
1312         **vinckei**. *PLoS One* 2011, **6**(10):e26689.

1313    38.    Gautret P, Deharo E, Chabaud AG, Ginsburg H, Landau I: **Plasmodium vinckei vinckei,**
1314           **P. v. lentum and P. yoelii yoelii: chronobiology of the asexual cycle in the blood**.
1315           *Parasite* 1994, **1**(3):235-239.
1316    39.    Chandra R, Kumar S, Puri SK: **Plasmodium vinckei: infectivity of arteether-sensitive**
1317           **and arteether-resistant parasites in different strains of mice**. *Parasitol Res* 2011,
1318           **109**(4):1143-1149.
1319    40.    Killick-Kendrick R, Peters W: **Rodent malaria**. London: Academic Press; 1978.
1320    41.    Kooij TW, Carlton JM, Bidwell SL, Hall N, Ramesar J, Janse CJ, Waters AP: **A**
1321           **Plasmodium whole-genome synteny map: indels and synteny breakpoints as foci**
1322           **for species-specific genes**. *PLoS Pathog* 2005, **1**(4):e44.
1323    42.    Carlton J, Angiuoli S, Suh B, Kooij T, Pertea M, Silva J, Ermolaeva M, Allen J, Selengut
1324           J, Koo H *et al*: **Genome sequence and comparative analysis of the model rodent**
1325           **malaria parasite Plasmodium yoelii yoelii**. *Nature* 2002, **419**(6906):512-519.
1326    43.    Liu SL, Sanderson KE: **Rearrangements in the genome of the bacterium Salmonella**
1327           **typhi**. *Proc Natl Acad Sci U S A* 1995, **92**(4):1018-1022.
1328    44.    Abkallo HM, Martinelli A, Inoue M, Ramaprasad A, Xangsayarath P, Gitaka J, Tang J,
1329           Yahata K, Zoungrana A, Mitaka H *et al*: **Rapid identification of genes controlling**
1330           **virulence and immunity in malaria parasites**. *PLoS Pathog* 2017, **13**(7):e1006447.
1331    45.    Akinosoglou KA, Bushell ES, Ukegbu CV, Schlegelmilch T, Cho JS, Redmond S, Sala K,
1332           Christophides GK, Vlachou D: **Characterization of Plasmodium developmental**
1333           **transcriptomes in Anopheles gambiae midgut reveals novel regulators of malaria**
1334           **transmission**. *Cell Microbiol* 2015, **17**(2):254-268.
1335    46.    Ukegbu CV, Giorgalli M, Tapanelli S, Rona LDP, Jaye A, Wyer C, Angrisano F,
1336           Blagborough AM, Christophides GK, Vlachou D: **PIMMS43 is required for malaria**
1337           **parasite immune evasion and sporogonic development in the mosquito vector**.
1338           *Proc Natl Acad Sci U S A* 2020, **117**(13):7363-7373.
1339    47.    Sultan AA, Thathy V, Frevert U, Robson KJ, Crisanti A, Nussenzweig V, Nussenzweig
1340           RS, Menard R: **TRAP is necessary for gliding motility and infectivity of plasmodium**
1341           **sporozoites**. *Cell* 1997, **90**(3):511-522.
1342    48.    Sanders PR, Gilson PR, Cantin GT, Greenbaum DC, Nebl T, Carucci DJ, McConville MJ,
1343           Schofield L, Hodder AN, Yates JR, 3rd *et al*: **Distinct protein classes including novel**
1344           **merozoite surface antigens in Raft-like membranes of Plasmodium falciparum**. *J*
1345           *Biol Chem* 2005, **280**(48):40169-40176.
1346    49.    Taechalertpaisarn T, Crosnier C, Bartholdson SJ, Hodder AN, Thompson J,
1347           Bustamante LY, Wilson DW, Sanders PR, Wright GJ, Rayner JC *et al*: **Biochemical and**
1348           **functional analysis of two Plasmodium falciparum blood-stage 6-cys proteins: P12**
1349           **and P41**. *PLoS One* 2012, **7**(7):e41937.
1350    50.    Ramaprasad A, Subudhi AK, Culleton R, Pain A: **A fast and cost-effective**
1351           **microsampling protocol incorporating reduced animal usage for time-series**
1352           **transcriptomics in rodent malaria parasites**. *Malar J* 2019, **18**(1):26.
1353    51.    Favaloro JM, Kemp DJ: **Sequence diversity of the erythrocyte membrane antigen 1**
1354           **in various strains of Plasmodium chabaudi**. *Mol Biochem Parasitol* 1994, **66**(1):39-
1355           47.
1356    52.    Miller JL, Murray S, Vaughan AM, Harupa A, Sack B, Baldwin M, Crispe IN, Kappe SH:
1357           **Quantitative bioluminescent imaging of pre-erythrocytic malaria parasite infection**
1358           **using luciferase-expressing Plasmodium yoelii**. *PLoS One* 2013, **8**(4):e60820.

1359    53.    Franke-Fayard B, Trueman H, Ramesar J, Mendoza J, van der Keur M, van der Linden
1360           R, Sinden RE, Waters AP, Janse CJ: **A Plasmodium berghei reference line that**
1361           **constitutively expresses GFP at a high level throughout the complete life cycle**. *Mol*
1362           *Biochem Parasitol* 2004, **137**(1):23-33.
1363    54.    Boundenga L, Ngoubangoye B, Ntie S, Moukodoum ND, Renaud F, Rougeron V,
1364           Prugnolle F: **Rodent malaria in Gabon: Diversity and host range**. *Int J Parasitol*
1365           *Parasites Wildl* 2019, **10**:117-124.
1366    55.    Palmer GH, Brayton KA: **Gene conversion is a convergent strategy for pathogen**
1367           **antigenic variation**. *Trends Parasitol* 2007, **23**(9):408-413.
1368    56.    Subudhi AK, O'Donnell AJ, Ramaprasad A, Abkallo HM, Kaushik A, Ansari HR, Abdel-
1369           Haleem AM, Ben Rached F, Kaneko O, Culleton R *et al*: **Malaria parasites regulate**
1370           **intra-erythrocytic development duration via serpentine receptor 10 to coordinate**
1371           **with host rhythms**. *Nat Commun* 2020, **11**(1):2763.
1372    57.    Voza T, Gautret P, Renia L, Gantier JC, Lombard MN, Chabaud AG, Landau I:
1373           **Variation in murid Plasmodium desequestration and its modulation by stress and**
1374           **pentoxifylline**. *Parasitol Res* 2002, **88**(4):344-349.
1375    58.    Clark IA, Cowden WB, Butcher GA, Hunt NH: **Possible roles of tumor necrosis factor**
1376           **in the pathology of malaria**. *Am J Pathol* 1987, **129**(1):192-199.
1377    59.    Yoeli M, Hargreaves BJ: **Brain capillary blockage produced by a virulent strain of**
1378           **rodent malaria**. *Science* 1974, **184**(4136):572-573.
1379    60.    Tang J, Templeton TJ, Cao J, Culleton R: **The Consequences of Mixed-Species Malaria**
1380           **Parasite Co-Infections in Mice and Mosquitoes for Disease Severity, Parasite**
1381           **Fitness, and Transmission Success**. *Front Immunol* 2019, **10**:3072.
1382    61.    Abkallo HM, Tangena JA, Tang J, Kobayashi N, Inoue M, Zoungrana A, Colegrave N,
1383           Culleton R: **Within-host competition does not select for virulence in malaria**
1384           **parasites; studies with Plasmodium yoelii**. *PLoS Pathog* 2015, **11**(2):e1004628.
1385    62.    Martinelli A, Cheesman S, Hunt P, Culleton R, Raza A, Mackinnon M, Carter R: **A**
1386           **genetic approach to the de novo identification of targets of strain-specific**
1387           **immunity in malaria parasites**. *Proc Natl Acad Sci U S A* 2005, **102**(3):814-819.
1388    63.    Culleton R, Martinelli A, Hunt P, Carter R: **Linkage group selection: rapid gene**
1389           **discovery in malaria parasites**. *Genome Res* 2005, **15**(1):92-97.
1390    64.    Pattaradilokrat S, Culleton RL, Cheesman SJ, Carter R: **Gene encoding erythrocyte**
1391           **binding ligand linked to blood stage multiplication rate phenotype in Plasmodium**
1392           **yoelii yoelii**. *Proc Natl Acad Sci U S A* 2009, **106**(17):7161-7166.
1393    65.    Knowles G, Sanderson A, Walliker D: **Plasmodium yoelii: genetic analysis of crosses**
1394           **between two rodent malaria subspecies**. *Exp Parasitol* 1981, **52**(2):243-247.
1395    66.    Letunic I, Bork P: **Interactive Tree Of Life (iTOL) v4: recent updates and new**
1396           **developments**. *Nucleic Acids Res* 2019, **47**(W1):W256-W259.
1397    67.    Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C,
1398           O'Malley R, Figueroa-Balderas R, Morales-Cruz A *et al*: **Phased diploid genome**
1399           **assembly with single-molecule real-time sequencing**. *Nat Methods* 2016,
1400           **13**(12):1050-1054.
1401    68.    Otto TD, Sanders M, Berriman M, Newbold C: **Iterative Correction of Reference**
1402           **Nucleotides (iCORN) using second generation sequencing technology**.
1403           *Bioinformatics* 2010, **26**(14):1704-1707.
1404    69.    Zerbino DR, Birney E: **Velvet: algorithms for de novo short read assembly using de**
1405           **Bruijn graphs**. *Genome Res* 2008, **18**(5):821-829.

70. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL: **Versatile and open software for comparing large genomes**. *Genome Biol* 2004, **5**(2):R12.

71. Carver T, Harris SR, Berriman M, Parkhill J, McQuillan JA: **Artemis: an integrated platform for visualization and analysis of high-throughput sequence-based experimental data**. *Bioinformatics* 2012, **28**(4):464-469.

72. Carver TJ, Rutherford KM, Berriman M, Rajandream MA, Barrell BG, Parkhill J: **ACT: the Artemis Comparison Tool**. *Bioinformatics* 2005, **21**(16):3422-3423.

73. Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP: **Integrative genomics viewer**. *Nat Biotechnol* 2011, **29**(1):24-26.

74. Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA: **Circos: an information aesthetic for comparative genomics**. *Genome Res* 2009, **19**(9):1639-1645.

75. Stanke M, Waack S: **Gene prediction with a hidden Markov model and a new intron submodel**. *Bioinformatics* 2003, **19 Suppl 2**:ii215-225.

76. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL: **TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions**. *Genome Biol* 2013, **14**(4):R36.

77. Campbell MS, Holt C, Moore B, Yandell M: **Genome Annotation and Curation Using MAKER and MAKER-P**. *Curr Protoc Bioinformatics* 2014, **48**:4 11 11-14 11 39.

78. Lagesen K, Hallin P, Rodland EA, Staerfeldt HH, Rognes T, Ussery DW: **RNAmmer: consistent and rapid annotation of ribosomal RNA genes**. *Nucleic Acids Res* 2007, **35**(9):3100-3108.

79. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G *et al*: **InterProScan 5: genome-scale protein function classification**. *Bioinformatics* 2014, **30**(9):1236-1240.

80. Krogh A, Larsson B, von Heijne G, Sonnhammer EL: **Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes**. *J Mol Biol* 2001, **305**(3):567-580.

81. Petersen TN, Brunak S, von Heijne G, Nielsen H: **SignalP 4.0: discriminating signal peptides from transmembrane regions**. *Nat Methods* 2011, **8**(10):785-786.

82. Sargeant TJ, Marti M, Caler E, Carlton JM, Simpson K, Speed TP, Cowman AF: **Lineage-specific expansion of proteins exported to erythrocytes in malaria parasites**. *Genome Biol* 2006, **7**(2):R12.

83. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler transform**. *Bioinformatics* 2009, **25**(14):1754-1760.

84. Morgulis A, Gertz EM, Schaffer AA, Agarwala R: **A fast and symmetric DUST implementation to mask low-complexity DNA sequences**. *J Comput Biol* 2006, **13**(5):1028-1040.

85. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M *et al*: **The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data**. *Genome Res* 2010, **20**(9):1297-1303.

86. Zhang Z, Li J, Zhao XQ, Wang J, Wong GK, Yu J: **KaKs_Calculator: calculating Ka and Ks through model selection and model averaging**. *Genomics Proteomics Bioinformatics* 2006, **4**(4):259-263.

1452 87. Li L, Stoeckert CJ, Jr., Roos DS: **OrthoMCL: identification of ortholog groups for**
1453 **eukaryotic genomes**. *Genome Res* 2003, **13**(9):2178-2189.
1454 88. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced time and**
1455 **space complexity**. *BMC Bioinformatics* 2004, **5**:113.
1456 89. Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T: **trimAl: a tool for automated**
1457 **alignment trimming in large-scale phylogenetic analyses**. *Bioinformatics* 2009,
1458 **25**(15):1972-1973.
1459 90. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis**
1460 **of large phylogenies**. *Bioinformatics* 2014, **30**(9):1312-1313.
1461 91. Larsson A: **AliView: a fast and lightweight alignment viewer and editor for large**
1462 **datasets**. *Bioinformatics* 2014, **30**(22):3276-3278.

1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474 # Tables

| Nuclear genome features | *P. vinckei vinckei* CY | *P. vinckei brucechwatti* DA | *P. vinckei lentum* DE | *P. vinckei petteri* CR | *P. vinckei subsp.* EL |
|---|---|---|---|---|---|
| Genome Size (Mb) | 18.34 | 19.17 | 19.31 | 19.39 | 19.50 |
| G+C content (%) | 22.94 | 23.17 | 23.33 | 23.19 | 23.16 |
| Gaps within assembly | 0 | 0 | 7 | 0 | 11 |
| Genes | 5,042 | 5,238 | 5,249 | 5,310 | 5,307 |
| Pseudogenes | 19 | 40 | 29 | 36 | 36 |
| *ema1* | 7 | 20 | 19 | 15 | 21 |
| *etramp* | 12 | 12 | 11 | 13 | 12 |
| *fam-a* | 87 | 187 | 201 | 188 | 207 |
| *fam-b* | 28 | 38 | 41 | 40 | 44 |
| *fam-c* | 20 | 59 | 63 | 64 | 65 |
| *fam-d* | 6 | 13 | 27 | 21 | 21 |
| *hdh* | 7 | 9 | 12 | 13 | 14 |
| *lpl* | 25 | 23 | 24 | 23 | 17 |
| *p235* | 5 | 8 | 7 | 6 | 8 |
| pir | 178 | 209 | 210 | 272 | 247 |

1475

1476 **Table 1. Genome assembly characteristics of five *Plasmodium vinckei***

1477 **reference genomes.** AT-rich *P. vinckei* genomes are 19.2 to 19.5 megabasepairs

1478 (Mbps) long except for *Pvv*CY which has a smaller genome size of 18.3 Mb, similar

1479 to *Plasmodium berghei.* PacBio long reads allowed for chromosomes to be

1480 assembled as gapless unitigs with a few exceptions. Number of genes include partial

1481 genes and pseudogenes. Copy numbers of the ten multigene families differ between

1482 the *P. vinckei* subspecies (*ema1*, erythrocyte membrane antigen 1, *etramp*, early

1483 transcribed membrane protein, *hdh*, haloacid dehalogenase-like hydrolase, *lpl*,

1484 lysophospholipases, *p235*, reticulocyte binding protein, *pir*, *Plasmodium* interspersed

1485 repeat protein).

1486

1487

1488 # Figure legends

1489 **Figure 1. *Plasmodium vinckei* parasites and their phenotypic characteristics.**

1490 A) Rodent malaria parasite species and subspecies and the geographical sites in

1491 sub-Saharan Africa where from which they were isolated (modified from [1]).

1492 *Plasmodium vinckei* is the only RMP species to have been isolated from five different

1493 locations. Inset: To date, several RMP isolates have been sequenced (black) to aid

1494 research with rodent malaria models. Additional RMP isolates have been sequenced

1495 in this study (red) to cover all subspecies of *P. vinckei* and further subspecies of

1496 *Plasmodium chabaudi* and *Plasmodium yoelii*. B) Morphology of different life stages

1497 of *P. vinckei baforti* EL. R: Ring, ET: early trophozoite, LT: Late trophozoite, S:

1498 Schizont, MG: Male gametocyte, FG: Female gametocyte, O: oocyst and Sp:

1499 Sporozoite. *Plasmodium vinckei* trophozoites and gametocytes are morphologically

1500     distinct from other RMPs due to their rich haemozoin content (brown pigment). C)

1501     Parasitaemia of ten *P. vinckei* isolates (split into two graphs for clarity) during

1502     infections in mice (n=5) for a 20-day duration. † denotes host mortality. *Plasmodium*

1503     *vinckei* isolates show significant diversity in their virulence phenotypes.

1504

1505     **Figure 2. Structural variations and genotypic diversity among *Plasmodium***

1506     ***vinckei* parasites.** A) Chromosomal rearrangements in *P. vinckei* parasites.

1507     Pairwise synteny was assessed between the five *P. vinckei* subspecies and

1508     *Plasmodium berghei* (to represent the earliest common RMP ancestor). The 14

1509     chromosomes of different RMP genomes are arranged as a Circos plot and the

1510     ribbons (grey) between them denote regions of synteny. Three reciprocal

1511     translocation events (red) and one inversion (blue) accompany the separation of the

1512     different *P. vinckei* subspecies. A pan-*vinckei* reciprocal translocation between

1513     chromosomes VIII and X was observed between *P. vinckei* and other RMP

1514     genomes. Within the *P. vinckei* subspecies, two reciprocal translocations, between

1515     chromosomes V and XIII, and between chromosomes V and VI, separate

1516     *Plasmodium vinckei petteri* and *P. v. baforti* from the other three subspecies. A small

1517     inversion of ~100 kb region in chromosome 14 has occurred in *Pvv*CY alone. B)

1518     Maximum likelihood phylogeny of different RMP species with high-quality reference

1519     genomes based on protein alignment of 3,920 one-to-one orthologs (bootstrap

1520     values of each node are shown). Genomes of three human malaria species-

1521     *Plasmodium falciparum*, *Plasmodium vivax* and *Plasmodium knowlesi* were included

1522     in the analysis as outgroups. C) Maximum likelihood phylogenetic tree of all

1523     sequenced RMP isolates based on 1,010,956 high-quality SNPs (bootstrap values of

1524     each node are shown). There exists significant genotypic diversity among the *P.*

57

1525    *vinckei* isolates compared to the other RMPs. All *P. vinckei* subspecies have begun

1526    to diverge from their common ancestor well before sub-speciation events within

1527    *Plasmodium yoelii* and *Plasmodium chabaudi*. Genetic diversity within *P. v. petteri*

1528    and *P. v. baforti* isolates are similar to those observed within *P. yoelii* and *P.*

1529    *chabaudi* isolates while *P. v. lentum* and *P. v. brucechwatti* isolates have

1530    exceptionally high and low divergences respectively. Genes with significantly high

1531    Ka/Ks ratios in different subspecies-wise comparisons (as indicated by connector

1532    lines), the gene's Ka/Ks ratio averaged across all indicated *P. vinckei* comparisons

1533    and geographical origin of the isolates are shown.

1534

1535    **Figure 3. Sub-telomeric multigene family expansions in *Plasmodium vinckei***

1536    **parasites.** A)

1537    Violin plots show sub-telomeric multigene family size variations among RMPs and

1538    *Plasmodium falciparum*. The erythrocyte membrane antigen 1 and *fam-c* multigene

1539    families are expanded in the non-Katangan *P. vinckei* parasites (red). Apart from

1540    these families, multigene families have expanded in *P. vinckei* similar to that in

1541    *Plasmodium chabaudi*. The Katangan isolate *Pvv*CY (purple) has a smaller number

1542    of family members compared to non-Katangan isolates (orange) except for

1543    lysophospholipases, *p235* and *pir* gene families. B) Maximum Likelihood phylogeny

1544    of 99 *ema1* (top) and 328 *fam-c* (bottom) genes in RMPs. Branch nodes with good

1545    bootstrap support (> 70) are marked in red. The first coloured band denotes the

1546    RMP species to which the particular gene taxon belongs to. The heatmap denotes

1547    the relative gene expression among rings, trophozoite, schizont and gametocyte

1548    stages in the RMPs for which data are publicly available. Orange denotes high

1549    relative gene expression and white denotes low relative gene expression, while grey

1550    denotes lack of information. Gene expression was classified into three categories

1551    based on FPKM level distribution- High (black) denotes the top 25% of ranked FPKM

1552    of all expressed genes (FPKM > 256), Low (light grey) is the lower 25% of all

1553    expressed genes (FPKM < 32) and Medium level expression (grey;32<FPKM<256).

1554    "P" denotes pseudogenes. Four *vinckei*-group (*P. chabaudi* and *P. vinckei*

1555    subspecies) specific clades (Clades I-IV; orange), two *vinckei*-specific clade (Clade

1556    IV and V; purple) and one non-Katangan-specific clade (Clade VI; blue) can be

1557    identified within *ema1* family with strong gene expression, maximal during ring

1558    stages. Rest of the family's expansion within non-Katangan *P. vinckei* isolates are

1559    mainly pseudogenes with weak transcriptional evidence. The fam-c gene phylogeny

1560    shows the presence of four distinctly distal clades (A) with robust basal support (96-

1561    100). Of the four clades, two are pan-RMP (grey) and two are *vinckei* group-specific

1562    (orange), each consisting of fam-c genes positionally conserved across the member

1563    subspecies. Most members of these clades show medium to high gene expression

1564    during asexual blood stages. Other well-supported clades can be classified as either

1565    *berghei* group-specific (two; green), *vinckei* group-specific (two; orange), *P. vinckei* -

1566    specific (two; purple) or non-Katangan *P. vinckei* -specific clades (three; blue). There

1567    is evidence of significant species-specific expansion with striking examples in

1568    *Plasmodium yoelii* (i), *P. chabaudi* (ii) and in *P. v. brucechwatti* (iii).

1569

1570    **Figure 4. Phenotypic variation and genetics in *Plasmodium vinckei* parasites.**

1571    A) Schematic of isolate-specific genetic markers detected in clonal line of *Pvs*EL X

1572    *Pvs*EH cross progeny by Sanger sequencing. Genetic markers from both EH (red)

1573    and EL (blue) isolates were detected in the crossed progeny proving successful

1574    genetic crossing. B) Schematic of homologous recombination-mediated insertion of a

1575    gfp-luciferase cassette into the p230p locus in *P. vinckei* CY. C) GFP expression in

1576   different blood stages of *Pvv*CY and luciferase expression of *Pvv*CY oocysts in

1577   mosquito midgut.

1578

## Additional files

1579   **Additional file 1. Summary of rodent malaria parasite isolates used and DNA**

1581   **and RNA sequencing performed in this study.**

1582   **Additional file 2. Infection profiles of ten *Plasmodium vinckei* isolates.** Changes

1583   in parasitaemia, host RBC density and host weight during *P. vinckei* infections. Error

1584   bars show standard deviation of the readings within five biological replicates. †

1585   denotes host mortality.

1586   **Additional file 3. Daily readings of parasitaemia, host RBC density and host**

1587   **weight during infections of *Plasmodium vinckei* isolates.**

1588   **Additional file 4. Assembly statistics of *Plasmodium vinckei* mitochondrial and**

1589   **apicoplast genomes.**

1590   **Additional file 5. Chromosomal synteny breakpoints among *Plasmodium***

1591   ***vinckei* genomes.**

1592   **Additional file 6. 3,920 one-to-one orthologous group used for genome-wide**

1593   **protein alignment-based phylogeny.**

1594   **Additional file 7.  Single nucleotide polymorphisms among RMP isolates and**

1595   **Ka/Ks ratios for various pair-wise comparisons of homologous protein-coding**

1596   **genes.**

1597   **Additional file 8. Copy number variations within multigene families and**

1598   **phylogenetic clade members.**

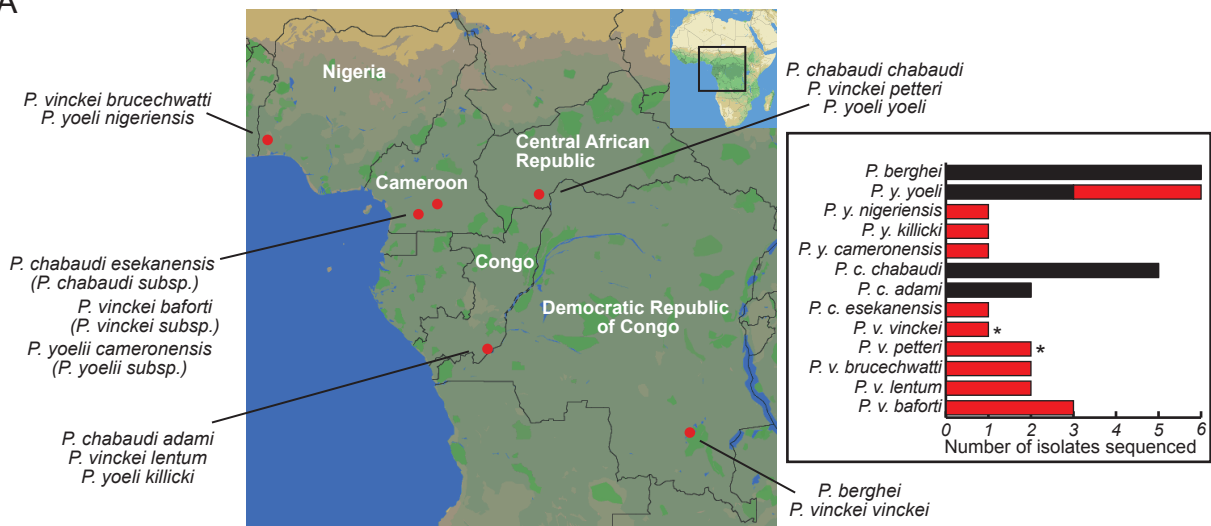1599   **Additional file 9. Maximum Likelihood trees for ten RMP multigene families.**

1600    **Additional file 10. Gene-wise RNA-seq FPKM values for *Plasmodium vinckei***

1601    ***petteri* CR, *P. v. lentum* DS (Rings, trophozoites and gametocyte stages), *P.v.***

1602    ***brucechwatti* DA and *P. v. baforti* EL (mixed blood stages).**

1603    **Additional file 11. Mosquito transmission and genetic cross experiments.**
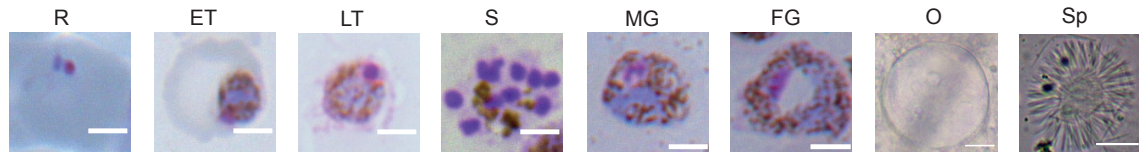
1604    **Additional file 12. A) Circos figure showing rearrangements among four RMP**

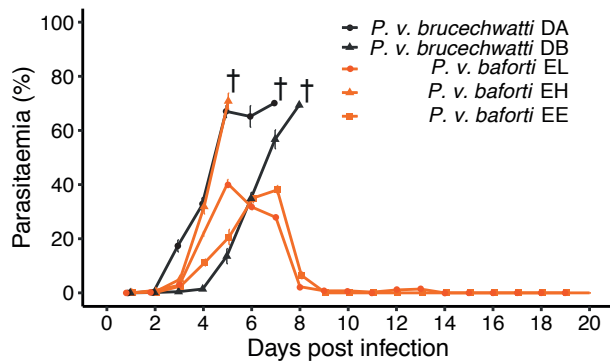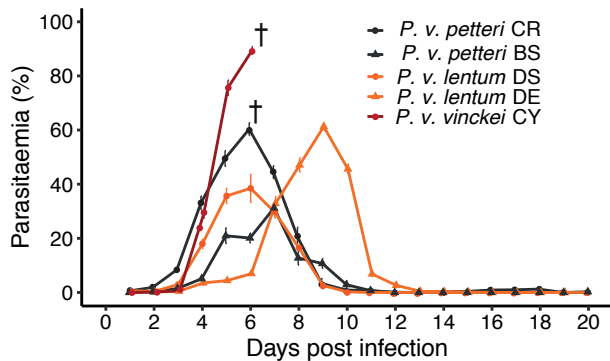1605    **species B) Gene alignment of pseudogenised ema1 genes.**
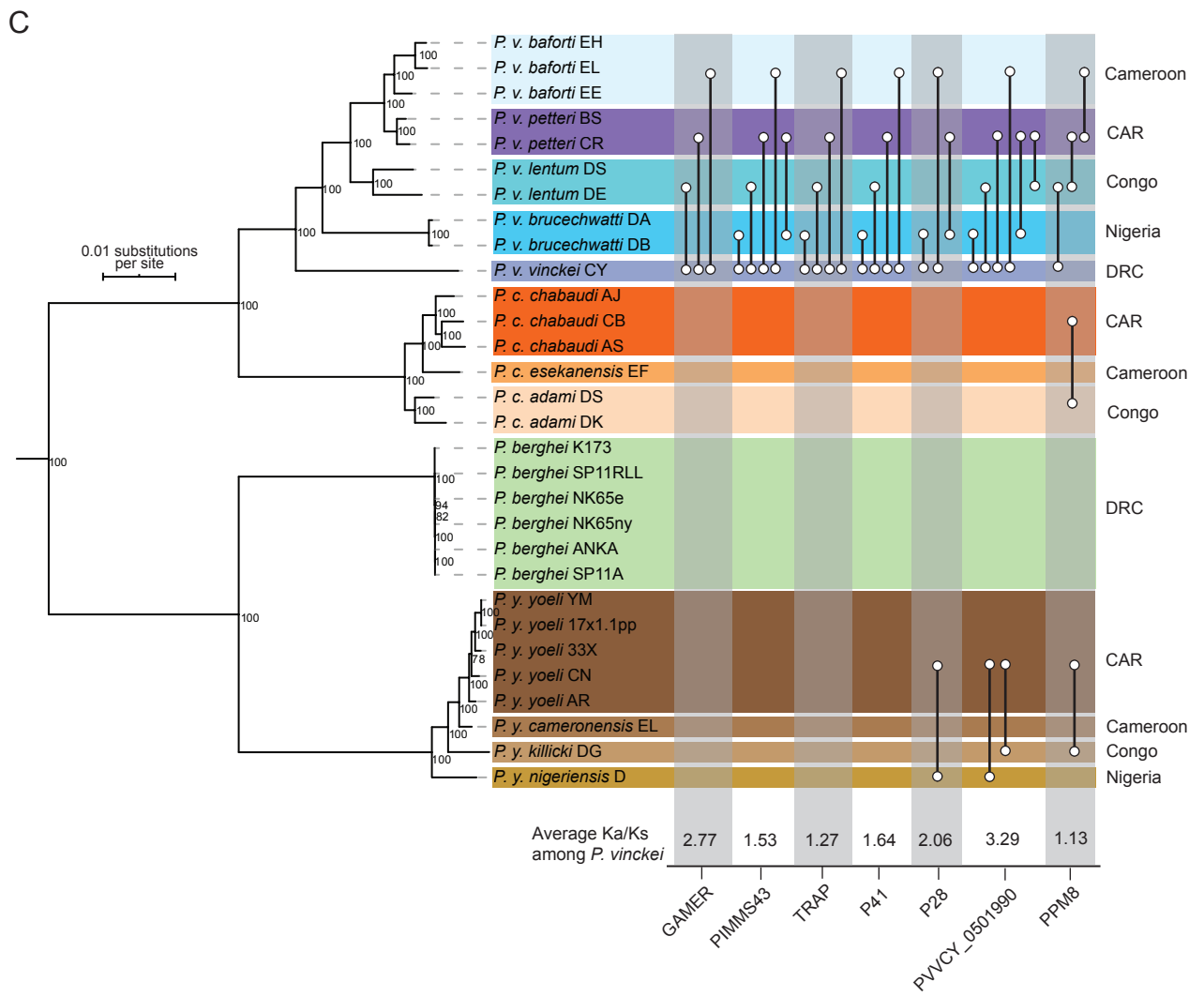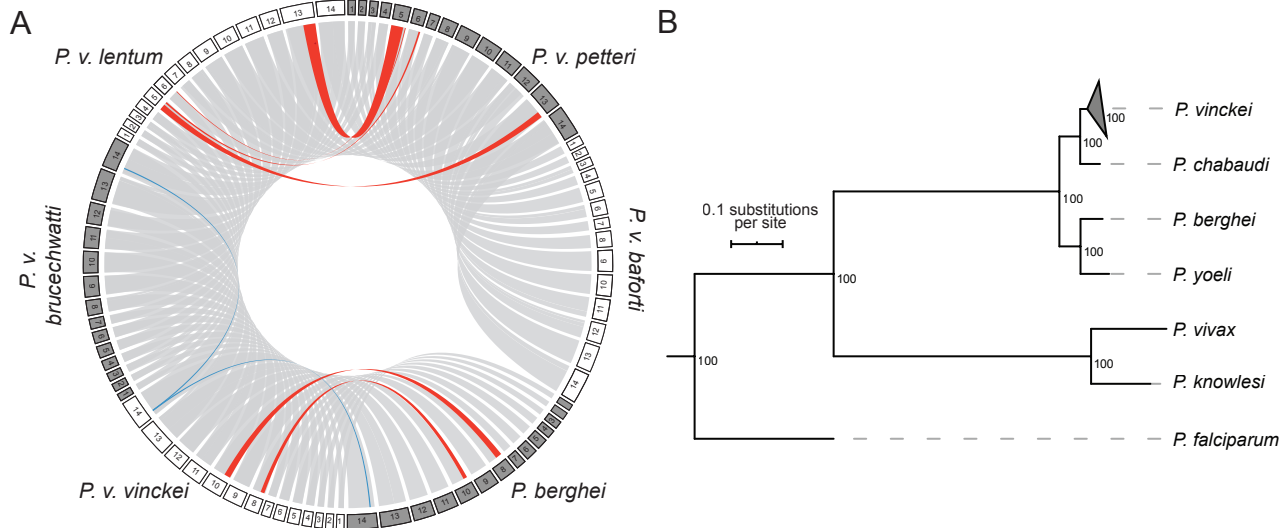
1606
1607

**A**

Nigeria

*P. vinckei brucechwatti*
*P. yoeli nigeriensis*

Cameroon

Central African Republic

*P. chabaudi chabaudi*
*P. vinckei petteri*
*P. yoeli yoeli*

*P. chabaudi esekanensis*
(*P. chabaudi subsp.*)

*P. vinckei baforti*
(*P. vinckei subsp.*)

*P. yoelii cameronensis*
(*P. yoelii subsp.*)

Congo

Democratic Republic
of Congo

*P. chabaudi adami*
*P. vinckei lentum*
*P. yoeli killicki*

*P. berghei*
*P. vinckei vinckei*

| Species | Number of isolates sequenced |
|---|---|
| *P. berghei* | |
| *P. y. yoeli* | |
| *P. y. nigeriensis* | |
| *P. y. killicki* | |
| *P. y. cameronensis* | |
| *P. c. chabaudi* | |
| *P. c. adami* | |
| *P. c. esekanensis* | |
| *P. v. vinckei* | * |
| *P. v. petteri* | * |
| *P. v. brucechwatti* | |
| *P. v. lentum* | |
| *P. v. baforti* | |

**B**

R   ET   LT   S   MG   FG   O   Sp

**C**

Days post infection — Parasitaemia (%)

- *P. v. petteri* CR
- *P. v. petteri* BS
- *P. v. lentum* DS
- *P. v. lentum* DE
- *P. v. vinckei* CY

- *P. v. brucechwatti* DA
- *P. v. brucechwatti* DB
- *P. v. baforti* EL
- *P. v. baforti* EH
- *P. v. baforti* EE
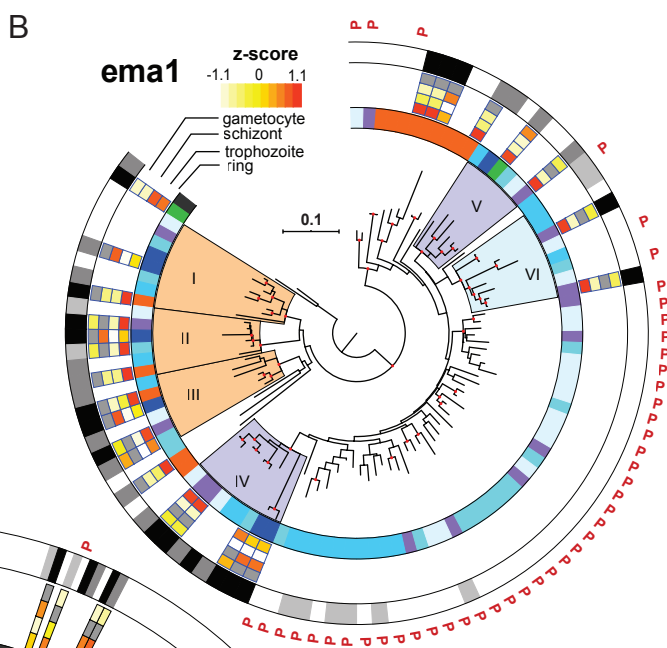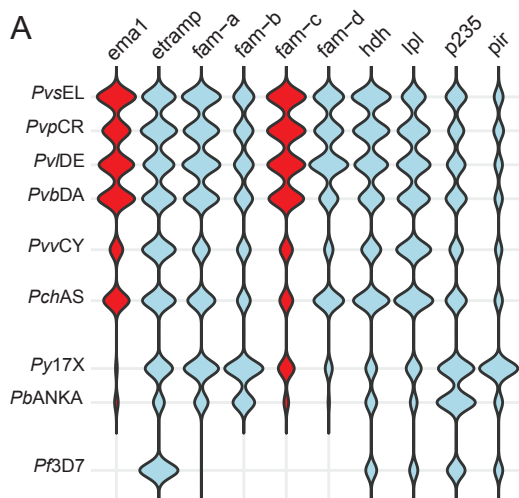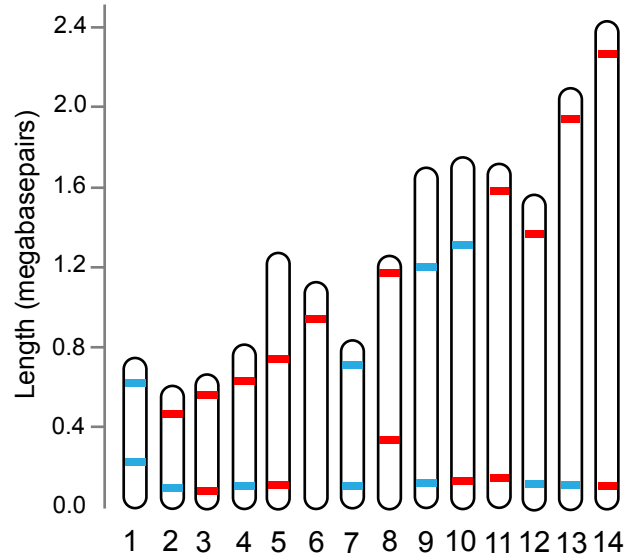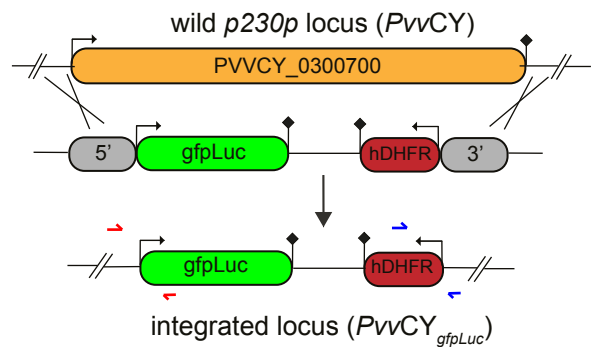
A

Gene derived from *P. vinckei baforti* EL
Gene derived from *P. vinckei baforti* EH

Length (megabasepairs)

B

wild *p230p* locus (*PvvCY*)

PVVCY_0300700

5'  gfpLuc  hDHFR  3'

integrated locus (*PvvCY_gfpLuc*)

gfpLuc  hDHFR

C

Bright field | DAPI | GFP-LUC | Merge

trophozoite

schizont

female gametocyte

oocyst