

A paradoxical population structure of *var* DBL α types in Africa

Mun Hua Tan^{1†}, Kathryn E. Tiedje^{1†}, Qian Feng², Qi Zhan³, Mercedes Pascual³, Heejung Shim², Yao-ban Chan², and Karen P. Day^{1*}

¹ Department of Microbiology and Immunology, The University of Melbourne, Bio21 Institute and Peter Doherty Institute, Melbourne, AU

² School of Mathematics and Statistics / Melbourne Integrative Genomics, The University of Melbourne, Melbourne, Australia

³ Department of Ecology and Evolution, University of Chicago; Chicago, Illinois, USA

[†]Co-first authors

ABSTRACT

The *var* multigene family encodes the *P. falciparum* erythrocyte membrane protein 1 (PfEMP1), which is important in host-parasite interaction as a virulence factor and major surface antigen of the blood stages of the parasite, responsible for maintaining chronic infection. Whilst important in the biology of *P. falciparum*, these genes (50 to 60 genes per parasite genome) are routinely excluded from whole genome analyses due to their hyperdiversity, achieved primarily through recombination. The PfEMP1 head structure almost always consists of a DBL α -CIDR tandem. Categorized into different groups (upsA, upsB, upsC), different head structures have been associated with different ligand-binding affinities and disease severities. We study how conserved individual DBL α types are at the country, regional, and local scales in Sub-Saharan Africa. Using publicly-available sequence datasets and a novel ups classification algorithm, *cUps*, we performed an *in silico* exploration of DBL α conservation through time and space in Africa. In all three ups groups, the population structure of DBL α types in Africa consists of variants occurring at rare, low, moderate, and high frequencies. Non-rare variants were found to be temporally stable in a local area in endemic Ghana. When inspected across different geographical scales, we report different

30 levels of conservation; while some DBL α types were consistently found in high frequencies in
31 multiple African countries, others were conserved only locally, signifying local preservation of
32 specific types. Underlying this population pattern is the composition of DBL α types within
33 each isolate DBL α repertoire, revealed to also consist of a mix of types found at rare, low,
34 moderate, and high frequencies in the population. We further discuss the adaptive forces and
35 balancing selection, including host genetic factors, potentially shaping the evolution and
36 diversity of DBL α types in Africa.

37

38 **1. INTRODUCTION**

39 Malaria parasites in endemic areas with high transmission undergo frequent outcrossing in
40 the vector (1,2). The diversification of parasites through meiotic, mitotic, and ectopic
41 recombination results in very high levels of genetic diversity in the *Plasmodium falciparum*
42 parasite population (3–7), particularly of the major variant surface antigens such as the *var*
43 multigene family that encode the *Plasmodium falciparum* erythrocyte membrane protein 1
44 (PfEMP1). PfEMP1 proteins are expressed on the surface of infected erythrocytes and can
45 bind to host receptors to mediate cytoadhesion and sequestration of infected cells (8–10).
46 Through clonal antigenic variation of *var* genes, whereby different genes are expressed
47 sequentially and exclusively during the blood stage, parasites are also able to effectively
48 evade immune detection, promoting chronic infection within a host (11,12). This high
49 diversity has served as the basis for *var* surveillance (13,14), population genetics (15–17), and
50 estimation of infection complexity (15). However, less work has been done to characterise
51 this diversity and the population structure of antigenic factors of and within these genes.

52

53 There are approximately 40 to 60 different *var* genes found across all 14 chromosomes of a
54 *P. falciparum* genome (18). Based on their upstream promoter sequences, *var* genes can be
55 divided into groups of A, B, C, and E, with a minority of genes grouped in two intermediate
56 groups of B/A or B/C (19). The three major ‘ups’ groups of upsA, upsB, and upsC are associated
57 with different chromosomal locations, transcriptional directions, and sequences (18–21).
58 Genes in the upsA and upsB groups are generally located at the subtelomeric regions whereas
59 genes in the upsC group are found in the central regions of specific chromosomes. UpsA genes
60 have been shown to be transcribed towards centromeres and conversely, upsB and upsC
61 genes are often found to be transcribed towards telomeres.

62

63 The extracellular N-terminal PfEMP1 head structure almost always consists of a Duffy-
64 binding-like alpha domain (DBL α) and a cysteine-rich interdomain region (CIDR) (i.e., a DBL α -
65 CIDR tandem). This head structure exists in different configurations of these DBL α and CIDR
66 domain subclasses (e.g., DBL α 0-2 with CIDR α 1-6, CIDR β 1-7, CIDR δ 1-12, CIDR γ 1,2) and can
67 influence ligand binding and disease pathogenicity. The prevailing understanding is that *var*
68 genes and DBL α variants (i.e., DBL α types) in the upsA group are generally more conserved
69 compared to those in upsB or upsC groups (i.e., non-upsA) (15,16,19,22). This has largely been
70 attributed to the association of upsA genes to severe malaria, including cerebral malaria, due
71 to their host receptor binding phenotypes (23). Expression of upsA *var* genes encoding the
72 DBL α +CIDR α 1 head structure mediate endothelial protein C receptor (EPCR)-binding and/or
73 intercellular adhesion molecule-1 (ICAM-1)-binding has been associated with severe malaria
74 and/or cerebral malaria (24–26). In addition, PfEMP1 in the upsA group containing the
75 DBL α +CIDR β / δ / γ head structure has been associated with rosetting with uninfected
76 erythrocytes (24). On the other hand, expression of upsB and upsC *var* genes (e.g.,

77 DBL α +CIDR α 2-6) have been associated with uncomplicated malaria, commonly mediated by
78 adhesion to the cluster differentiation 36 (CD36) receptor (27,28), and may be more active in
79 establishing chronic infections (29).

80

81 With the exception of one specific *var* gene involved in pregnancy-associated malaria (i.e.,
82 *var2csa*), all other *var* genes encode for a DBL α domain, which is one of the most diverse
83 domains and has been shown to be immunogenic to variant-specific epitopes, recognised
84 serologically in an age-dependent manner (30). Multiple studies have noted that there exists
85 a minority of DBL α types that are highly conserved over various spatial scales (16,31–35).
86 These studies typically focussed on DBL α types or *var* genes found within the highest
87 percentiles (e.g., (35)), at very high frequencies (e.g., (16,31)), or those found globally
88 conserved and prevalent (e.g., (34,35)). Understandably, looking for the most common DBL α
89 or *var* is instinctive in the search for the elusive vaccine candidate targeting the most
90 important group of variant surface antigens of the parasite. However, *var* gene sharing among
91 isolates, particularly those living in high transmission, has been shown to be minimal (33).
92 Given that multiple variants of *var* genes with the same binding phenotype exist, this pre-
93 occupation over the most globally-common types or genes risks overlooking the genetic
94 patterns underlying a local population.

95

96 In the same way that ‘severe malaria’ must be defined by different malaria epidemiologies,
97 patterns of conservation must also be interpreted within the context of an area’s local
98 epidemiology. In low-transmission areas in South America and Asia, conservation found
99 across countries and continents likely relates to small population sizes due to founder effects,
100 in which *var* genes have not yet diversified (36–38). In moderate transmission, profiles may

101 exhibit bias toward more moderate frequency classes, in conjunction with greater diversity
102 within the area. Conservation in high-transmission areas would be most interesting, as these
103 areas possess the epidemiological and genetic characteristics to generate vast diversity. In
104 such a system of high prevalence/incidence, high outcrossing rates, and high genetic diversity
105 of parasites, finding conservation will provide insights into factors constraining and shaping
106 diversity in a highly dynamic system.

107

108 Equipped with publicly-available DBL α data from several populations in Africa and a novel ups
109 classification algorithm (**cUps**) introduced in this paper, this study sought to identify DBL α
110 type conservation beyond those reported in the highest frequencies within and among
111 populations globally. By categorising DBL α types into broad frequency classes of rare to
112 highly-frequent types, we identified two kinds of conservation within all ups groups in our
113 study area in Bongo, Ghana: (1) Conservation of specific types across isolates in a time point
114 and through time (i.e., survey), (2) Conservation of type frequencies (i.e., types were found
115 at relatively stable frequencies through time). We show that these patterns are maintained
116 through the composition of DBL α types within each isolate repertoire, revealed to consist of
117 a mix of types found at rare, low, moderate, and high frequencies in the population. In a
118 spatial analysis, in addition to identifying DBL α types conserved at the continent level, we
119 noted that there are DBL α types conserved at the local and/or regional levels but not
120 necessarily prevalent across wider geographical scales, prompting a discussion on the
121 adaptive forces potentially driving balancing selection and shaping the population structure
122 of DBL α types in Africa.

123

124 **2. RESULTS**

125 **2.1 Description of time-series cross-sectional surveys in Bongo, Ghana**

126 This study analysed publicly-available DBL α tag sequence data from an interrupted time-
127 series study design (i.e., Malaria Reservoir Study (MRS)) (Figure I in Data S1) (13,15,39). This
128 MRS dataset consists of sampling at seven time points from 2012 to 2017 at the end of wet
129 or dry seasons. Each time point represented an age-stratified cross-sectional survey of
130 approximately 2,000 asymptomatic participants per survey (ages from 1 to 97 years old) from
131 two proximal catchment areas (Vea/Gowrie and Soe, with a sampling area ~ 60 km²) in Bongo
132 District in Northern Ghana. Surveyed participants (i.e., isolates) represented approximately
133 15% of the total population that reside in the two catchment areas in Bongo District at a time
134 (Table I in Data S1). This area is characterised by high, seasonal malaria transmission and has
135 undergone several types of interventions, including long-lasting insecticide-treated nets
136 (LLINs) and indoor residual spraying (IRS) that reduced transmission (13,15), and seasonal
137 malaria chemoprevention (SMC) that reduced the burden of infection in children younger
138 than 5 years old (15). Clustering of DBL α tag sequences from seven surveys (S1 to S7) and
139 further post-processing of the dataset (see Methods) resulted in 62,168 representative DBL α
140 types found in 3,166 isolates for this study of DBL α conservation in Bongo (Table I in Data S1).

141

142 In a high-transmission setting, the asymptomatic “population” typically consists of “isolates”
143 infected by one or more unique parasite “genomes”. This complexity of infections is indicated
144 by multiplicity of infection (MOI), where an isolate with MOI = 1 would represent a single
145 unique parasite genome. Hence, at MOI = 1, an isolate’s DBL α repertoire is synonymous to a
146 parasite’s DBL α repertoire (i.e., the DBL α repertoire in a single parasite genome). Conversely,
147 at MOI > 1, an isolate’s DBL α repertoire would encompass > 1 parasites’ DBL α repertoire.

148

149 **2.2 A novel ups classification algorithm based on *var* DBL α tags**

150 This study introduces ***cUps***, a novel algorithm for classifying DBL α types into the different
151 groups of upsA, upsB, and upsC (Data S2, Figure 1A). At the isolate level, the average isolate
152 repertoire consists of 20.9%, 48.6%, and 30.5% of upsA, upsB, and upsC DBL α types,
153 respectively (Figure II in Data S1). These proportions differ from those reported in (22) that
154 estimated higher proportions of upsB and lower proportions of upsC in isolate repertoires,
155 based on the average of seven genomes. The algorithm shows a tendency to classify more
156 upsB types as upsC types. This is in line with validation results on the algorithm's specificity
157 and sensitivity presented Data S2. A reduced analysis that involved DBL α types with higher
158 confidence in classification (threshold = 8, see Data S2) yielded similar trends and patterns of
159 observation. Genetic similarity by pairwise type sharing (PTS) remains extremely low for all
160 ups groups (median PTS: 4.55% (upsA), 1.00% (upsB), and 2.15% (upsC)) (Figure 1B). The
161 62,168 representative DBL α types from the seven combined MRS surveys were classified into
162 upsA (5.4%), upsB (56.6%), and upsC (37.9%) groups (Table I in Data S1). This points to a
163 highest DBL α richness for the upsB group (35,215 types), followed by the upsC group (23,583
164 types) and upsA group (3,370 types) in combined surveys, and this hierarchy of richness is
165 similarly observed for individual surveys (Figure 1A). The differences in proportions of ups
166 groups at the isolate versus population levels is attributed to the negative relationship
167 between PTS and richness, as a higher level of upsA type sharing will result in lower proportion
168 of unique representative DBL α types in the population.

169

170 **2.3 *Var* DBL α types are conserved in the local Bongo population and through time**

171 The frequency of a DBL α type is defined as the proportion of the population a DBL α type is
172 found in and is calculated in the context of individual surveys (i.e., "survey-specific

173 frequency”), or of the averaged frequencies of the seven surveys (i.e., “survey-averaged
174 frequency”, see Methods for details on frequency calculation and normalisation). To explore
175 the conservation of DBL α types, we further categorised these frequencies into four classes of
176 **low** (0%, 1%), **moderate** [1%, 5%), **high** [5%, 10%), and **very high** [10%, 100%] frequencies.
177 While it is well reported that DBL α types in the upsA group are generally more conserved
178 relative to types in the upsB and upsC groups, this study identified conservation of DBL α types
179 in all three ups groups. We describe patterns of DBL α conservation observed in Bongo at the
180 population level in the following subsections.

181

182 **2.3.1. Conservation of var DBL α types in a surveyed time point**

183 DBL α types are considered conserved in a surveyed time point when found in multiple isolates
184 sampled in a same survey (i.e., moderate-to-high survey-specific frequencies). Overall,
185 distributions of survey-averaged and survey-specific frequencies were found to be strongly
186 and positively skewed, with most DBL α types occurring at low frequencies in <1% of isolates.
187 For all three ups groups, when categorised into frequency classes, the second largest subsets
188 of DBL α types are shown to occur at moderate frequencies between 1% to 5% in each survey,
189 followed by smaller subsets of DBL α types found at higher frequencies exceeding 5% and/or
190 10% frequencies (Figure 2A). The most frequent DBL α types in the upsA, upsB, and upsC
191 groups were detected at survey-specific frequencies of 61.1%, 42.9%, and 62.0%,
192 respectively. In the different surveys, this study identified hundreds to thousands of
193 moderate-to-highly conserved DBL α types (i.e., $\geq 1\%$ survey-specific frequencies) in all three
194 ups groups (upsA: 525 to 790 types per survey, upsB: 539 to 1,917 types per survey, upsC:
195 365 to 1,121 types per survey). This translates into different proportions of DBL α types in each
196 ups group, owing to the higher DBL α type richness of upsB and upsC groups (upsA: 29.2% to

197 43.9% per survey, upsB: 5.2% to 15.5% per survey, upsC: 5.6% to 14.3% per survey) (Table II
198 in Data S1).

199

200 **2.3.2. Conservation of *var* DBL α type frequencies across multiple surveys and through time**

201 Frequencies of DBL α types are considered conserved if a same DBL α type occurs at stable or
202 similar frequencies in multiple surveys. In all ups groups, a strong positive correlation
203 between survey-specific frequencies and survey-averaged frequencies of DBL α types is
204 shown, indicating that DBL α types occurring at high survey-averaged frequencies were also
205 generally found at high frequencies in specific surveys and likely in multiple surveys (Figure
206 2C, Figure III in Data S1). Likewise, most DBL α types found at moderate frequencies also
207 showed consistent maintenance of frequencies across multiple surveys. Furthermore, DBL α
208 types are considered conserved through time when found in multiple surveys. This study
209 shows a correlation between DBL α type frequencies (per-survey and survey-averaged) and
210 the number of surveys the types are found in, with most of the relatively conserved DBL α
211 types in all three ups groups also persisting through time (Figure 2B, Figure III and IV in Data
212 S1). All DBL α types occurring at high or very high survey-averaged and/or survey-specific
213 frequencies were seen in all seven surveys. The majority of those with moderate survey-
214 averaged or survey-specific frequencies were found in all seven surveys, with a smaller subset
215 found in four to six surveys. On the other hand, the remainder of DBL α types found at low
216 frequencies (< 1%) can be found in the range of one to seven surveys. This was observed in
217 all three ups groups.

218

219 **2.4 Isolate repertoires consist of a mix of conserved and rare *var* DBL α types**

220 While isolate DBL α repertoires in high-transmission populations have been reported to be
221 unrelated and largely non-overlapping (15,17), there has not been a detailed exploration of
222 the composition of DBL α types and their respective frequency classes at the isolate repertoire
223 level (i.e., ‘per-isolate frequency profiles’). This study showed that all isolate repertoires
224 consist of DBL α types from all three ups groups. More interestingly, these per-isolate
225 frequency profiles show consistent proportions of rare, moderately-frequent, and highly-
226 frequent DBL α types in every isolate, based on survey-specific frequencies (Figure 3, Figure V
227 in Data S1). In all ups groups, these frequency profiles were consistent across isolates within
228 the same survey, regardless of isolates’ infection complexities (i.e., MOI). Importantly, the
229 observation of these frequency patterns in MOI = 1 isolates (i.e., isolate non-upsA repertoire
230 size of approximately ≤ 45 in Figure 3) indicates that these per-isolate frequency profiles are
231 a consequence of similar repertoire composition within actual parasite genomes.

232

233 Of the three groups, isolates’ upsA frequency profiles exhibit largest proportions of DBL α
234 types found at moderate survey-specific frequencies and smaller proportions of DBL α types
235 in the lowest survey-specific frequency class. In contrast, isolates’ upsB and upsC frequency
236 profiles both consist of largest proportions of DBL α types in the lowest frequency class (i.e.,
237 (0%, 1%)), followed by those in the moderate frequency class. The introduction of
238 interventions did not perturb these frequency profiles, which maintained the composition of
239 different frequency classes, albeit in different proportions. In surveys of the population
240 affected by interventions (e.g., S4 and S5 during and after indoor residual spraying), frequency
241 profiles trended toward a generally larger proportion of low-frequency DBL α types and
242 smaller proportions of higher frequency DBL α types within each isolate (Figure 3). We
243 confirmed this observation of per-isolate frequency profiles using an independent DBL α

244 sequence dataset (14,33,40), extracted from *var* genes of isolates sampled from Navrongo in
245 Ghana, situated ~30 km adjacent to Bongo district (Figure VI in Data S1).

246

247 As expected, genetic similarity between pairwise isolate repertoires (i.e., PTS values)
248 increases when DBL α types of lower frequencies are excluded (Figure VII in Data S1). PTS
249 values range from 0 to 1 representing unrelated to identical isolate repertoires, respectively.
250 Interestingly, even when considering only DBL α types found at very high frequencies, median
251 PTS values remained generally low (median PTS of 0.02, 0.05, 0.13, and 0.19 when considering
252 DBL α types at >0%, \geq 1%, \geq 5%, \geq 10% survey-averaged frequencies, respectively, across all
253 surveys). This indicates that, while every isolate repertoire contains sets of types that are
254 conserved in the population, identical sets of conserved DBL α types are rare. Shifts in PTS
255 distributions were more substantial when exclusively evaluating DBL α types in the upsA or
256 upsC groups relative to the upsB group, consistent with the lower DBL α richness in the two
257 former groups (i.e., the less variants there are in the population, the higher the probability of
258 overlaps).

259

260 **2.5 Global and local preservation of *var* DBL α types in Africa**

261 Further, a separate spatial study of DBL α conservation in multiple African countries (i.e.,
262 “locations”) representing West Africa (Ghana, Gabon), Central Africa (Malawi) and East Africa
263 (Uganda) was conducted based on 82,027 DBL α types found in 4,783 isolates (Table I in Data
264 S3, Figure I in Data S3). Similarly, this spatial study showed that the majority of DBL α types
265 were found at low frequencies with smaller proportions seen at higher frequencies.
266 Comparison of DBL α types and frequencies in the four locations showed conservation of the
267 same upsA DBL α types at moderate-to-high frequencies in all locations (i.e., a highly-frequent

268 DBL α type in Ghana was also found at moderate-to-high frequencies in other analysed
269 locations) (Figure 4). While this was also observed for some highly-frequent DBL α types in the
270 upsB and upsC groups (i.e., non-upsA groups), this study additionally identified some highly-
271 frequent DBL α types in these groups that were present predominantly in a single location,
272 suggesting local selection and preservation of DBL α types in these ups groups (Figure 4, Figure
273 II and III in Data S3). Additionally, as was also observed for isolates in the MRS study, per-
274 isolate frequency profiles in these different locations also consisted of a mix of rare to
275 conserved DBL α types (Figure IV in Data S3). It is worth reminding that these conserved DBL α
276 types make up the minority of all DBL α types in every ups group. An exploration of the
277 relationship between DBL α types and *var* exon 1 sequences revealed that these conserved
278 DBL α types are associated with multiple different *var* exon 1 sequences (Figure V in Data S3),
279 indicating that other parts of the gene were still diversifying even though the DBL α types were
280 maintained in the population. For some of these DBL α types with 1-to-many DBL α -*var*
281 relationships, pairwise nucleotide identity between *var* exon 1 sharing the same DBL α type
282 suggest that some of these *var* sequences could be alleles of a same gene (14). However, the
283 majority of these *var* exon 1 sequences exhibit low shared identity and therefore appear to
284 represent actual different genes (Figure V in Data S3). In a highly-dynamic system where the
285 DBL α domain has been shown *in vitro* to exhibit the highest recombination rate (6), the
286 maintenance of specific DBL α types at high frequencies and through extensive durations
287 could suggest that selection for adaptive advantages.

288

289 **2.6 Factors driving the conservation of DBL α types remain unknown**

290 The spatial study considered a few possible factors to explain these conserved DBL α types,
291 focusing specifically on 51 and 17 DBL α types in the high and very high location-averaged

292 frequency classes, respectively (i.e., a total of 68 DBL α types with location-averaged
293 frequencies of $\geq 5\%$) (Figure 5). Firstly, conservation of *var* genes on specific *P. falciparum*
294 chromosomes 4, 6, 7, and 8 has been previously reported and potentially attributed to
295 selective sweep events associated with antimalarial drug resistance (33). While positional
296 information is unavailable for the DBL α types analysed in this study, sequence alignments
297 show that only six of the 68 highly-frequent DBL α types are homologous to *var* genes on these
298 chromosomes. Furthermore, some of these highly-frequent DBL α types were identified as
299 homologs to the DBL α tags of five *var* genes in *P. praefalciparum*, a Plasmodium species that
300 naturally infects gorillas and is the closest living sister species of *P. falciparum* (41,42).
301 Homologs to the DBL α tag of one *var* gene of another ancestral Laverania species, *P.*
302 *reichenowi*, was identified but occurring at relatively low-to-moderate frequencies
303 (frequencies range from 0.57% to 2.34%). No homologs to the DBL α tags of *var* genes of *P.*
304 *gaboni* were identified. Hence, it is clear that while some of these factors can explain the
305 reason a few of these sequences are conserved, the majority of these conserved DBL α types
306 are still unaccounted for.

307

308 Homologs to other published globally-conserved DBL α types and *var* gene (PF3D7_0617400)
309 are shown in Figure 5. Tonkin-Hill et al. (35) reported a set of 100 most frequent DBL α types
310 conserved in their analysis of ten countries across diverse global regions in Africa,
311 Asia/Oceania, and South America. Homologs to 84 of these globally-conserved DBL α types
312 were identified in this study, with 26 and 11 of these types found in the two highest location-
313 averaged frequency classes. Furthermore, in the context of general prevalence in the
314 analysed African locations, 30 of these 37 DBL α types were found in all four locations, six in
315 three locations, and only one was found in a single location.

316

317 A conserved *P. falciparum* var gene (PF3D7_0617400) was also recently reported in a
318 Gabonese parasite isolate and characterised (34). The homolog to the DBL α tag of this var
319 gene was found to occur at high frequencies (ranging from 7.0% to 20.2% in different
320 locations) and present in all locations except, strangely, in Gabon itself. Notably, this
321 conserved PF3D7_0617400 var gene is located on chromosome 6, coinciding with the
322 previous reports of haplotypes in linkage disequilibrium on the same chromosome (43,44),
323 though this var gene is located outside of this region's cluster. A possible explanation for the
324 absence of this homolog in the Gabon dataset used in this study may be that the isolates were
325 sampled relatively early in the timeline (year 2000), which precedes the switch to artemisinin
326 (ART)-based combination therapies (ACT) in Africa (45,46), suggesting that the selection for
327 this specific type may have still been in progress and may not yet have risen in frequency to
328 result in observed fixation in the population at the time. Additionally, we also checked for
329 possible conservation of a *P. falciparum* var gene (PF3D7_0809100) that has been reported
330 to be expressed in sporozoites and potentially play a role in hepatocyte infection (47), but
331 homologs to the DBL α tag of this gene were found at only <1% frequency in three locations.

332

333 3. DISCUSSION

334

335 Extensive DBL α type diversity is reported in areas with high malaria transmission, generated
336 by meiotic and mitotic recombination (3–7) with DBL α repertoire diversity driven by frequent
337 outcrossing in the mosquito vector (1,2), such that we would not expect conservation of
338 types. However, a closer inspection of the population structure of DBL α types reveals
339 conservation of DBL α types *beyond* sequences found at very high frequencies or within the

340 highest percentiles. Instead, conservation also encompasses types that are seen stable
341 through time and can be found in a population at various frequencies, be it low, moderate,
342 or high. This study observed the conservation of DBL α types in the three major ups groups
343 within a large natural parasite population in a local area in Bongo. Frequencies of these DBL α
344 types at the population level were shown to be temporally stable over at least five years and
345 through wet and dry seasons. In addition to conservation at the continent level, spatial
346 analysis observed local conservation of specific high-frequency upsB and upsC types in
347 individual countries in Africa, despite the high genetic diversity typically reported for these
348 groups. Global analyses such as (35) and (34) would have uncovered the most conserved types
349 and genes globally but could have missed out on much of local signatures of conservation,
350 which this study has shown to exist within different frequency classes.

351

352 The key result of this analysis is that the frequency pattern of DBL α types that make up every
353 isolate repertoire not only underlies these local population structures but will maintain them.
354 Looking at individual DBL α types found in every isolate repertoire and the corresponding
355 frequency at which each type occurs in the population, per-isolate frequency profiles revealed
356 that every isolate repertoire consists of a mix of low-, moderate-, and high-frequency types,
357 in proportions consistent across all isolates. This presents a paradox in the population
358 structure of DBL α types, where there is a very high diversity of DBL α types found in a
359 population, but each isolate still maintains a combination of low- to moderate- to high-
360 frequency types in its repertoire. Even more interestingly, this paradoxical structure, both at
361 the level of population and isolate, was observed for types in all three groups of upsA, upsB,
362 and upsC and is maintained despite the expectation of frequent outcrossing in these endemic
363 areas.

364

365 In high transmission, despite the high rates of outcrossing and recombination, the consistency
366 in these per-isolate frequency profiles suggest a level of constraint on the modularity of each
367 isolate's repertoire, i.e., each isolate repertoire must have a combination of common and rare
368 types while maintaining limited overlaps with other isolates in the population overall. This
369 pattern was observed with both DBL α field data and DBL α encoding sequences identified in
370 assembled *var* exon 1 sequences, thereby excluding any biases from genotyping methods.
371 DBL α -*var* relationships revealed that high-frequency DBL α types are likely to be associated
372 with multiple distinct *var* exon 1 sequences (i.e., 1-to-many), though high sequence
373 similarities were estimated for some pairwise *var* sequences, suggesting that a proportion of
374 these are alleles of a same gene. We make clear that this study describes the conservation of
375 DBL α types, not necessarily the conservation of *var* genes. Conservation of *var* genes can be
376 better studied if we can properly define alleles.

377

378 Assuming that there are biological advantages conferred by these conserved types, why do
379 we observe stable presence at different frequencies but not fixation of most of these types in
380 the population? Furthermore, why are these per-isolate frequency profiles maintained? One
381 hypothesis relates to balancing selection as a result of co-evolution between the parasite and
382 the human host population it is infecting, which can occur at the local, regional, and continent
383 levels. The role of PfEMP1 in evading recognition by the host innate immune system would
384 select for its variation, and the DBL α domain has been shown to be immunogenic to variant-
385 specific epitopes and serologically recognised in an age-dependent manner (30). On the other
386 hand, its role in virulence, such as the need to bind to specific host endothelial cells for
387 cytoadhesion or blood cell receptors for rosetting, could select for some level of conservation.

388 This tension between the dual different roles, both of which relate to local host genetics,
389 could be creating the observed paradoxical pattern.

390

391 Stochastic simulations and network analyses have provided clear evidence for a role of
392 immune selection or negative frequency-dependent selection resulting from specific immune
393 memory, which is a form of balancing selection, in shaping antigenic diversity within natural
394 populations. As antibody-mediated immunity plays a significant role in recognition of PfEMP1
395 variants, we hypothesise that another possible driver of balancing selection is the arm's race
396 between the parasite PfEMP1 variants and host HLA class II haplotypes (48–51). Similar to our
397 finding of local signatures of DBL α type conservation against a highly-diverse background,
398 there are also geographic differences in HLA class II alleles across the African continent (52–
399 54). Immune evasion of common local HLA class II alleles could drive DBL α types or *var* genes
400 containing these types to persist through time at stable low-to-moderate frequencies.
401 Alternatively, genetic variation in a parasite's DBL α or *var* repertoires may have been shaped
402 by underlying differences in host receptors of varying spatial niches and if not, these types
403 could be in linkage disequilibrium with other proximal domains (e.g., CIDR) or genes vital to
404 these roles.

405

406 The consistency of these per-isolate frequency profile patterns is striking and suggests that
407 maintaining such frequency profiles within a parasite repertoire is advantageous to the
408 parasite. Having a range of rare to common types may allow malaria parasites to adapt to
409 host factors in order to persist through dynamics and competition within and between hosts.
410 The translational implication of this work suggests that breaking this pattern to what is seen

411 in low transmission i.e., high relatedness of *var* repertoires and clonality could be a target of
412 elimination efforts.

413

414 **4. MATERIALS AND METHODS**

415 **4.1 Data sources and types**

416 Conservation analyses were performed on a small ~450bp region of a *var* gene that encodes
417 a portion of the DBL α domain of PfEMP1 (i.e., DBL α tags) (55,56). DBL α tag sequences
418 included in this study were either generated from targeted amplicon sequencing (15–17) or
419 extracted from assembled *var* gene sequences (14,33,40). This made available DBL α tag
420 datasets of varying sizes from Africa and Asia, which were clustered to generate
421 representative DBL α types (see Section 4.2). However, the scope of this study on DBL α
422 conservation was limited to African locations only, with higher transmission, because lower
423 transmission areas may present a different context underlying conservation (e.g., clonality or
424 smaller population sizes). Data in Africa were available from West Africa (Senegal, Gambia,
425 Guinea, Mali, Ghana, Gabon), Central Africa (Congo, Malawi) and East Africa (Uganda, Kenya)
426 (Table I in Data S1, Table I in Data S3). However, most of these African countries were
427 excluded due to limited dataset sizes (number of isolates < 100), resulting in a final analysis
428 from four locations in Africa (i.e., Ghana, Gabon, Malawi, Uganda). Sources and methods that
429 the different studies used to generate these DBL α tag datasets are described in the following
430 subsections.

431

432 **4.1.1 DBL α tags from targeted amplicon sequencing data**

433 Published DBL α tag datasets from three locations were generated from targeted amplicon
434 sequencing (Table I in Data S1, Table I in Data S3). Amplicon sequencing of DBL α tag

435 sequences involves PCR amplification of a small sequence region encoding the DBL α domain
436 of PfEMP1 with degenerate primers (55,56), followed by high throughput sequencing on
437 either the Illumina MiSeq platform (Ghana, (13,15,39)) or on the 454 sequencing platform
438 (Gabon, (17)), Uganda, (16)). These include sequences from:

439 (I) One area (Bongo) in Ghana: dataset spans seven time points (surveys) from 2012
440 to 2017 through sampling of asymptomatic individuals through multiple dry and
441 wet seasons.

442 (II) Six areas (Apac, Arua, Jinja, Kanungu, Kyenjojo, Tororo) in Uganda: dataset
443 included sampling of clinical isolates over two years.

444 (III) One area (Bakoumba) in Gabon: dataset included sampling of asymptomatic
445 children in one year.

446

447 **4.1.2 DBL α tags from assembled *var* gene sequences**

448 Published *var* gene sequences (from isolates in Africa and Asia) were downloaded from the
449 ‘Full Dataset’ published by (33). DBL α tag sequences were identified and extracted from *var*
450 gene sequences (regardless of *var* gene completeness) as described in (14). Briefly, domain
451 annotations provided by (33) were used to extract nucleotide sequences encoding the DBL α
452 domain. These extracted sequences were further translated into the best reading frames and,
453 using *hmmsearch* (57), the resulting amino acid sequences were further searched against
454 positions 189 to 430 of the PFAM profile alignment (PF05424_seed.txt) to identify the ‘tag’
455 region (domain score cut-off of 60 and ≥ 100 aligned positions) and to ultimately extract the
456 DBL α tag sequence that would have been amplified with degenerate primers (55,56).

457

458 **4.2 Clustering of DBL α tags into DBL α types**

459 DBL α tags (Africa and Asia) were translated into amino acid sequences and any untranslatable
460 sequences (i.e., stop codons in reading frame) were excluded. The remaining DBL α tags were
461 combined and clustered with *clusterDBL α* (58), using a 96% nucleotide identity threshold (31)
462 to produce representative DBL α types. This also generated a binary matrix detailing the
463 presence/absence matrix of each DBL α type in each isolate.

464

465 **4.3 Classification of DBL α types into domain classes and ups groups**

466 The *classifyDBL α* pipeline (16) was used to classify DBL α types into DBL α domain classes of
467 DBL α 0, DBL α 1, or DBL α 2, in order to confirm that sequences were indeed those encoding the
468 DBL α domain of PfEMP1. In addition, a novel algorithm (*cUps*) described in this study was
469 used to classify DBL α types into the most probable ups group (i.e., upsA, upsB, or upsC),
470 accompanied by assignment probability values. For each DBL α type, ups groups were assigned
471 according to the prediction with the highest assignment probability. We describe this novel
472 classification algorithm below as well as in Data S2). An implementation of the algorithm is
473 available at <https://github.com/qianfeng2/cUps>.

474

475 Through the alignment and clustering of 2kb sequences upstream of *var* genes, followed by
476 the classification *var* genes into ups groups by Neighbour-joining (NJ) and Markov clustering
477 (MCL) methods (trees available in Data S2), a reference dataset of DBL α tag sequences was
478 generated from 846 *var* genes from 16 *P. falciparum* genomes ((59) and NCBI). We begin with
479 this reference database of DBL α tag sequences with ups groups and DBL α subclasses known.
480 For each category (ups group/DBL α subclass combination), we align the reference sequences
481 in the category using Clustal Omega v1.2.4 (60), then fit a profile hidden Markov model (61)
482 using HMMER v3.2.1 (57) with default settings.

483

484 For a given query sequence (representing a DBL α type), we calculate the likelihood of the
485 query sequence being drawn from the profile HMM of each category, using the forward
486 algorithm. The posterior probability for each category is then calculated using Bayes'
487 Theorem, with the prior probabilities of each category calculated from the reference
488 database. Summing over DBL α subclasses gives the posterior probability for each ups group
489 (i.e., assignment probability). The query sequence can be classified to the ups group with the
490 highest assignment probability. Although we do not do so in this paper, a threshold may
491 optionally be applied, so that sequences with highest assignment probability below the
492 threshold are categorised as 'unclassified'. Alternatively, a summary statistic may weight each
493 ups group by the assignment probability. This method is described in much more detail, with
494 verification (Feng, submitted).

495

496 **4.4 Exclusion of DBL α types, isolates, and populations from the final DBL α type dataset**

497 Only the DBL α types that were successfully classified into a DBL α domain class (i.e., DBL α 0,
498 DBL α 1, or DBL α 2) were retained in the final dataset. Subsequently, isolates with < 20 DBL α
499 types were also removed from dataset to ensure robust analyses downstream (Table I in Data
500 S1, Table I in Data S3). Specifically for the time-series dataset from the Malaria Reservoir Study
501 (MRS) in Bongo, Ghana (13,15,39), submicroscopic or symptomatic isolates were additionally
502 excluded from the dataset. Further, using *blastn* ($\geq 96\%$ nucleotide identity, $\geq 95\%$ query
503 coverage) (62), DBL α types with homology to isolate-transcendent *var1*, *var2csa*, and *var3*
504 sequences (sequences from (22,33)) were excluded to remove putative DBL α types previously
505 reported as isolate-transcendent (20,21). Finally, given that frequency classes and profiles
506 were calculated based on proportional frequencies, only locations with datasets of ≥ 100

507 isolates were retained. This resulted in the exclusion of six African countries from this study
508 (“*” in Table I in Data S3).

509

510 **4.5 Genetic similarity between pairwise isolate repertoires**

511 The pairwise type sharing metric (PTS) (31) is used to estimate the overlap between pairwise
512 isolate repertoires (e.g., isolates i & j). Specifically:

513

$$514 \quad PTS = \frac{2 * shared_{ij}}{Size_i + Size_j}$$

515

516 where $shared_{ij}$ is the number of shared DBL α types between repertoires of isolates i and j ,
517 and $Size_i$ and $Size_j$ are the total number of DBL α types (i.e., repertoire sizes) of isolates i and j ,
518 respectively. A value of 0 indicates the absence of sharing between two isolates whereas a
519 value of 1 indicates completely identical isolate repertoires.

520

521 **4.6 Calculation of DBL α type frequencies and assignments into frequency classes**

522 Depending on the analysis, a population can be the collection of isolates sampled at a specific
523 survey or time point in the time-series analyses (i.e., by year or survey in the MRS data) or the
524 collection of isolates sampled from a specific region or location/country in the spatial
525 analyses. Raw frequencies of DBL α types were defined at the survey or location level in counts
526 (i.e., number of isolates with a particular DBL α type in each survey or location). Raw (count)
527 frequencies were converted into proportional frequencies through division of count
528 frequencies by the total number of isolates at a corresponding time point or location, leading
529 to “survey-specific frequencies” or “location-specific frequencies”. Subsequently, these

530 frequencies were further categorised into frequency classes of 0%, **low** (0%, 1%), **moderate**
531 [1%, 5%), **high** [5%, 10%), and **very high** [10%, 100%].

532

533 Given the substantial differences in dataset sizes across surveys or locations (e.g., 499 isolates
534 for Uganda *versus* 176 isolates for Gabon), simply summing isolates across datasets of
535 multiple surveys or locations would bias total frequencies to reflect those of larger datasets.
536 Hence, averaged frequencies were used instead as a means to normalise total frequencies by
537 isolate counts in each survey or location. For example, a DBL α type found in 10 out of 100
538 isolates for location A and 10 out of 500 isolates for location B would be reported to have 10%
539 and 2% frequencies for locations A and B, respectively. This would yield a crude total
540 frequency of 3.33% (20 of 600 isolates), which is more reflective of the frequency observed
541 in location B even though the DBL α type was found at relatively high frequency at location A.
542 In this instance, with normalisation, an averaged frequency of 6% would be estimated (12 of
543 200 isolates), reducing the bias towards larger dataset sizes. Given the focus of this study on
544 conserved DBL α types, this normalisation method provides a less biased approach in
545 identifying DBL α types that are found at high frequencies but not necessarily uniformly across
546 all datasets.

547

548 **4.7 Determination of DBL α -*var* relationships**

549 For two locations (Ghana and Malawi), *var* gene sequences were available from assemblies
550 performed by (33). DBL α -*var* relationships were determined using complete *var* exon 1
551 sequences that are bounded by an N-terminal segment (NTS) and a transmembrane region
552 (TM) on the 5' and 3' ends of exon 1, respectively (14). Briefly, using *vsearch* (63), DBL α types
553 were globally aligned to *var* exon 1 sequences from the same location (e.g., Malawi DBL α

554 types to Malawi *var* exon 1). Given that DBL α types were generated from clustering at a 96%
555 nucleotide identity threshold, these types were aligned to *var* exon 1 sequences at the same
556 threshold of 96% identity, calculated over the alignment length and excluding terminal gaps
557 (*--iddef 2*). The relationship between a DBL α type and distinct *var* exon 1 was determined
558 based on the number of unique *var* exon 1 sequences sharing a same DBL α type (e.g., a 1-to-
559 *n* DBL α -*var* relationship is defined as a DBL α type found in *n* unique *var* exon 1).

560

561 For each group of *var* exon 1 that share a same DBL α type, an all vs all sequence alignment of
562 *var* exon 1 sequences in the group was performed using the *allpairs_global* option within
563 *vsearch* (63) and set to include all pairwise alignments (*--acceptall*). Pairwise nucleotide
564 identities were estimated based on calculations over whole alignment lengths, including
565 terminal gaps (*--iddef 1*), to account for differences in pairs of *var* exon 1 of variable lengths.

566

567 **4.8 Search for homology to other DBL α types or *var* genes**

568 **4.8.1 *Var* genes in association with selective sweeps on select chromosomes**

569 Published work reported conserved *var* genes on chromosomes 4, 6, 7 and 8 associated with
570 selective sweep events, potentially due to drug resistance or other factors. Accession
571 numbers of these conserved genes were obtained from the author (33,59) and used as
572 reference. Using *blastn* (62), DBL α types were searched against these reference sequences
573 and hits from alignments were reported ($\geq 96\%$ nucleotide identity, $\geq 95\%$ query coverage).

574

575 **4.8.2 *Var* genes in primate *Plasmodium* species**

576 *Var* genes from three *Plasmodium* species, *P. praefalciparum*, *P. reichenowi* and *P. gaboni*,
577 were downloaded from PlasmoDB (41) and used as reference. Using *blastn* (62), DBL α types

578 were searched against these reference sequences and hits from alignments were reported
579 ($\geq 96\%$ nucleotide identity, $\geq 95\%$ query coverage).

580

581 **4.8.3 Globally-conserved DBL α types or *var* genes**

582 The 100 most frequent DBL α sequences reported in the global analysis by (35) was used as
583 reference. Using *blastn* (62), DBL α types were searched against these reference sequences
584 and hits from alignments were reported ($\geq 96\%$ nucleotide identity, $\geq 95\%$ query coverage).

585 The same search parameters and thresholds were applied in searching for homologs to *var*
586 gene sequence PF3D7_0617400, a conserved 3D7 *var* gene reported by (34). Homologs to *var*
587 gene sequence PF3D7_0809100, shown by (47) to be expressed at the sporozoite stage, were
588 also searched for.

589

590 **ACKNOWLEDGMENTS**

591 This publication uses *var* gene data assembled from data generated from the MalariaGEN
592 *Plasmodium falciparum* Community Project. This research was supported by The University
593 of Melbourne's Research Computing Services and the Petascale Campus Initiative.

594

595 **REFERENCES**

596

- 597 1. Babiker HA, Ranford-Cartwright LC, Currie D, Charlwood JD, Billingsley P, Teuscher T,
598 et al. Random mating in a natural population of the malaria parasite *Plasmodium*
599 *falciparum*. *Parasitology*. 2009/04/06. 1994;109(4):413–21.
- 600 2. Paul REL, Packer MJ, Walmsley M, Lagog M, Ranford-Cartwright LC, Paru R, et al.
601 Mating Patterns in Malaria Parasite Populations of Papua New Guinea. *Science* (80-).
602 1995;269(5231):1709–11.
- 603 3. Freitas-Junior LH, Bottius E, Pirrit LA, Deitsch KW, Scheidig C, Guinet F, et al. Frequent
604 ectopic recombination of virulence factor genes in telomeric chromosome clusters of
605 *P. falciparum*. *Nature*. 2000;407(6807):1018–22.
- 606 4. Duffy MF, Byrne TJ, Carret C, Ivens A, Brown G V. Ectopic Recombination of a Malaria
607 *var* Gene during Mitosis Associated with an Altered *var* Switch Rate. *J Mol Biol*.

- 608 2009;389(3):453–69.
- 609 5. Bopp SER, Manary MJ, Bright AT, Johnston GL, Dharia N V, Luna FL, et al. Mitotic
610 Evolution of *Plasmodium falciparum* Shows a Stable Core Genome but Recombination
611 in Antigen Families. *PLOS Genet.* 2013;9(2):e1003293.
- 612 6. Claessens A, Hamilton WL, Kekre M, Otto TD, Faizullabhoj A, Rayner JC, et al.
613 Generation of Antigenic Diversity in *Plasmodium falciparum* by Structured
614 Rearrangement of Var Genes During Mitosis. *PLOS Genet.* 2014 Dec
615 18;10(12):e1004812.
- 616 7. Zhang X, Alexander N, Leonardi I, Mason C, Kirkman LA, Deitsch KW. Rapid antigen
617 diversification through mitotic recombination in the human malaria parasite
618 *Plasmodium falciparum*. *PLOS Biol.* 2019;17(5):e3000271.
- 619 8. Baruch DI. Adhesive receptors on malaria-parasitized red cells. *Best Pract Res Clin*
620 *Haematol.* 1999;12(4):747–61.
- 621 9. Newbold C, Craig A, Kyes S, Rowe A, Fernandez-Reyes D, Fagan T. Cytoadherence,
622 pathogenesis and the infected red cell surface in *Plasmodium falciparum*. *Int J*
623 *Parasitol.* 1999;29(6):927–37.
- 624 10. Chen Q, Schlichtherle M, Wahlgren M. Molecular Aspects of Severe Malaria. *Clin*
625 *Microbiol Rev.* 2000;13(3):439–50.
- 626 11. Smith JD, Chitnis CE, Craig AG, Roberts DJ, Hudson-Taylor DE, Peterson DS, et al.
627 Switches in expression of *Plasmodium falciparum* var genes correlate with changes in
628 antigenic and cytoadherent phenotypes of infected erythrocytes. *Cell.* 1995;82(1):101–
629 10.
- 630 12. Su X, Heatwole VM, Wertheimer SP, Guinet F, Herrfeldt JA, Peterson DS, et al. The large
631 diverse gene family var encodes proteins involved in cytoadherence and antigenic
632 variation of *Plasmodium falciparum*-infected erythrocytes. *Cell.* 1995;82(1):89–100.
- 633 13. Tiedje KE, Oduro AR, Bangre O, Amenga-Etego L, Dadzie SK, Appawu MA, et al. Indoor
634 residual spraying with a non-pyrethroid insecticide reduces the reservoir of
635 *Plasmodium falciparum* in a high-transmission area in northern Ghana. *PLOS Glob*
636 *Public Heal.* 2022 May 18;2(5):e0000285.
- 637 14. Tan MH, Shim H, Chan Y, Day KP. Unravelling var complexity: Relationship between
638 DBL α types and var genes in *Plasmodium falciparum*. Vol. 1, *Frontiers in Parasitology.*
639 2023.
- 640 15. Tiedje KE, Zhan Q, Ruybal-Pesántez S, Tonkin-Hill G, He Q, Tan MH, et al. Measuring
641 changes in *Plasmodium falciparum* census population size in response to sequential
642 malaria control interventions. *medRxiv.* 2023;
- 643 16. Ruybal-Pesántez S, Tiedje KE, Tonkin-Hill G, Rask TS, Kanya MR, Greenhouse B, et al.
644 Population genomics of virulence genes of *Plasmodium falciparum* in clinical isolates
645 from Uganda. *Sci Rep.* 2017;7(1):11810.
- 646 17. Day KP, Artzy-Randrup Y, Tiedje KE, Rougeron V, Chen DS, Rask TS, et al. Evidence of
647 strain structure in *Plasmodium falciparum* var gene repertoires in children from Gabon,
648 West Africa. *Proc Natl Acad Sci.* 2017;114(20):E4103–11.
- 649 18. Gardner MJ, Hall N, Fung E, White O, Berriman M, Hyman RW, et al. Genome sequence
650 of the human malaria parasite *Plasmodium falciparum*. *Nature.* 2002;419(6906):498–
651 511.
- 652 19. Lavstsen T, Salanti A, Jensen ATR, Arnot DE, Theander TG. Sub-grouping of *Plasmodium*
653 *falciparum* 3D7 var genes based on sequence analysis of coding and non-coding
654 regions. *Malar J.* 2003;2(1):27.

- 655 20. Kyes SA, Kraemer SM, Smith JD. Antigenic Variation in *Plasmodium falciparum*: Gene
656 Organization and Regulation of the var Multigene Family. *Eukaryot Cell*. 2007 Sep
657 1;6(9):1511–20.
- 658 21. Kraemer SM, Kyes SA, Aggarwal G, Springer AL, Nelson SO, Christodoulou Z, et al.
659 Patterns of gene recombination shape var gene repertoires in *Plasmodium falciparum*:
660 comparisons of geographically diverse isolates. *BMC Genomics*. 2007;8(1):45.
- 661 22. Rask TS, Hansen DA, Theander TG, Gorm Pedersen A, Lavstsen T. *Plasmodium*
662 *falciparum* Erythrocyte Membrane Protein 1 Diversity in Seven Genomes – Divide and
663 Conquer. *PLOS Comput Biol*. 2010 Sep 16;6(9):e1000933.
- 664 23. Walker IS, Rogerson SJ. Pathogenicity and virulence of malaria: Sticky problems and
665 tricky solutions. *Virulence*. 2023 Dec 31;14(1):2150456.
- 666 24. Smith JD, Rowe JA, Higgins MK, Lavstsen T. Malaria's deadly grip: cytoadhesion of
667 *Plasmodium falciparum*-infected erythrocytes. *Cell Microbiol*. 2013 Dec
668 1;15(12):1976–83.
- 669 25. Lau CKY, Turner L, Jespersen JS, Lowe ED, Petersen B, Wang CW, et al. Structural
670 Conservation Despite Huge Sequence Diversity Allows EPCR Binding by the PfEMP1
671 Family Implicated in Severe Childhood Malaria. *Cell Host Microbe*. 2015 Jan
672 14;17(1):118–29.
- 673 26. Turner L, Lavstsen T, Berger SS, Wang CW, Petersen JE V, Avril M, et al. Severe malaria
674 is associated with parasite binding to endothelial protein C receptor. *Nature*.
675 2013;498(7455):502–5.
- 676 27. Hsieh F-L, Turner L, Bolla JR, Robinson C V, Lavstsen T, Higgins MK. The structural basis
677 for CD36 binding by the malaria parasite. *Nat Commun*. 2016;7(1):12837.
- 678 28. Robinson BA, Welch TL, Smith JD. Widespread functional specialization of *Plasmodium*
679 *falciparum* erythrocyte membrane protein 1 family members to bind CD36 analysed
680 across a parasite genome. *Mol Microbiol*. 2003 Mar 1;47(5):1265–78.
- 681 29. Smith JD. The role of PfEMP1 adhesion domain classification in *Plasmodium falciparum*
682 pathogenesis research. *Mol Biochem Parasitol*. 2014;195(2):82–7.
- 683 30. Barry AE, Trieu A, Fowkes FJI, Pablo J, Kalantari-Dehaghi M, Jasinskas A, et al. The
684 Stability and Complexity of Antibody Responses to the Major Surface Antigen of
685 *Plasmodium falciparum* Are Associated with Age in a Malaria Endemic Area. *Mol Cell*
686 *Proteomics*. 2011;10(11).
- 687 31. Barry AE, Leliwa-Sytek A, Tavul L, Imrie H, Migot-Nabias F, Brown SM, et al. Population
688 Genomics of the Immune Evasion (var) Genes of *Plasmodium falciparum*. *PLOS Pathog*.
689 2007 Mar 16;3(3):e34.
- 690 32. Rougeron V, Tiedje KE, Chen DS, Rask TS, Gamboa D, Maestre A, et al. Evolutionary
691 structure of *Plasmodium falciparum* major variant surface antigen genes in South
692 America: Implications for epidemic transmission and surveillance. *Ecol Evol*. 2017 Nov
693 1;7(22):9376–90.
- 694 33. Otto TD, Assefa SA, Böhme U, Sanders MJ, Kwiatkowski DP, Null N, et al. Evolutionary
695 analysis of the most polymorphic gene family in *falciparum* malaria [version 1; peer
696 review: 1 approved, 2 approved with reservations]. *Wellcome Open Res*. 2019;4(193).
- 697 34. Dimonte S, Bruske EI, Enderes C, Otto TD, Turner L, Kremsner P, et al. Identification of
698 a conserved var gene in different *Plasmodium falciparum* strains. *Malar J*.
699 2020;19(1):194.
- 700 35. Tonkin-Hill G, Ruybal-Pesántez S, Tiedje KE, Rougeron V, Duffy MF, Zakeri S, et al.
701 Evolutionary analyses of the major variant surface antigen-encoding genes reveal

- 702 population structure of *Plasmodium falciparum* within and between continents. *PLOS*
703 *Genet.* 2021 Feb 25;17(2):e1009269.
- 704 36. Albrecht L, Merino EF, Hoffmann EHE, Ferreira MU, de Mattos Ferreira RG, Osakabe
705 AL, et al. Extense variant gene family repertoire overlap in Western Amazon
706 *Plasmodium falciparum* isolates. *Mol Biochem Parasitol.* 2006;150(2):157–65.
- 707 37. Albrecht L, Castiñeiras C, Carvalho BO, Ladeia-Andrade S, Santos da Silva N, Hoffmann
708 EHE, et al. The South American *Plasmodium falciparum* var gene repertoire is limited,
709 highly shared and possibly lacks several antigenic types. *Gene.* 2010;453(1):37–44.
- 710 38. Ruybal-Pesántez S, Sáenz FE, Deed SL, Johnson EK, Larremore DB, Vera-Arias CA, et al.
711 Molecular epidemiology of continued *Plasmodium falciparum* disease transmission
712 after an outbreak in Ecuador. Vol. 4, *Frontiers in Tropical Diseases.* 2023.
- 713 39. Tiedje KE, Oduro AR, Agongo G, Anyorigiya T, Azongo D, Awine T, et al. Seasonal
714 Variation in the Epidemiology of Asymptomatic *Plasmodium falciparum* Infections
715 across Two Catchment Areas in Bongo District, Ghana. *Am Soc Trop Med Hyg.*
716 2017;97(1):199–212.
- 717 40. MalariaGen, Ahouidi A, Ali M, Almagro-Garcia J, Amambua-Ngwa A, Amaratunga C, et
718 al. An open dataset of *Plasmodium falciparum* genome variation in 7,000 worldwide
719 samples [version 2; peer review: 2 approved]. *Wellcome Open Res.* 2021;6:42.
- 720 41. Otto TD, Gilabert A, Crellen T, Böhme U, Arnathau C, Sanders M, et al. Genomes of all
721 known members of a *Plasmodium* subgenus reveal paths to virulent human malaria.
722 *Nat Microbiol.* 2018;3(6):687–97.
- 723 42. Sharp PM, Plenderleith LJ, Hahn BH. Ape Origins of Human Malaria. *Annu Rev*
724 *Microbiol.* 2020 Sep 8;74(1):39–63.
- 725 43. Amambua-Ngwa A, Park DJ, Volkman SK, Barnes KG, Bei AK, Lukens AK, et al. SNP
726 Genotyping Identifies New Signatures of Selection in a Deep Sample of West African
727 *Plasmodium falciparum* Malaria Parasites. *Mol Biol Evol.* 2012;29(11):3249–53.
- 728 44. Amambua-Ngwa A, Danso B, Worwui A, Ceesay S, Davies N, Jeffries D, et al.
729 Exceptionally long-range haplotypes in *Plasmodium falciparum* chromosome 6
730 maintained in an endemic African population. *Malar J.* 2016;15(1):515.
- 731 45. Eastman RT, Fidock DA. Artemisinin-based combination therapies: a vital tool in efforts
732 to eliminate malaria. *Nat Rev Microbiol.* 2009;7(12):864–74.
- 733 46. Bosman A, Mendis. KN. A Major Transition in Malaria Treatment: The Adoption and
734 Deployment of Artemisinin-Based Combination Therapies. In: Breman JG, Alilio MS,
735 White NJ, editors. *Defining and Defeating the Intolerable Burden of Malaria III: Progress*
736 *and Perspectives: Supplement to Volume 77(6) of American Journal of Tropical*
737 *Medicine and Hygiene.* Northbrook (IL): American Society of Tropical Medicine and
738 Hygiene; 2007.
- 739 47. Zanghì G, Vembar SS, Baumgarten S, Ding S, Guizetti J, Bryant JM, et al. A Specific
740 PfEMP1 Is Expressed in *P. falciparum* Sporozoites and Plays a Role in Hepatocyte
741 Infection. *Cell Rep.* 2018 Mar 13;22(11):2951–63.
- 742 48. Hill AVS, Allsopp CEM, Kwiatkowski D, Anstey NM, Twumasi P, Rowe PA, et al. Common
743 West African HLA antigens are associated with protection from severe malaria. *Nature.*
744 1991;352(6336):595–600.
- 745 49. Hill AVS, Yates SNR, Allsopp CEM, Gupta S, Gilbert SC, Lalvani A, et al. Human Leukocyte
746 Antigens and Natural Selection by Malaria. *Philos Trans Biol Sci.* 1994 Jul
747 12;346(1317):379–85.
- 748 50. Lima-Junior J da C, Pratt-Riccio LR. Major Histocompatibility Complex and Malaria:

- 749 Focus on Plasmodium vivax Infection. Vol. 7, *Frontiers in Immunology*. 2016.
- 750 51. Peterson TA, Bielawny T, Kimani M, Ball TB, Plummer FA, Luo M, et al. Diversity and
751 frequencies of HLA class I and class II genes of an East African population. 2014;
- 752 52. Hill AVS, Elvin J, Willis AC, Aidoo M, Allsopp CEM, Gotch FM, et al. Molecular analysis
753 of the association of HLA-B53 and resistance to severe malaria. *Nature*.
754 1992;360(6403):434–9.
- 755 53. Goeury T, Creary LE, Brunet L, Galan M, Pasquier M, Kervaire B, et al. Deciphering the
756 fine nucleotide diversity of full HLA class I and class II genes in a well-documented
757 population from sub-Saharan Africa. *HLA*. 2018 Jan 1;91(1):36–51.
- 758 54. Sanchez-Mazas A. African diversity from the HLA point of view: influence of genetic
759 drift, geography, linguistics, and natural selection. *Hum Immunol*. 2001;62(9):937–48.
- 760 55. Taylor HM, Kyes SA, Harris D, Kriek N, Newbold CI. A study of var gene transcription in
761 vitro using universal var gene primers. *Mol Biochem Parasitol*. 2000;105(1):13–23.
- 762 56. Bull PC, Berriman M, Kyes S, Quail MA, Hall N, Kortok MM, et al. Plasmodium falciparum
763 Variant Surface Antigen Expression Patterns during Malaria. *PLOS Pathog*. 2005 Nov
764 18;1(3):e26.
- 765 57. Eddy SR. Accelerated Profile HMM Searches. *PLOS Comput Biol*. 2011;7(10):e1002195.
- 766 58. He Q, Pilosof S, Tiedje KE, Ruybal-Pesántez S, Artzy-Randrup Y, Baskerville EB, et al.
767 Networks of genetic similarity reveal non-neutral processes shape strain structure in
768 Plasmodium falciparum. *Nat Commun*. 2018;9(1):1817.
- 769 59. Otto TD, Böhme U, Sanders MJ, Reid AJ, Bruske EI, Duffy CW, et al. Long read
770 assemblies of geographically dispersed Plasmodium falciparum isolates reveal highly
771 structured subtelomeres [version 1; peer review: 3 approved]. *Wellcome Open Res*.
772 2018;3(52).
- 773 60. Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, et al. Fast, scalable generation
774 of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst*
775 *Biol*. 2011 Jan 1;7(1):539.
- 776 61. Eddy SR. Profile hidden Markov models. *Bioinformatics*. 1998 Jan 1;14(9):755–63.
- 777 62. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
778 architecture and applications. *BMC Bioinformatics*. 2009;10(1):421.
- 779 63. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool
780 for metagenomics. *PeerJ*. 2016;4:e2584.

781

782

783 FUNDING

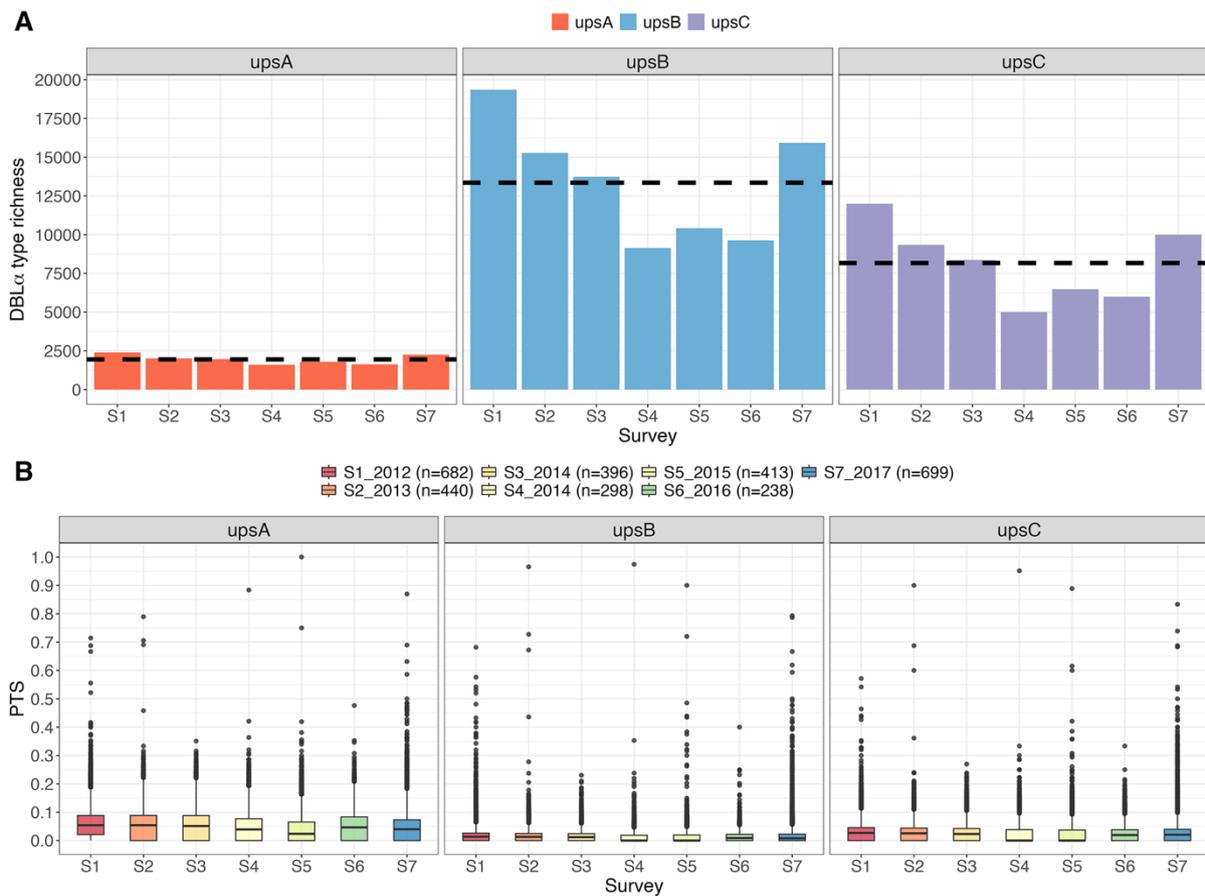
784 This study was funded by the National Institute of Allergy and Infectious Diseases, National
785 Institutes of Health through the joint NIH-NSF-NIFA Ecology and Evolution of Infectious
786 Disease award R01-AI149779 to KPD. Salary support for MT was provided by R01-AI149779.

787

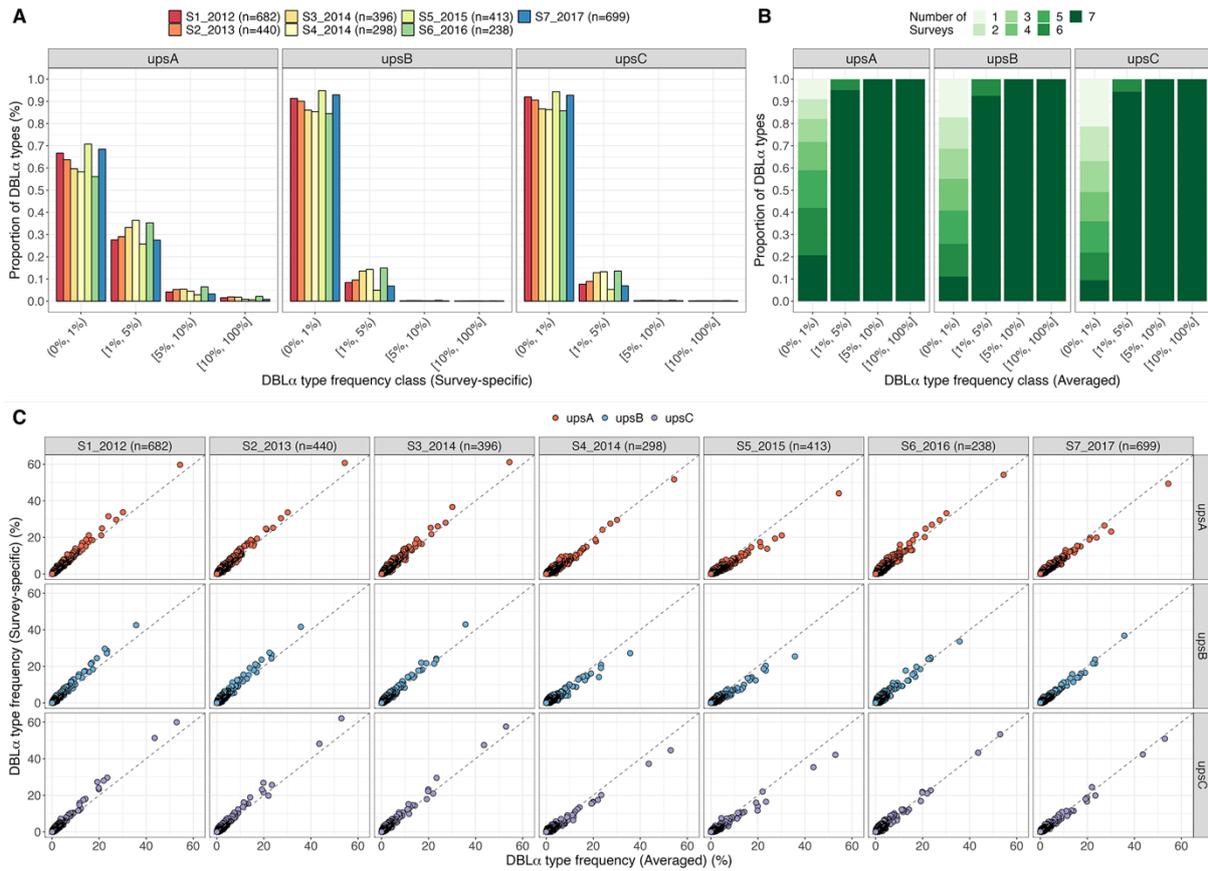
788

789

790 **FIGURES**

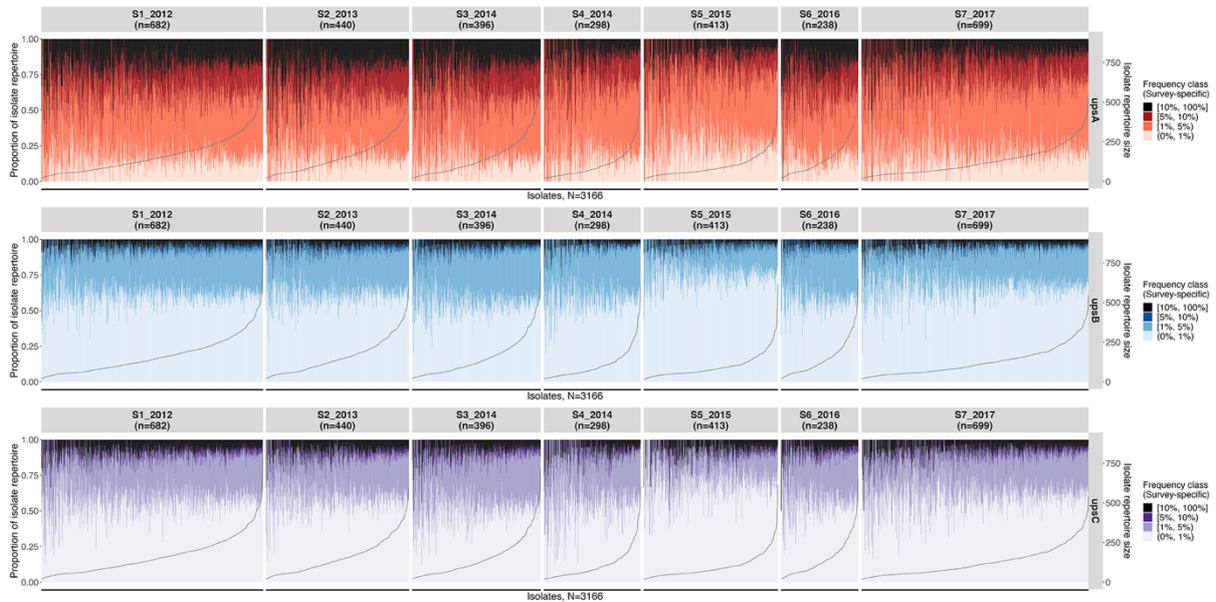


791
792 **Figure 1. Classification of DBL α types into ups groups (upsA, upsB, upsC) [Malaria Reservoir**
793 **Study (MRS)].** Figure shows (A) DBL α richness (horizontal dashed lines show mean richness
794 per ups group) and (B) genetic similarity (i.e., overlaps in isolate repertoire by pairwise type
795 sharing (PTS)), in each of the seven MRS surveys in Bongo, Ghana.
796



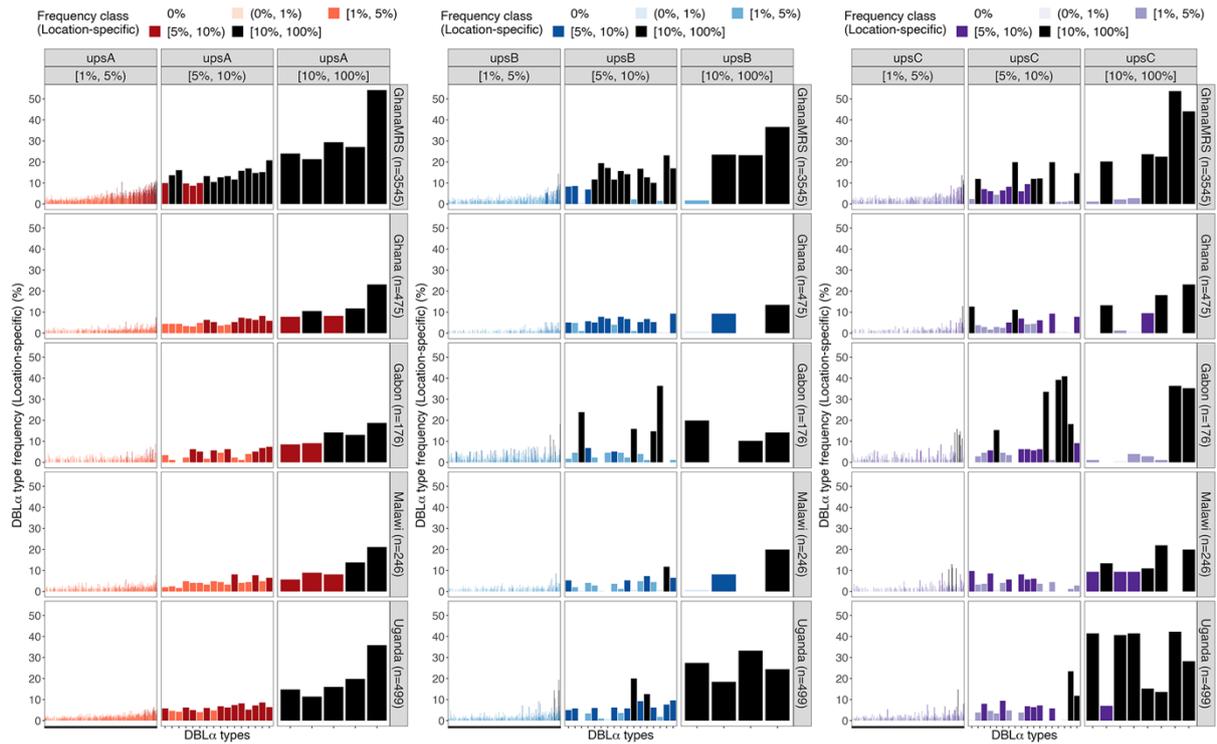
797
798
799
800
801
802
803
804
805
806
807
808
809

Figure 2. Conservation of DBL α types and frequencies in a local population and through time [Malaria Reservoir Study (MRS)]. (A) Distribution of survey-specific frequencies of DBL α types, binned into categorical frequency classes. (B) Number of surveys DBL α types were observed in, showing that DBL α types found at $\geq 1\%$ survey-averaged frequencies were seen to also persist through most surveys. This is also true for DBL α types found at $\geq 1\%$ survey-specific frequencies (Figure IV in Data S1). (C) Survey-specific frequencies (y-axis) of DBL α types are plot against survey-averaged frequencies (x-axis), showing positive correlation between both frequencies. Points represent individual DBL α types, coloured by ups groups.



810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825

Figure 3. Per-isolate frequency profiles show the composition of survey-specific frequency classes in every isolate repertoire [Malaria Reservoir Study (MRS)]. Vertical bars represent individual isolates, showing the composition of frequency classes in each isolate (left y-axis) by ups group (horizontal panels). Isolates are sorted by isolate repertoire size in increasing order, with isolate repertoire sizes indicated by the grey line (right y-axis).



826

827

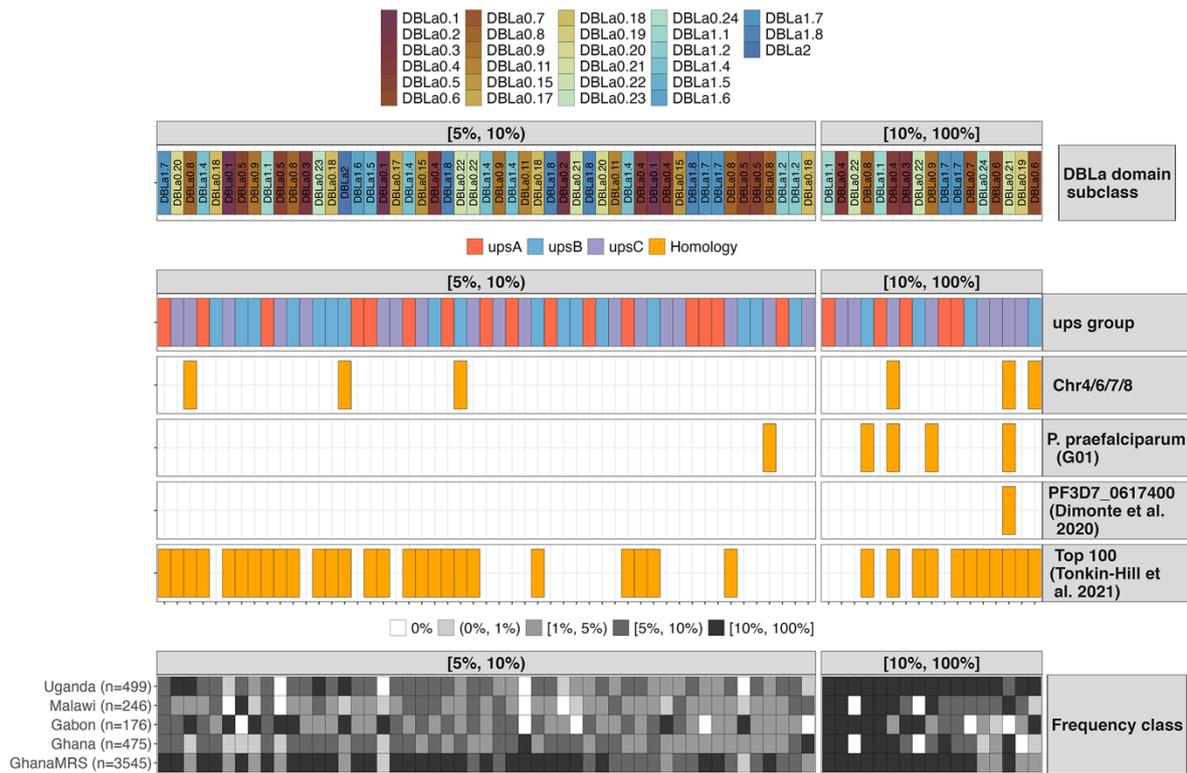
Figure 4. Conservation of DBL α types and frequencies at local and continent levels [spatial analysis]. Location-specific frequencies of individual DBL α types by ups groups (left to right: upsA, upsB, upsC). Shown here are DBL α types with $\geq 1\%$ location-averaged frequencies, ordered in increasing frequencies.

831

832

833

834



835
 836
 837
 838
 839
 840

Figure 5. Annotation of possible factors maintaining conserved DBL α types [spatial analysis]. Vertical panels indicate DBL α types with high or very high frequencies ($\geq 5\%$ location-averaged frequencies) and sequence homology (presence/absence) to published DBL α tag and *var* sequence data.