1    26 July 2019

2

3

4    **High-resolution micro-epidemiology of parasite spatial and temporal dynamics in a high**

5    **malaria transmission setting in Kenya**

6

7

8    Cody S. Nelson PhD[1]*, Kelsey M. Sumner MSPH[2,3], Elizabeth Freedman[3], Joseph W. Saelens

9    PhD[3], Andrew A. Obala PhD[4], Judith N. Mangeni MPH[5], Steve M. Taylor MD MPH[1,2,3+], Wendy

10   P. O'Meara PhD[1,3+]

11

12   [1]Duke Global Health Institute, Duke University, Durham, North Carolina, USA

13   [2]Department of Epidemiology, Gillings School of Global Public Health, University of North

14   Carolina, Chapel Hill, North Carolina, USA

15   [3]Division of Infectious Diseases, Duke University School of Medicine, Durham, North Carolina,

16   USA

17   [4]School of Medicine, Moi University College of Health Sciences, Eldoret, Kenya

18   [5]School of Nursing, Moi University College of Health Sciences, Eldoret, Kenya

19   [+]Authors contributed equally to this work.

20

21   *Please address correspondence to Cody S. Nelson (cody.nelson@duke.edu)

22

23

24   **Abstract WC:** 149 (150 max)

25   **Main Text WC:** 4,830 (5,000 max)

26   **Methods WC:** 2,450 (3,000 max)

**ABSTRACT**

Novel interventions that leverage the heterogeneity of parasite transmission are needed to push malaria further towards elimination. To better understand spatial and temporal dynamics of transmission, we applied amplicon NGS of two polymorphic gene regions (*csp* and *ama1*) to a cohort identified via reactive case detection in a high-transmission setting in western Kenya. From 4/2013–6/2014, we enrolled 442 symptomatic children with malaria, 442 matched controls, and all household members of both groups. We evaluated genetic similarity between infected individuals using three novel indices: sharing of parasite haplotypes on binary and proportional scales and the L1 norm. Symptomatic children more commonly shared haplotypes with their own household members. Furthermore, we identified robust temporal structuring of parasite genetic similarity that we exploited to identify the molecular signature of an outbreak. These findings of both micro- and macro-scale organization of parasite populations might be harnessed to inform next-generation malaria control measures.

**INTRODUCTION**

The global burden of malaria has decreased considerably: from 2000 until 2015, cases declined by 41% and malaria deaths fell by 62%[1]. This improvement is broadly associated with the adoption of core control measures, principally the use of long-lasting insecticide-treated bednets, improved case management with rapid diagnostic tests (RDTs), and treatment with artemisinin-combination therapies (ACTs). Yet in some areas such as Bungoma county in western Kenya, *Plasmodium falciparum* transmission has failed to decline in proportion to control efforts, underscoring the need for new strategies[2,3]. Owing to this, there is renewed interest in developing, testing, and deploying supplemental strategies to accelerate malaria elimination[4], including enhanced surveillance, mass drug administration, and active community-based case detection[3]. Some of these strategies have been deployed in low/moderate transmission settings with variable success[5,6], but have been less frequently trialed in high-transmission settings of sub-Saharan Africa. Apart from the operational challenges of such settings, a major impediment to the application of these tools is a limited understanding of the fine-scale heterogeneity in malaria risk and transmission. Gaps in knowledge include: the extent to which time and geographic space structure parasite populations, the introduction or propagation of novel parasite strains, and the ability to identify hosts related through discrete parasite transmission chains. A greater understanding of the dynamics of natural infections and their impact on parasite transmissibility could enable rational implementation of control measures to reduce the malaria disease burden in high-transmission settings.

Next-generation sequencing (NGS) technologies may enable both local- and population-level tracking of parasite transmission[7]. However, *P. falciparum* molecular surveillance efforts have been hindered by: 1) the size and complexity of the parasite genome[8], 2) the high prevalence of polygenomic infections in high-transmission areas[9,10], and 3) the complex life cycle involving sexual recombination[11]. Collectively, these challenges limit the ability to define sequence identity of parasites *in vivo*, which is essential for identifying transmission between hosts. One potential

66    solution is NGS of PCR amplicons, which enables resolution of parasite variants within a

67    population into an array of genetically-distinct haplotypes[12-14]. Using control parasite DNA

68    mixtures, we[14-16] and others[13] have previously demonstrated that amplicon sequencing results in

69    high-fidelity haplotype output with identity and frequency directly proportional to the input genetic

70    material. The high sensitivity of this technique coupled with the ability to parse out multiple

71    genotypes in polyclonal infection is highly appealing for the use of amplicon NGS data for malaria

72    molecular epidemiology[7].

73         Here we describe an investigation of malaria transmission across temporal and spatial

74    scales using amplicon NGS of two highly-polymorphic *P. falciparum* gene targets: those encoding

75    circumsporozoite protein (*csp*) and apical membrane antigen-1 (*ama1*). We applied this

76    genotyping approach to parasites collected during a case-control study of malaria in western

77    Kenya in 2013-14. Over a 15-month period, infected (case) and uninfected (control) index children

78    along with all household members of both groups were enrolled and tested with malaria RDTs.

79    Previously, we observed clustering of RDT-positive individuals, noting that infections were 2.5

80    times more common among the household members of cases compared to controls[2]. Though

81    entomology in this study was quite limited, we also identified clustering of larval sites and bloodfed

82    anopheline mosquitoes in case households; however, these relatively weak associations suggest

83    that vector proximity is not a primary driver of disease risk. Following genotyping of these same

84    samples, we subsequently utilized three novel metrics to assess interhost parasite genetic

85    similarity: binary haplotype sharing (any haplotypes in common), proportional haplotype sharing

86    (percentage of haplotypes in common), and the L1 norm (sequence-based distance). We

87    hypothesized that this genotyping approach coupled with metrics of genetic similarity would yield

88    information regarding spatial and temporal scales of malaria transmission that might be

89    harnessed to inform next-generation malaria control measures. Specifically, in this study we

90    predicted that parasite populations in case children would be more genetically similar to

91    asymptomatically-infected members of their own households than to infected members of other

4

92    households, and that the overall likelihood of genetically-similar haplotypes between any two

93    sampled individuals would be inversely related to geographic and temporal distance.

94    **RESULTS**

95    *Haplotype metrics and read coverage.*

96    Of 5,353 total study participants from across the study area (**Fig. 1**), 1,050 were RDT+ for

97    *P. falciparum*. A total of 966 RDT+ infections were submitted for amplicon NGS of *csp* and *ama1*

98    loci (**Fig. 1**). After haplotype assignment[17] and quality filtering of reads, we identified 120 unique

99    *csp* haplotypes across 655 participants and 180 *ama1* haplotypes across 667 participants (**Fig.**

100    **1**). In total, 617 samples (64% of infections initially submitted for sequencing) were assigned both

101    *csp* and *ama1* haplotypes. Compared to un-genotyped samples, the median parasite density was

102    nearly 2 orders of magnitude higher for samples successfully assigned *csp* (2.36 vs. 0.71,

103    p<0.001, Mann-Whitney U test) or *ama1* haplotypes (2.31 vs. 0.74, p<0.001, Mann-Whitney U

104    test) (**Supplementary Table 1**). In addition, a greater percentage of symptomatic children

105    enrolled by passive case detection ("case children" – CC) than asymptomatic household members

106    were successfully assigned *csp*/*ama1* haplotypes (**Supplementary Table 1**), though this is likely

107    a consequence of high parasite density in CC. Of participants successfully assigned *csp*

108    haplotypes, 43.6% were CC, 38.4% case household members (CHM), and 17.9% control

109    household members (**Table 1**). Very similar data were obtained for samples assigned *ama1*

110    haplotypes (**Table 1**). Furthermore, those assigned *csp* and *ama1* haplotypes were representative

111    of the overall population of RDT+ study participants. The number of reads per participant was

112    strongly correlated with $log_{10}$ parasite density for both *csp* (**Supplementary Fig. 2a**; $\rho$=0.54,

113    p<0.001, Spearman Rank test) and *ama1* (**Supplementary Fig. 2b**; $\rho$=0.47, p<0.001, Spearman

114    Rank test). Across all successfully-genotyped infections, the median read coverage was 13,369

115    for *csp* and 11,392 for *ama1* (**Supplementary Fig. 2c,d**).

116

117    *Multiplicity of infection*.

118    Most study participants had polygenomic infections: single haplotypes were detected in

119    only 34.7% (227/655; *csp*) and 33.9% (226/667; *ama1*) of genotyped infections (**Supplementary**

120    **Fig. 3a,b**). The median number of haplotypes detected at each locus per study participant was 2

121    (**Table 1, Supplementary Fig. 3a,b**), with maxima of 16 (*csp*) and 14 (*ama1*). Overall intrahost

122    nucleotide diversity was high, dominated by nonsynonymous sequence polymorphisms

123    (**Supplementary Fig. 4c,d**). Of 227 participants with monogenomic infection at the *csp* locus,

124    58.1% also had only a single *ama1* haplotype (and vice versa – 58.4%) (**Supplementary Fig.**

125    **3c**). Additionally, within individual participants, the number of *csp* and *ama1* haplotypes was highly

126    correlated (**Supplementary Fig. 3d**; *ρ*=0.68, p<0.001, Spearman Rank test). However, there was

127    no consistent difference between the MOI detected in CC vs. CHM within a single household

128    (p=ns, Wilcoxon Sign-Rank test).

129

130    *Case and haplotype distribution over time.*

131         The distribution of participants with *csp* (**Fig. 2a**) and *ama1* (**Supplementary Fig. 5a**)

132    haplotypes by month (April 2013 through June 2014) indicates year-round malaria transmission,

133    though notable seasonal variation with case incidence peaking during the rainy season

134    (approximately April through June). We observed heterogeneous haplotype persistence over time

135    among the 120 *csp* (**Fig. 2b**) and 180 *ama1* haplotypes (**Supplementary Fig. 5b**): 4 *csp* and 2

136    *ama1* haplotypes were detected in at least 14 of 15 months, and a large proportion (48% of *csp*

137    and 57% of *ama1*) appeared in only a single month (**Table 2**). The remaining haplotypes were

138    detected intermittently over the study period, not necessarily during consecutive months.

139         We tested if haplotype presence was impacted by age, because parasite density (and

140    thereby haplotype detection sensitivity) often depends upon host age in areas of endemic

141    transmission[18]. To do so, we computed the prevalence difference of each haplotype ($PD_H$)

142    between young children (≤5y) and older children/adults (>5y). However, we observed no

143    consistent difference in haplotype prevalence between the ≤5y and >5y populations

144    (**Supplementary Fig. 6**).

145

7

146

147    *Macro-level parasite genetic similarity and temporospatial structuring*

148    We next investigated the overall temporal and spatial structuring of *csp* (**Fig. 3**) and *ama1*

149    (**Supplementary Fig. 7**) haplotypes. Visual inspection of the spatial distribution of haplotypes

150    with variable duration (**Table 2** – 'persistent,' 'intermittent,' and 'sporadic') indicates that unique

151    variants are dispersed across geographic space during the time period in which they are detected

152    (**Fig. 3a-c; Supplementary Fig. 7a-c**). Furthermore, categorization of infections by administrative

153    locations during the high-transmission season in 2013 (April-June) revealed that individual

154    haplotypes were well mixed across the study area (**Fig. 3d**). Collectively, these findings suggest

155    a lack of spatial structuring of haplotypes.

156    We subsequently employed binary sharing, proportional sharing, and the L1 norm, novel

157    metrics which describe population-level genetic similarity, to examine *csp* and *ama1* genetic

158    similarity over time (**Fig. 3e-g; Supplementary Fig. 7d-f**) and between administrative locations

159    (**Fig. 3h-m; Supplementary Fig. 7g-l**). Binary sharing calculates whether a pair of infections

160    share any haplotype, while proportional sharing expresses the pairwise comparison of haplotype

161    sets as a proportion of the total number of haplotypes shared by infected individuals; the L1 norm

162    is a unitless index of the diversity of sequences between any two populations, which were here

163    defined as parasites constituting a single infection (see methods for details of calculation,

164    **Supplementary Fig. 8** for correlation of metrics). We computed each of these metrics for both

165    *csp* and *ama1* haplotypes between all study months by iterative, random sampling of individuals

166    from each month. Intriguingly, we identified a higher degree of parasite genetic similarity for

167    individuals sampled within a single month compared with those sampled during different months

168    (**Fig. 3e-g; Supplementary Fig. 7d-f)**, which was statistically-significant for binary sharing

169    (p=0.029), proportional sharing (p<0.001), and the L1 norm (p=0.013) (**Table S2**). Analogous

170    results were obtained for the temporal comparison of *ama1* haplotypes (**Table S2**). However,

171    when genetic similarity was assessed for each combination of sampled administrative locations

172   within a constrained window of time (highest malaria transmission season), we observed no

173   consistent differences between metrics for pairs collected in the same location compared to pairs

174   to those collected in different locations (**Table S3**). Overall, we observe spatially-homogenous

175   binary and proportional haplotype sharing (**Fig. 3h-k; Supplementary Fig. 7g-j)** though a

176   heterogenous L1 norm (**Fig. 3l,m; Supplementary Fig. 7k,l).** Collectively, these findings indicate

177   temporal structuring (though no clear spatial structuring) of parasite populations as defined by all

178   metrics of genetic similarity at both *csp* and *ama1* loci.

179

180   *Micro (household)-level parasite genetic similarity*

181        To investigate whether index cases and household members have similar parasite

182   populations, and thus whether asymptomatic parasitemia among household members may be a

183   risk factor for new clinical cases, we assessed genetic similarity between CC and asymptomatic

184   case household members (CHM) from their household of origin. As a comparison, we also

185   measured parasite genetic similarity between CC and members of an unrelated case household

186   (URCHM), which was selected based on time of sampling (given the strong evidence for temporal

187   macro-structuring of parasite populations). We evaluated CC:CHM and CC:URCHM genetic

188   similarity using binary sharing, proportional sharing, and the L1 norm. Binary sharing was

189   enhanced between CC and origin household CHM compared to binary sharing between CC and

190   URCHM when computed using both *csp* (**Fig. 4a**; p=0.02, Wilcoxon Signed-Rank test) and *ama1*

191   haplotypes (**Fig. 4b**; p=0.03, Wilcoxon Signed-Rank test). Similarly, proportional share scores

192   were higher for CC with their origin CHM for both *csp* and *ama1* (**Fig. 4c,d**; *csp* p=0.04, *ama1*

193   p=0.01, Wilcoxon Signed-Rank test). Finally, the L1 norm, which measures genetic distance

194   between parasite populations, was statistically reduced for CC to origin CHM compared to

195   URCHM at the *csp* (**Fig. 4e;** p=0.03) but not *ama1* locus (**Fig. 4f**). These results were robust to

196   alternative matching algorithms, including a combination of CC age, household location, and time

197   of sampling. Lastly, if we lift the restriction that a household be comprised of 3+ individuals to be

198     included in this analysis, the findings hold at the *csp* though not the *ama1* locus, possibly owing

199     to a greater diversity of *ama1* haplotypes overall (180) compared to *csp* (120) and the resulting

200     lower probability of observing exact matches.

201         We next examined the ability of genetic similarity metrics to predict the correct origin

202     household for each CC from among houses with 3 or more infected household members (n=38

203     households). To do so, we calculated aggregated *csp* and *ama1* binary sharing, proportional

204     sharing, and the L1 norm between each pairwise combination of CC and CHM, yielding 1,444

205     genetic relatedness values. We surmised that if binary sharing, proportional sharing, and the L1

206     norm are highly-predictive indicies, the calculated CC:CHM value should be greatest for the

207     comparison of CC with their own household members thus accurately identifying the CC

208     household of origin. Overall, binary sharing was the most predictive of the origin household of a

209     CC, with maximal binary share score correctly predicting the CC origin household 18% of the time

210     (7/38) compared with 16% (6/38) for proportional sharing and 11% (4/38) for the L1 norm (**Table**

211     **3**). However, none of these were significantly different compared to what would be expected by

212     random sampling alone (p=0.06, 0.11, and 0.20 respectively, Fisher's Exact test). We also

213     evaluated the ability of each sharing index to predict the time and geographic location of CC

214     infections by comparing CC:CHM scores across time and space. All three genetic similarity

215     metrics were highly predictive of CC temporal malaria acquisition, identifying the correct position

216     in time of the CC $\pm$ 30 days (~15% maximal temporal distance) for approximately 50% of CCs

217     (**Table 3**). All three indices were generally less predictive of CC spatial position, pinpointing the

218     CC origin household within 2.25 km (~15% maximal geographic distance) for approximately 25%

219     of CCs (**Table 3**).

220

221     *Outbreak: temporally-restricted unique combination of haplotypes among case children*

222         To investigate the temporal and spatial relationship between parasite populations in CC,

223     we computed pairwise binary sharing, proportional sharing, and L1 norm metrics for all CC

10

224    infections (*csp* n=283; *ama1* n=288), then calculated the temporal (**Fig. 5**) and geographic

225    (**Supplementary Fig. 9**) distance between CC pairings. When plotted against time between

226    enrollment, we can clearly discern that *csp*/*ama1* binary and proportional sharing as well as

227    parasite population sequence divergence (L1 norm) is highly dependent upon temporal distance

228    between CC (**Fig. 5**). Furthermore, among geographically-proximal infections (**Supplementary**

229    **Fig. 7**), there is some enhanced binary sharing at both *csp* and *ama1* loci (**Supplementary Fig.**

230    **9a,d**) though no clear overall trend of increased genetic similarity.

231        Owing to the apparent temporal structuring of CC parasite haplotypes, we tested for

232    specific clusters of parasite haplotypes in time among CC by computing the prevalence difference

233    (PD) of each haplotype between CC and CHM (PD$_H$) (**Fig. 6a,b**) during each month of study

234    enrollment. Comparing the monthly PD$_H$ of each haplotype, we determined that 4 haplotypes (*csp*

235    H8, H48, and H54 as well as *ama1* H13) were significantly more common in CC than CHM during

236    June 2013 (each p < 0.0001 by Fisher Exact test) (**Fig. 6a,b**). We examined CC from 5/6/2013–

237    7/29/2013, noting all those infected with parasites bearing *csp* H8, H48, and H54 also had

238    evidence of *csp* H1 (**Fig. 6c**). Likewise, from 5/6/2013–7/29/2013 all CC in which *ama1* H13 was

239    detected also had *ama1* H5 and H8 (**Fig. 6d**). In total we identified 26 CC with *csp* H1 + H8 + H48

240    + H54 and 27 with *ama1* H5 + H8 + H13 (**Fig. 6c,d**). Intriguingly, we observed substantial overlap

241    of this set of haplotypes in case children: 23 CC had evidence of all haplotypes combined (*csp*

242    H1 + H8 + H48 + H54 and *ama1* H5 + H8 + H13) (**Fig. 6e**). In comparison, from 5/6/2013–

243    7/29/2013 no CHM had either of these unique haplotype combinations (**Fig. 6f,g**). Thus, by

244    comparing haplotype prevalences in CC and CHM, we identified a combination of haplotypes that

245    co-occur within a temporally-restricted window and are associated with clinical disease.

246        Interestingly, we observe that this unique combination of *csp* and *ama1* haplotypes largely

247    occurred in CC during a 3-week period from 6/17/2013 to 7/8/2013, peaking during the week of

248    June 24th (**Fig. 6h**). We defined an outbreak 'case' as the presence of 5 or more of the 7 outbreak

249    haplotypes (*csp* H1/H8/H48/H54 and *ama1* H5/H8/H13), comprising more than 98% of the reads

11

250    detected in an individual. Employing this definition, we identified a total of 29 outbreak cases and

251    48 non-outbreak cases among the 77 total CC between 5/6/2013 and 7/29/2013.  We see a clear

252    peak of genetically-homogenous outbreak cases that accounts for nearly all cases during this 3-

253    week period and is book-ended with endemic transmission of non-outbreak cases (**Fig. 6i**).

254    Intriguingly, this outbreak event was not geographically confined, but rather was widely

255    disseminated across the study area (**Fig. 6j**). The seemingly random geographic distribution was

256    supported by a test for spatial structure of the outbreak haplotype combination in SaTScan, which

257    identified 0 high clusters and 1 low cluster (Relative Risk: 0.00, p=0.71) – a null result.

**DISCUSSION**

In this study, we utilized amplicon deep sequencing of two *P. falciparum* polymorphic gene targets in parasites collected from households in western Kenya to investigate the geographic and temporal structuring of parasite populations. Our analyses indicate that temporally-proximate infections have enhanced genetic similarity, suggesting that parasite populations are at least partially structured by time. Furthermore, we report that children with malaria are more likely to share parasite haplotypes with asymptomatically-infected members of their own household compared with members of time-matched households. Finally, we utilized temporal structuring of parasite populations to identify the genetic signature of an outbreak of parasite genotypes manifest as enrichment of specific parasite haplotype combinations during a single month in 2013. Collectively, our analysis identifies both micro and macro-level organization of parasite populations, enhancing our understanding of malaria transmission heterogeneity.

To our knowledge, this investigation presents the first genetic data directly linking parasites causing asymptomatic infections among household members with those causing symptomatic disease. Similar to prior studies[19-23], we previously noted that CC in this cohort were 2.5 times more likely to have RDT+ household members than were control children[2] indicating household-level hotspots of high-risk individuals. Herein, we extended this association using our parasite genotyping approach to quantify the degree of genetic similarity between infected individuals. For example, symptomatic children had a much higher median probability of sharing a *csp* haplotype with asymptomatic members of their own household (50%) than with those of a matched household (0%; p=0.02, Wilcoxon Sign-Rank test), indicating that parasites infecting symptomatic children are more genetically similar to those in their own household than to the population at large (**Fig. 4**). These findings expand upon the prior recognition of household-clustering of malaria risk by providing direct evidence that symptomatic children are participating in the same parasite transmission network as their household members. What remains unknown is the direction of this transmission, because we cannot assess whether specific haplotype-

13

284    identical parasites were transmitted from infected household members to CC, or whether both CC

285    and household members acquired haplotype-identical parasites from the same source. Better

286    understanding this phenomenon, through further high-resolution investigations into transmission

287    networks within individual households, will be critical to discern whether treatment of

288    asymptomatic infections might reduce the burden of clinical malaria disease.

289        Our amplicon NGS approach detected measurable temporal structuring of parasite

290    populations in this geographically-restricted, high-transmission setting, though no spatial

291    structuring beyond the household unit. First, we observed temporal clusters of increased genetic

292    similarity, with enhanced average similarity between the infections in any given month

293    (**Supplementary Table 2**, **Fig. 3, Supplementary Fig. 7**). Furthermore, genetic similarity metrics

294    assessed between CC and household members was reliably more predictive of household

295    temporal proximity (**Table 3**). Lastly, among CC we observed a strong inverse relationship

296    between temporal distance and binary/proportional sharing (direct relationship between temporal

297    distance and the L1 norm) (**Fig. 5**). Thus, all investigations point to a strong time-dependency of

298    parasite transmission in this endemic setting, particularly among case children. Intriguingly,

299    individual haplotypes (**Fig. 3, Supplementary Fig. 7**) and even unique haplotype combinations

300    (**Fig. 6**) appear evenly distributed across the study area at any single point in time. The

301    mechanism of this long-range transmission (cases up to 20km apart) of genetically-identical

302    infections is perplexing, since the flight range of unfed *Anopheles gambiae* has been measured

303    at a maximum of 3km.[24] Nonetheless, this robust temporal structuring of haplotypes has profound

304    implications for malaria elimination in high-transmission settings, suggesting that strategies

305    focused purely on cluster-based control/prevention may fail to prevent forward parasite

306    transmission and disease.

307        Notably, this is the first study to detect the genetic signature of an outbreak amidst

308    endemic malaria transmission. While several investigations have recognized genetic

309    homogeneity in small-scale outbreaks[25,26] and traced the spread of drug-resistance alleles[27], we

14

310 identified strong genetic evidence of a malaria outbreak nested within the temporal structure of

311 haplotype-sharing among CC, namely the appearance in symptomatic children of parasites with

312 identical *csp* and *ama1* combinations of haplotypes nearly simultaneously across the study area.

313 As has been suggested[7] these ensembles of parasites, when detected in linked individuals,

314 provide compelling evidence of a shared origin of infections. Some outbreak-associated

315 haplotypes (*csp* H1/H8 and *ama1* H5/H8) were exceedingly common throughout the study period

316 and present in both CC and CHM during the majority of months, while other outbreak haplotypes

317 (*csp* H48/H52 and *ama1* H13) are rare and appear almost exclusively among CC from May-July

318 2013. One CC infected with all 7 *csp*/*ama1* outbreak haplotypes reported a travel history to

319 another malaria-endemic region within the past 3 months, so we might hypothesize that rare

320 haplotypes (*csp* H48/H54 and *ama1* H13) were imported from outside the study area. The reason

321 for the co-occurrence of this unique combination of haplotypes among symptomatic children is

322 unclear, and two major questions remain unanswered by this investigation: 1) How did the

323 outbreak spread nearly simultaneously across a relatively large geographic area and 2) why was

324 this combination only detected among CC? The geographic co-occurrence of outbreak cases may

325 be contingent upon unconventional and/or undescribed vector movement and biting behavior or

326 cryptic human movements within the study area. Strain-specific immune response have been

327 associated with protection from disease[28]– including those directed against a vaccine-elicited *csp*

328 haplotype[13,29]– and therefore, it is possible that rare haplotypes represent new, antigenically-

329 diverse malaria strains to the study area (not previously encountered by CC) leading to enhanced

330 clinical disease. Alternatively, these rare haplotypes could be genomically-linked with novel

331 virulence factors. Though we cannot discern the origin of vector transmission and/or cause of

332 pathogenesis with any degree of certainty using our amplicon NGS data, the ability to detect and

333 monitor genetically-identical polyclonal infections over time demonstrates the power of this

334 amplicon NGS sequencing approach for malaria molecular epidemiology.

335    Our sampling frame and genetic analytic approaches enabled us to detect signatures of

336    parasite population structure in a high-transmission setting. While malaria cases are recognized

337    to be clustered spatiotemporally[30-36], prior molecular epidemiology studies have generally failed

338    to observe parasite population genetic structure (principally by geographic location) using a

339    variety of genotyping methods, sampling schemes, and analytic tools[37-42]; one recent notable

340    exception found a significant decay in genetic relatedness as a function of increasing distance

341    using microsatellite genotyping of polygenomic infections, but reported from a low-transmission

342    setting in Namibia[43]. In our high-transmission setting, we observed that parasite populations were

343    structured 1) by household and 2) by time more than geography by the application of three novel

344    metrics (binary sharing, proportional sharing, and the L1 norm) to systematically interrogate the

345    genetic similarity of polygenomic parasite populations present within individual hosts expressed

346    as haplotypes of unlinked parasite genes *csp* and *ama1*. The concordant results from parallel

347    analyses using these unlinked gene targets extends the credibility and utility of our novel metrics

348    as tools.  Two of these metrics, binary sharing (any shared haplotypes) and proportional sharing

349    (percentage of shared haplotypes), express the genetic relationship between two infections as a

350    function of the presence or absence of identical individual haplotypes within the pair. The premise

351    of these metrics is that haplotypes that are identical by state represent parasites identical by

352    descent; although this assumption would not be reasonable over large distances or time periods

353    owing to a variety of biological and epidemiologic factors, the constrained temporal and spatial

354    scales and dense sampling of our study render this a more plausible assumption. The third metric,

355    the L1 norm, has been used in viral genomic epidemiology as a measure of the sequence-based

356    genetic divergence of two pathogen populations[44]. The polygenomic nature of most infections in

357    high-transmission settings has precluded the use of numerous traditional tools of population

358    genetics, but the use of our metrics enabled us to exploit these polygenomic infections, by

359    capturing and then linking diverse haplotypes in separate hosts. We suggest that binary and

360    proportional sharing metrics produce highly similar results in our analyses, whereas the L1 norm

16

361   results are somewhat distinct and heterogeneous (**Fig. 3, Supplementary Figs. 7,8**). We

362   anticipate that appropriate use of these metrics will depend on local epidemiology, namely

363   parasite genetic diversity, transmission intensity, and prevalence of parasitemias. Thus, we

364   propose these metrics ought to be applied to diverse datasets to define the context of their utility.

365   Nevertheless, we hypothesize these high-resolution genetic metrics will enable investigators to

366   identify connectivity between polygenomic infections on more granular temporal and geographic

367   scales.

368   What are the implications of our findings for transmission-reducing strategies? The

369   identification of similar parasite populations between case children and their household of origin

370   supports notions in the malaria transmission field that symptomatic malaria infection may be

371   partially fueled by asymptomatic infection of household members. Yet we also present strong

372   evidence for temporal structuring of parasite populations and episodic transmission of genetically-

373   identical parasite haplotypes, unbounded by household or geography, suggesting that intra-

374   household sharing is not the only source of new infections. At face value these findings are

375   contradictory, which emphasizes the complexity of local and population-level transmission

376   networks. Nonetheless, our findings are most relevant to reactive programs, in which index cases

377   trigger either ring-testing (as in RACD) or ring-treatment (as in ring focal drug administration) in

378   an effort to mitigate foci of transmission. These programs are popular and literature supports their

379   efficacy at identifying additional cases, though there are insufficient data regarding cost-

380   effectiveness and the impact of transmission intensity or diagnostic test sensitivity for these

381   methods to be fully embraced. This investigation provides direct evidence that clinically-silent

382   parasite transmission chains within households are an important risk factor (but not the exclusive

383   source) of new infections, which supports the rationale for employing reactive strategies to

384   interrupt household-level transmission.  Yet, our data also suggests that parasite populations are

385   structured more by time than space, and therefore that household-level interventions *may not*

386   have measurable effects on community-level risk. Thus, whether transmission foci extend into

17

387   surrounding households, and to what extent mitigating them with reactive strategies contributes

388   to a reduction in aggregate community transmission, remains to be tested by future studies.

389        The greatest limitation of this investigation is the lack of longitudinal sampling of

390   cases/household members. We anticipate that a longitudinal dataset would enrich our

391   understanding of parasite transmission dynamics, including the directionality and time scale of

392   transmission, temporal fluctuations in haplotype frequency and parasite density, and the impact

393   of parasite density upon the probability of onward transmission (and thus the clinical import of

394   low-density infections). Another limitation is the study protocol resulted in a dataset that is 1)

395   biased towards a young age range (given the tendency for infected children to have higher

396   parasite density) and 2) is dominated by CC + direct household members (CC neighbors and

397   community members were excluded apart from control households) and thus we could not test

398   for fine-scale decay in sharing by distance or empirically define a distance threshold for the

399   observed genetic sharing. Additionally, our analyses of haplotype sharing between polygenomic

400   infections require acceptance of the assumption that parasites identical by descent can be

401   inferred from haplotypes identical by state, and this can be undermined by convergent evolution

402   on sequences by balancing selection and constrained haplotype frequencies across a dataset,

403   which collectively produce identical matches by chance. These risks, however, were mitigated in

404   our study owing to the fine temporal and spatial scales of sampling and the abundance of unique

405   haplotypes at each locus; in addition, the concordant findings in analyses using haplotypes

406   derived from unlinked gene targets further suggests that matching by chance was common.

407   Finally, as with any study requiring PCR amplification and DNA sequencing, there is the potential

408   for systematic error (primer bias, contamination, etc) to impact the experimental findings. We

409   reduced the risk of systematic error by enforcing strict sequencing read quality criteria, using

410   validated haplotype inference tools, and measuring effects using unlinked polymorphic targets in

411   orthogonal analyses of *csp* and *ama1*.

412      In conclusion, this is the first study to develop and apply indices of genetic similarity

413      between infected individuals to explore dynamics of parasite transmission at both household and

414      population levels. Collectively, our data suggest that *Pf* haplotypes are structured more clearly by

415      time than space, but also that, at the household level, children with malaria share parasite

416      genotypes with asymptomatically-infected household members who may constitute a risk factor

417      for childhood malaria. Subsequent longitudinal, high-resolution studies ought to investigate how

418      parasite populations change over time as well as the origin of symptomatic infections – are

419      parasites associated with malaria disease transmitted from household members or is their origin

420      external to the household? These definitive findings could have widespread implications for the

421      next-generation of malaria control efforts to combat heterogeneous transmission and ultimately

422      to eliminate malaria disease.

423    **METHODS**

424    *Study design, enrollment, and sample collection*

425    Participants were enrolled between 18-April-2013 and 5-June-2014 as part of a case-

426    control study as previously reported[2]. Briefly, children admitted to the Webuye County Hospital

427    with a diagnosis of malaria (confirmed by SD Bioline Pf HRP2 RDT) who resided within the six

428    administrative sublocations immediately surrounding the hospital were eligible for enrollment as

429    case children (CC). Household members of CC were tested by RDT and provided a dried blood

430    spot (DBS) within 1-7 days of enrolling the CC. At that time, CC were matched by age and village

431    to an RDT-negative control child and all of the members of the control child's household were

432    similarly tested and provided a DBS. A household was defined in this study as all family members

433    and individuals residing under a single shared roof. All household members of case and control

434    children were RDT-tested and DBS obtained at a single point in time immediately following child

435    enrollment in the study. While case households were matched to control based on geographic

436    proximity, neighbors and community members residing in close proximity to the enrolled

437    household were not necessarily tested or sampled. All DBS were stored at 4C with desiccant,

438    then shipped to Duke University (Durham, NC, USA) for subsequent analysis. Consent was

439    obtained from all participants prior to enrollment. The study protocol and consent procedures were

440    reviewed and approved by the Moi University Institutional Research and Ethics Committee

441    (IREC/2013/13) and the Duke University Institutional Review Board (Pro00044098).

442

443    *DNA isolation and quantitative PCR*

444    A single 6mm diameter punch from each DBS was deposited into a unique well of a deep

445    96-well plates *Genomic* DNA (gDNA) was extracted from each using a Chelex-100 protocol[45] and

446    reconstituted in 100 $\mu$L nuclease free water. We estimated parasite densities in all RDT-positive

447    infections using a real-time PCR assay targeting the *pfr364* motif in the parasite genome[46]. gDNA

448    from each sample was tested in duplicate, and densities were estimated from standard curves on

449     each reaction plate computed from amplification of a series of quantitative standards ranging from

450     1 to 2,000 parasites/uL of whole blood.

451

452     *csp/ama1 target amplification, library preparation, and sequencing*

453         Parasites were genotyped from each DBS obtained from RDT+ participants. Dual-indexed

454     libraries for each target were prepared using a nested PCR strategy and then pooled for NGS on

455     an Illumina MiSeq platform. Polymorphic segments of *P. falciparum* genes encoding

456     circumsporozoite protein (*csp;* 288bp) and apical membrane antigen 1 (*ama1*; 300bp) were

457     amplified in separate reactions from gDNA using primers that each included an identical overhang

458     sequence (**Supplementary Dataset 1**). PCR1 reactions consisted of 3 $\mu$L of template gDNA, 1.5

459     $\mu$M of each primer, 2 mM of MgCl$_2$, 300 $\mu$M each dNTP, 0.5 units of KAPA HiFi HotStart Taq

460     (Roche), and nuclease-free water to a total reaction volume of 25 $\mu$L; cycling conditions were 95C

461     x 3' →(98C x 20s →62C x 15s →72C x 25s) x 35 →72C x 1'. PCR1 products were used as the

462     template for PCR2 reactions, which used forward and reverse primers that annealed to PCR1

463     overhang sequences and also contained a MiSeq adaptor and a unique 8-mer index sequence

464     (**Supplementary Dataset 2**). PCR2 reactions consisted of 1.5 $\mu$L of template, 200 nM of each

465     primer, 1.5 mM of MgSO4, 200 $\mu$M each dNTP, 0.1 units of Platinum Taq High Fidelity (Thermo-

466     Fisher Scientific), and nuclease-free water to a total reaction volume of 25 $\mu$L; cycling conditions

467     were 94C x 2'→(94C x 15s →72C x 10s →68C x 30s) x 15 68C x 5'. PCR2 products were

468     amplicons including requisite adaptors for MiSeq sequencing as well as unique combinations of

469     2 MiSeq indices to enable disaggregation of reads by sample after sequencing. Separate library

470     pools were prepared for *csp* and *ama1* targets by combining an equal volume of PCR2 products

471     from each reaction. Libraries were purified and concentrated with ethanol, electrophoresed in a

472     2% agarose gel, gel-purified using QIAquick gel extraction kit (Qiagen), and subsequently cleaned

473     using Ampure XP beads (Agencourt). Finally, *csp* and *ama1* pools were combined in equimolar

21

474    fashion into a single sequencing pool. The resulting single pool was divided between two MiSeq

475    (v3 300-cycle PE) runs.

476

477    *Haplotype inference*

478    *csp* and *ama1* haplotypes were inferred from Illumina sequence reads. First, using

479    Trimmomatic[47], CutAdapt[48], and BBmap[49], reads were mapped to the *P. falciparum* strain 3D7

480    reference sequences for *csp* or *ama1*, then primer sequences (and nucleotides directly adjacent)

481    were trimmed from the sequences. For quality filtering of mapped reads, we used a sliding window

482    to remove reads if the average Phred quality score over 4 adjacent nucleotides was < 15.  These

483    quality-filtered reads were input into DADA2 (version 1.8) in order to join paired-end reads, call

484    haplotypes, and remove chimeras[17]. Within DADA2, read quality was further enforced prior to

485    haplotype inference using the Phred quality score for each read to model error frequency,

486    removing reads with greater than the predicted errors. To ensure only high-quality haplotypes

487    were included in the analysis data set, we censored haplotypes supported by fewer than 100

488    reads or present in less than 1% of the data.

489

490    *Amplicon variation and nucleotide diversity*

491    Variation across the amplicon was assessed by Shannon Entropy ($Hn$), which is a unitless

492    measure of variability calculated as:

493    $$Hn = -\sum_{i=0}^{n} p_i \, \log_e p_i$$

494    where n is the number of possible nucleotides at each position (A,C,T,G), and p the frequency of

495    the variant nucleotides being compared. The sum of Shannon entropy at each nucleotide position

496    represents the composite entropy score for the full amplicon. Nucleotide diversity ($\pi$) was

497    computed as the average distance between each possible pair of sequences using[50]:

498    $$\pi = \frac{\sum_i^H \sum_{j \leq i}^H d_{ij} f_i f_j}{L * N(N-1)/2}$$

499     where $L$ is sequence length in nucleotides for $\pi$, $N$ is total number of reads per participant, $d_{ij}$ is

500     number of nucleotide differences between haplotype $i$ and $j$, and $f_i/f_j$ is the number of reads

501     belonging to haplotype $i$. $\pi_S$ and $\pi_N$ were calculated as the average $dS$ and $dN$ between pairs of

502     haplotypes weighted by the haplotype's abundance:

503

$$\pi_{S,N} = \frac{\sum_i^H \sum_{j \le i}^H d_{S_{ij}, N_{ij}} f_i f_j}{L * N(N-1)/2}$$

504     where $d_{Sij}$ is $dS$ between haplotype $i$ and $j$ sequences and $d_{Nij}$ is $dN$ between haplotype $i$ and $j$

505     sequences. $\pi$, $\pi_S$, and $\pi_N$ were compared by Friedman test + posthoc Wilcoxon Signed-Rank test

506     (two-tailed). Of note, *csp* amplicon sequence variability (**Supplementary Fig. 4a,b**) and the

507     proportion of nucleotide diversity ($\Pi$) that occurs at nonsynonymous nucleotide positions ($\Pi_N$) are

508     consistent with previous reports[12].

509

510     *Development of genetic similarity metrics*

511         To quantify the degree of genetic similarity between individuals or groups, we defined two

512     novel metrics termed "binary sharing" and "proportional sharing". Binary sharing (possible values:

513     1 = true, 0 = false) was defined as the presence in any two sampled individuals of at least one

514     identical haplotype. Proportional sharing (continuous; possible values 0.0 to 1.0) was defined as

515     the proportion of identical haplotypes shared between any two sampled individuals (A and B):

516

$$PHS_{AB} = \frac{A \cap B}{A \cup B}$$

517     Complementary to these, to assess the genetic distance across samples, we utilized the unitless

518     L1 and L2 norm[44]:

519

$$L_1 = \sum_{k=0}^N \sum_{i=0}^n |p_i - q_i|$$

520

$$L_2 = \sum_{k=0}^N \sqrt{\sum_{i=0}^n (p_i - q_i)^2}$$

521     where $N$ is the sequence length in nucleotides, $n$ the number of possible nucleotides at each

522     position (A,C,T,G), and $p_i/q_i$ the frequencies of the variant nucleotides being compared. To

523    calculate genetic similarity metrics between two groups, bootstrap iterations of each metric were

524    performed upon a random sampling of pairs of individuals selected from each group (with

525    replacement) and the average computed (10,000 iterations for binary/proportional sharing and

526    100 iterations for L1/L2 given the computational requirements). L2 was not utilized for analysis in

527    this study because L1 and L2 are directly proportional and nearly perfectly correlated at both *csp*

528    (**Supplementary Fig. 8a**; $\rho$=1.0, Spearman Rank test) and *ama1* (**Supplementary Fig. 8b**;

529    $\rho$=0.99, Spearman Rank test). Binary/proportional share scores and the L1 norm were validated

530    by assessing the mean values for each household (3+ members), sub-location, location, and the

531    overall population for both *csp* and *ama1* (**Supplementary Fig. 8c-h**). On average, as genetic

532    similarity metrics are compared between larger and larger groupings of study participants, the

533    median decreases (increases for L1 norm) and there is enhanced central tendency as the range

534    narrows to approximate the computed metric for the overall population. The calculated genetic

535    similarity metrics for each group were compared by Kruskal Wallis test + posthoc Mann-Whitney

536    U test (two-tailed).

537

538    *Population-level temporal and spatial structure of genetic similarity*

539    To assess overall temporal structure of genetic similarity, we computed mean binary

540    sharing, proportional sharing, and the L1 norm metrics between each pairing of study months

541    (n=15, yielding 120 pairwise comparisons). These metrics were computed by iterative random

542    sampling (with replacement) of individuals between months as described above. Likewise, we

543    assessed spatial structure by calculating each genetic similarity metrics between each pairing of

544    administrative locations (n=5, yielding 15 pairwise comparisons) for a restricted time window

545    (April-June 2013 or April-June 2014) by repetitive random sampling of individuals from each

546    administrative location. A restricted time window was utilized to limit the impact of time as a

547    confounding factor that might obscure spatial structuring of haplotypes. We compared metrics

548    computed between pairs of months or pairs of administrative locations using a nonparametric

549    Mann-Whitney U test, with a p-value<0.05 considered to be significant.

550          We limited assessment of spatial structure to inter-household distance owing to: 1) prior

551    mapping of larval sites around a subset of these households, which indicate that larval sites are

552    numerous, small, and transient[2]; 2) vector behavior in western Kenya, which occurs after 9pm[51]

553    and thus renders the household as epicenter of parasite transmission; and 3) the absence in the

554    very circumscribed study site (~100 km$^2$) of other candidate features (e.g. rivers, lakes, mountain

555    ranges).

556

557    *Genetic similarity between case children and households*

558          We limited the assessment of genetic similarity between case child (CC) and case

559    household members (CHM) to the subset of households with parasite sequence data from 3 or

560    more infected members such that CC sharing could be assessed against at least 2 CHM for each

561    household. We enforced this constraint to make the analysis more conservative by mitigating the

562    risk of sampling error resulting in spurious findings regarding genetic similarity between CC and

563    a single household member. This subset consisted of 41 households for *csp* and 45 households

564    for *ama1*. We matched each of these households to one other household based on sampling date

565    using a weighted scale (5 points for ± 5 days, 4 points for ± 10 days, 3 points for ± 15 days, 2

566    points for ± 30 days, and 1 point for ± 60 days). In the event of a tie matching score, we randomly

567    selected a single household. All households were successfully matched to at least one other

568    household +/- 30 days. Furthermore, alternative matching algorithms employing case child age

569    and household geographic location were also trialed and produced similar results.

570          Next, for each CC (n = 41 for *csp* and n=45 for *ama1*) we computed the genetic similarity

571    between CC and CHM from their *household of origin* as well as CHM from the *matched*

572    *household*. Each comparison relied upon random, iterative sampling of CHM as previously

573    described. The outputs of this process were pairwise (between CC and CHM) estimates of the

25

574    likelihood of CC parasite genetic similarity with members of their own household in comparison

575    to members of a matched household. Finally, we compared haplotype genetic similarity metrics

576    of CC with origin CHM to those of CC with matched household members by Wilcoxon Signed-

577    Rank test (two-tailed).

578

579    *Temporal and spatial predictive value of CC:household genetic similarity*

580        We estimated the value of these genetic similarity metrics to predict temporal and

581    geographic order in the dataset within the 38 households with *csp* and *ama1* haplotypes

582    representing 3+ household members. *csp* and *ama1* genetic similarity metrics (binary sharing,

583    proportional sharing, L1 norm) were calculated for the pairing of every case with household

584    members from each case household. Values for *csp* and *ama1* were normalized and combined

585    into binary sharing, proportional sharing, and L1 norm composite scores. On the basis of these

586    composite scores, we rank-ordered households for each case child. Subsequently, we computed

587    the geographic and temporal distance between each case child and each rank-ordered

588    household. Of note, the Fisher's Exact test comparison of observed vs. expected frequency in

589    **Table 3** is based upon the assumption that cases are distributed evenly throughout time and

590    across geographic space.

591

592    *Genetic similarity between infections in case children*

593        We analyzed the genetic relatedness between infections in all passively-detected case

594    children by computing pairwise genetic similarity metrics for each unique pairing of infections

595    (n=283 for *csp* yielding 40,044 pairwise comparisons; n=288 for ama1 yielding 41,472 pairwise

596    comparisons). We analyzed each pairwise estimate as a function of temporal distance (expressed

597    as the interval in year between the dates of collection) or geographic distance, (expressed in

598    kilometers), and modeled the relationships between these metrics using locally-estimated

599    scatterplot smoothing (LOESS).

600

601     *Interrogation of haplotypes for enhanced prevalence by age and in case children*

602     We computed the prevalence of each unique *csp* (n=120) and *ama1* (n=180) haplotype

603     by month between individuals ≤5 years vs. >5 years old, defined as $\mathrm{age\ PD}_H$:

$$csp\ age\ PD_{H_i}\ \forall_{1<i<120} = P_{H_{i,\ \leq 5y}} - P_{H_{i,\ >5y}}$$

$$ama1\ age\ PD_{H_j}\ \forall_{1<j<180} = P_{H_{j,\ \leq 5y}} - P_{H_{j,\ >5y}}$$

606     where $P_{H_{i,\ \leq 5y}}$ is the prevalence of *csp* haplotype *i* during a given month in individuals ≤5y and

607     $P_{H_{i,\ >5y}}$ is the prevalence of that *csp* haplotype *i* during a given month in >5y (with analogous

608     calculations for *ama1* haplotype *j*). Furthermore, we assessed haplotype prevalence in case

609     children and in asymptomatic case household members, and used these to calculate the

610     prevalence difference of each haplotype as $\mathrm{CC:CHM\ PD}_H$:

$$csp\ CC{:}CHM\ PD_{H_i}\ \forall_{1<i<120} = P_{H_{i,CC}} - P_{H_{i,CHM}}$$

$$ama1\ CC{:}CHM\ PD_{H_j}\ \forall_{1<j<180} = P_{H_{j,CC}} - P_{H_{j,CHM}}$$

613     where $P_{H_{i,CC}}$ is the prevalence of *csp* haplotype *i* during a given month in CC and $P_{H_{i,CHM}}$ is the

614     prevalence of that *csp* haplotype *i* during a given month in CHM (with analogous calculations for

615     *ama1* haplotype *j*). We compared within each month the PD$_H$ for each haplotype using a Fisher's

616     Exact test (two-tailed); for these comparisons, we applied a Bonferroni correction for the alpha as

617     p = 0.05 / 505 observations ~ 1e-4. We used these PD$_H$ values to identify haplotypes that were

618     significantly over-represented in specific age groups or in CC compared to CHM. This

619     investigation served as a basis to empirically define outbreak cases of parasites bearing either

620     *csp* or *ama1* haplotypes. We assessed spatial structure of outbreak cases vs. non-outbreak

621     children using a purely spatial Bernoulli model in SaTScan (version 9.4.4)"[52].

622

623     *Statistical analysis*

624      All statistical tests were two-tailed. Details of individual statistical tests are described in

625    relevant sections throughout the methods.

626

627    *Data and code availability*

628      Haplotype sequences are available in NCBI GenBank under accessioning numbers

629    MK933826 – MK933945 (*csp*) and MK933946 – MK934125 (*ama1*). All analyses were completed

630    in R version 3.5.1. Source data and code for analyses is available for download from GitHub:

631    https://github.com/codynelson08/MESAseq_compiled_code. Raw analysis results available upon

632    request.

**AUTHOR CONTRIBUTIONS**

W.P.O., A.A.O, and J.M. designed and carried out cohort study; B.F. processed and sequenced samples; C.S.N., K.M.S., and J.S. analyzed data; S.M.T. and W.P.O. contributed expertise; and C.S.N., S.M.T., and W.P.O. wrote the paper.

**COMPETING INTERESTS**

The authors have no conflicts of interest to declare.

**REFERENCES**

1    World Malaria Report 2016. (World Health Organization, Geneva, 2016).
2    Obala, A. A. *et al.* What Is Threatening the Effectiveness of Insecticide-Treated Bednets? A Case-Control Study of Environmental, Behavioral, and Physical Factors Associated with Prevention Failure. *PloS one* **10**, e0132778, doi:10.1371/journal.pone.0132778 (2015).
3    A Framework for Malaria Elimination. (World Health Organization, Geneva, 2017).
4    The Global Technical Strategy for Malaria 2016-2030. (World Health Organization, Geneva, 2015).
5    von Seidlein, L. *et al.* The impact of targeted malaria elimination with mass drug administrations on falciparum malaria in Southeast Asia: A cluster randomised trial. *PLoS medicine* **16**, e1002745, doi:10.1371/journal.pmed.1002745 (2019).
6    Larsen, D. A. *et al.* Malaria surveillance in low-transmission areas of Zambia using reactive case detection. *Malar J* **14**, 465, doi:10.1186/s12936-015-0895-9 (2015).
7    Wesolowski, A. *et al.* Mapping malaria by combining parasite genomic and epidemiologic data. *BMC Med* **16**, 190, doi:10.1186/s12916-018-1181-9 (2018).
8    Murray, L. *et al.* Microsatellite genotyping and genome-wide single nucleotide polymorphism-based indices of Plasmodium falciparum diversity within clinical infections. *Malar J* **15**, 275, doi:10.1186/s12936-016-1324-4 (2016).
9    Greenhouse, B. *et al.* Validation of microsatellite markers for use in genotyping polyclonal Plasmodium falciparum infections. *Am J Trop Med Hyg* **75**, 836-842 (2006).
10   Happi, C. T. *et al.* Molecular analysis of Plasmodium falciparum recrudescent malaria infections in children treated with chloroquine in Nigeria. *Am J Trop Med Hyg* **70**, 20-26 (2004).
11   Chang, H. H. *et al.* Malaria life cycle intensifies both natural selection and random genetic drift. *Proceedings of the National Academy of Sciences of the United States of America* **110**, 20129-20134, doi:10.1073/pnas.1319857110 (2013).
12   Early, A. M. *et al.* Host-mediated selection impacts the diversity of Plasmodium falciparum antigens within infections. *Nature communications* **9**, 1381, doi:10.1038/s41467-018-03807-7 (2018).
13   Neafsey, D. E. *et al.* Genetic Diversity and Protective Efficacy of the RTS,S/AS01 Malaria Vaccine. *The New England journal of medicine* **373**, 2025-2037, doi:10.1056/NEJMoa1505819 (2015).
14   Levitt, B. *et al.* Overlap Extension Barcoding for the Next Generation Sequencing and Genotyping of Plasmodium falciparum in Individual Patients in Western Kenya. *Scientific reports* **7**, 41108, doi:10.1038/srep41108 (2017).
15   Patel, J. C. *et al.* Increased risk of low birth weight in women with placental malaria associated with P. falciparum VAR2CSA clade. *Scientific reports* **7**, 7768, doi:10.1038/s41598-017-04737-y (2017).
16   Taylor, S. M. *et al.* Absence of putative artemisinin resistance mutations among Plasmodium falciparum in Sub-Saharan Africa: a molecular epidemiologic study. *The Journal of infectious diseases* **211**, 680-688, doi:10.1093/infdis/jiu467 (2015).
17   Callahan, B. J. *et al.* DADA2: High-resolution sample inference from Illumina amplicon data. *Nat Methods* **13**, 581-583, doi:10.1038/nmeth.3869 (2016).
18   Baird, J. K. Host age as a determinant of naturally acquired immunity to Plasmodium falciparum. *Parasitol Today* **11**, 105-111 (1995).
19   Sturrock, H. J. *et al.* Targeting asymptomatic malaria infections: active surveillance in control and elimination. *PLoS medicine* **10**, e1001467, doi:10.1371/journal.pmed.1001467 (2013).
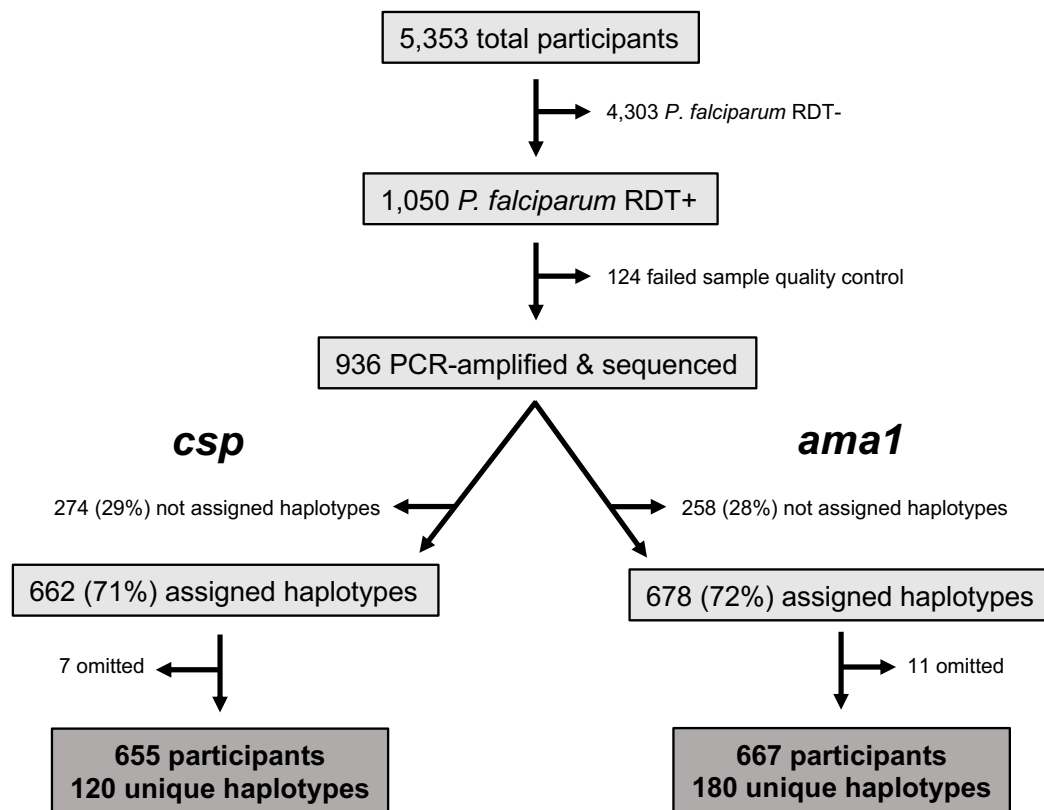
20  Njama-Meya, D., Kamya, M. R. & Dorsey, G. Asymptomatic parasitaemia as a risk factor for symptomatic malaria in a cohort of Ugandan children. *Trop Med Int Health* **9**, 862-868, doi:10.1111/j.1365-3156.2004.01277.x (2004).

21  Gahutu, J. B. *et al.* Prevalence and risk factors of malaria among children in southern highland Rwanda. *Malar J* **10**, 134, doi:10.1186/1475-2875-10-134 (2011).

22  Lindblade, K. A., Steinhardt, L., Samuels, A., Kachur, S. P. & Slutsker, L. The silent threat: asymptomatic parasitemia and malaria transmission. *Expert Rev Anti Infect Ther* **11**, 623-639, doi:10.1586/eri.13.45 (2013).

23  Stresman, G. H. *et al.* A method of active case detection to target reservoirs of asymptomatic malaria and gametocyte carriers in a rural area in Southern Province, Zambia. *Malar J* **9**, 265, doi:10.1186/1475-2875-9-265 (2010).

24  Kaufmann, C. & Briegel, H. Flight performance of the malaria vectors Anopheles gambiae and Anopheles atroparvus. *J Vector Ecol* **29**, 140-153 (2004).

25  Baldeviano, G. C. *et al.* Molecular Epidemiology of Plasmodium falciparum Malaria Outbreak, Tumbes, Peru, 2010-2012. *Emerg Infect Dis* **21**, 797-803, doi:10.3201/eid2105.141427 (2015).

26  Obaldia, N., 3rd *et al.* Clonal outbreak of Plasmodium falciparum infection in eastern Panama. *The Journal of infectious diseases* **211**, 1087-1096, doi:10.1093/infdis/jiu575 (2015).

27  Amato, R. *et al.* Origins of the current outbreak of multidrug-resistant malaria in southeast Asia: a retrospective genetic study. *The Lancet. Infectious diseases* **18**, 337-345, doi:10.1016/S1473-3099(18)30068-9 (2018).

28  Cheesman, S., Raza, A. & Carter, R. Mixed strain infections and strain-specific protective immunity in the rodent malaria parasite Plasmodium chabaudi chabaudi in mice. *Infection and immunity* **74**, 2996-3001, doi:10.1128/IAI.74.5.2996-3001.2006 (2006).

29  Fluck, C. *et al.* Strain-specific humoral response to a polymorphic malaria vaccine. *Infection and immunity* **72**, 6300-6305, doi:10.1128/IAI.72.11.6300-6305.2004 (2004).

30  Woolhouse, M. E. *et al.* Heterogeneities in the transmission of infectious agents: implications for the design of control programs. *Proceedings of the National Academy of Sciences of the United States of America* **94**, 338-342 (1997).

31  Bejon, P. *et al.* A micro-epidemiological analysis of febrile malaria in Coastal Kenya showing hotspots within hotspots. *Elife* **3**, e02130, doi:10.7554/eLife.02130 (2014).

32  Bjorkman, A. *et al.* Spatial Distribution of Falciparum Malaria Infections in Zanzibar: Implications for Focal Drug Administration Strategies Targeting Asymptomatic Parasite Carriers. *Clin Infect Dis* **64**, 1236-1243, doi:10.1093/cid/cix136 (2017).

33  Mogeni, P. *et al.* Effect of transmission intensity on hotspots and micro-epidemiology of malaria in sub-Saharan Africa. *BMC Med* **15**, 121, doi:10.1186/s12916-017-0887-4 (2017).

34  Gaudart, J. *et al.* Space-time clustering of childhood malaria at the household level: a dynamic cohort in a Mali village. *BMC Public Health* **6**, 286, doi:10.1186/1471-2458-6-286 (2006).

35  Ernst, K. C., Adoka, S. O., Kowuor, D. O., Wilson, M. L. & John, C. C. Malaria hotspot areas in a highland Kenya site are consistent in epidemic and non-epidemic years and are associated with ecological factors. *Malar J* **5**, 78, doi:10.1186/1475-2875-5-78 (2006).

36  Stresman, G. H. *et al.* Do hotspots fuel malaria transmission: a village-scale spatio-temporal analysis of a 2-year cohort study in The Gambia. *BMC Med* **16**, 160, doi:10.1186/s12916-018-1141-4 (2018).

37  Ingasia, L. A., Cheruiyot, J., Okoth, S. A., Andagalu, B. & Kamau, E. Genetic variability and population structure of Plasmodium falciparum parasite populations from different

753     malaria ecological regions of Kenya. *Infect Genet Evol* **39**, 372-380,
754         doi:10.1016/j.meegid.2015.10.013 (2016).
755  38  Duffy, C. W. *et al.* Population genetic structure and adaptation of malaria parasites on
756         the edge of endemic distribution. *Mol Ecol* **26**, 2880-2894, doi:10.1111/mec.14066
757         (2017).
758  39  Carrel, M. *et al.* The geography of malaria genetics in the Democratic Republic of
759         Congo: A complex and fragmented landscape. *Social science & medicine* **133**, 233-241,
760         doi:10.1016/j.socscimed.2014.10.037 (2015).
761  40  Pacheco, M. A. *et al.* Limited differentiation among Plasmodium vivax populations from
762         the northwest and to the south Pacific Coast of Colombia: A malaria corridor? *PLoS Negl*
763         *Trop Dis* **13**, e0007310, doi:10.1371/journal.pntd.0007310 (2019).
764  41  Zhu, X. *et al.* Analysis of Pvama1 genes from China-Myanmar border reveals little
765         regional genetic differentiation of Plasmodium vivax populations. *Parasit Vectors* **9**, 614,
766         doi:10.1186/s13071-016-1899-1 (2016).
767  42  Omedo, I. *et al.* Geographic-genetic analysis of Plasmodium falciparum parasite
768         populations from surveys of primary school children in Western Kenya. *Wellcome Open*
769         *Res* **2**, 29, doi:10.12688/wellcomeopenres.11228.2 (2017).
770  43  Tessema, S. *et al.* Using parasite genetic and human mobility data to infer local and
771         cross-border malaria connectivity in Southern Africa. *Elife* **8**, doi:10.7554/eLife.43510
772         (2019).
773  44  Poon, L. L. *et al.* Quantifying influenza virus diversity and transmission in humans. *Nat*
774         *Genet* **48**, 195-200, doi:10.1038/ng.3479 (2016).
775  45  Plowe, C. V., Djimde, A., Bouare, M., Doumbo, O. & Wellems, T. E. Pyrimethamine and
776         proguanil resistance-conferring mutations in Plasmodium falciparum dihydrofolate
777         reductase: polymerase chain reaction methods for surveillance in Africa. *Am J Trop Med*
778         *Hyg* **52**, 565-568, doi:10.4269/ajtmh.1995.52.565 (1995).
779  46  Demas, A. *et al.* Applied genomics: data mining reveals species-specific malaria
780         diagnostic targets more sensitive than 18S rRNA. *J Clin Microbiol* **49**, 2411-2418,
781         doi:10.1128/JCM.02603-10 (2011).
782  47  Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
783         sequence data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170 (2014).
784  48  Martin, M.  Vol. 17  (EMBnet, Administration Office, 2011).
785  49  Bushnell, B. BBMap short-read aligner, and other bioinformatics tool.  (2015).
786  50  Nelson, C. W. & Hughes, A. L. Within-host nucleotide diversity of virus populations:
787         insights from next-generation sequencing. *Infect Genet Evol* **30**, 1-7,
788         doi:10.1016/j.meegid.2014.11.026 (2015).
789  51  Wamae, P. M., Githeko, A. K., Otieno, G. O., Kabiru, E. W. & Duombia, S. O. Early biting
790         of the Anopheles gambiae s.s. and its challenges to vector control using insecticide
791         treated nets in western Kenya highlands. *Acta Trop* **150**, 136-142,
792         doi:10.1016/j.actatropica.2015.07.008 (2015).
793  52  Kulldorff, M., Huang, L. & Konty, K. A scan statistic for continuous data based on the
794         normal probability model. *Int J Health Geogr* **8**, 58, doi:10.1186/1476-072X-8-58 (2009).
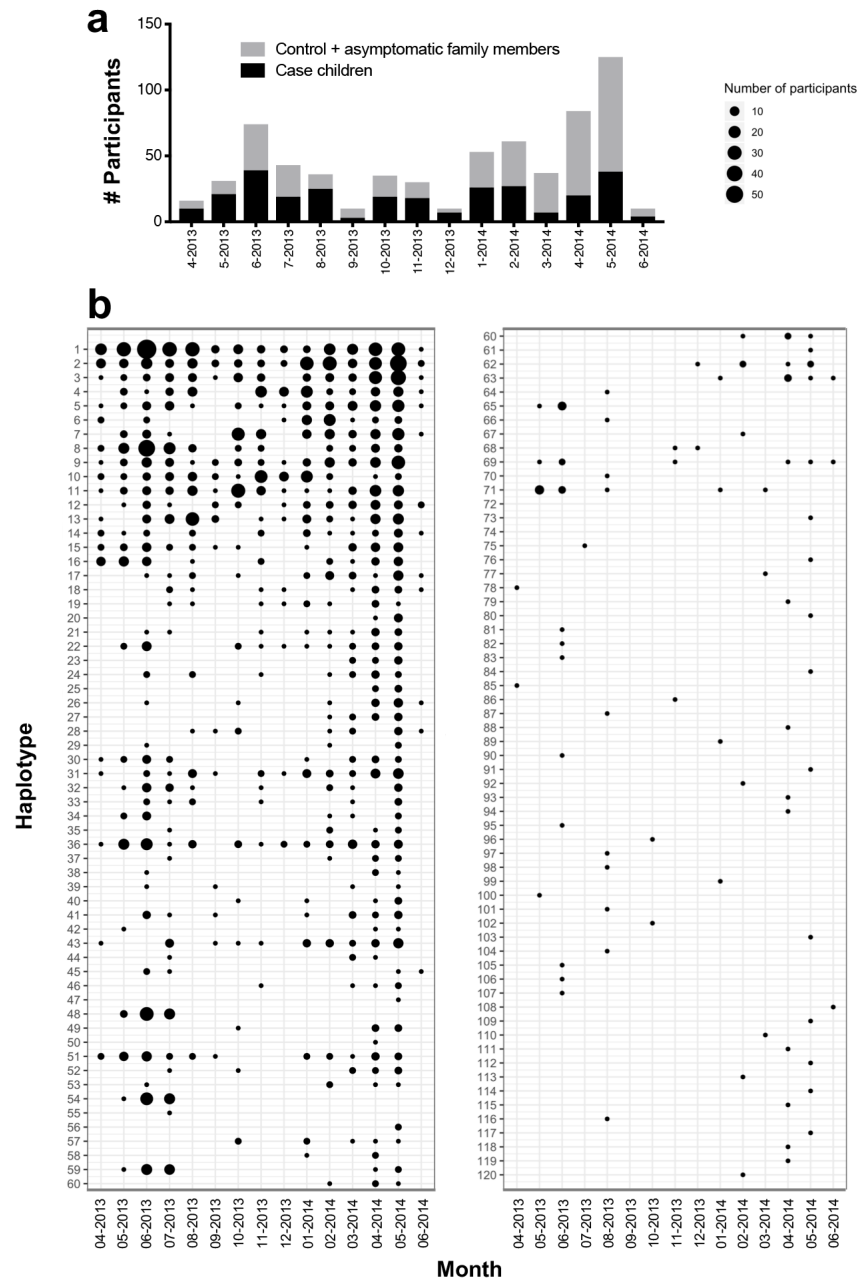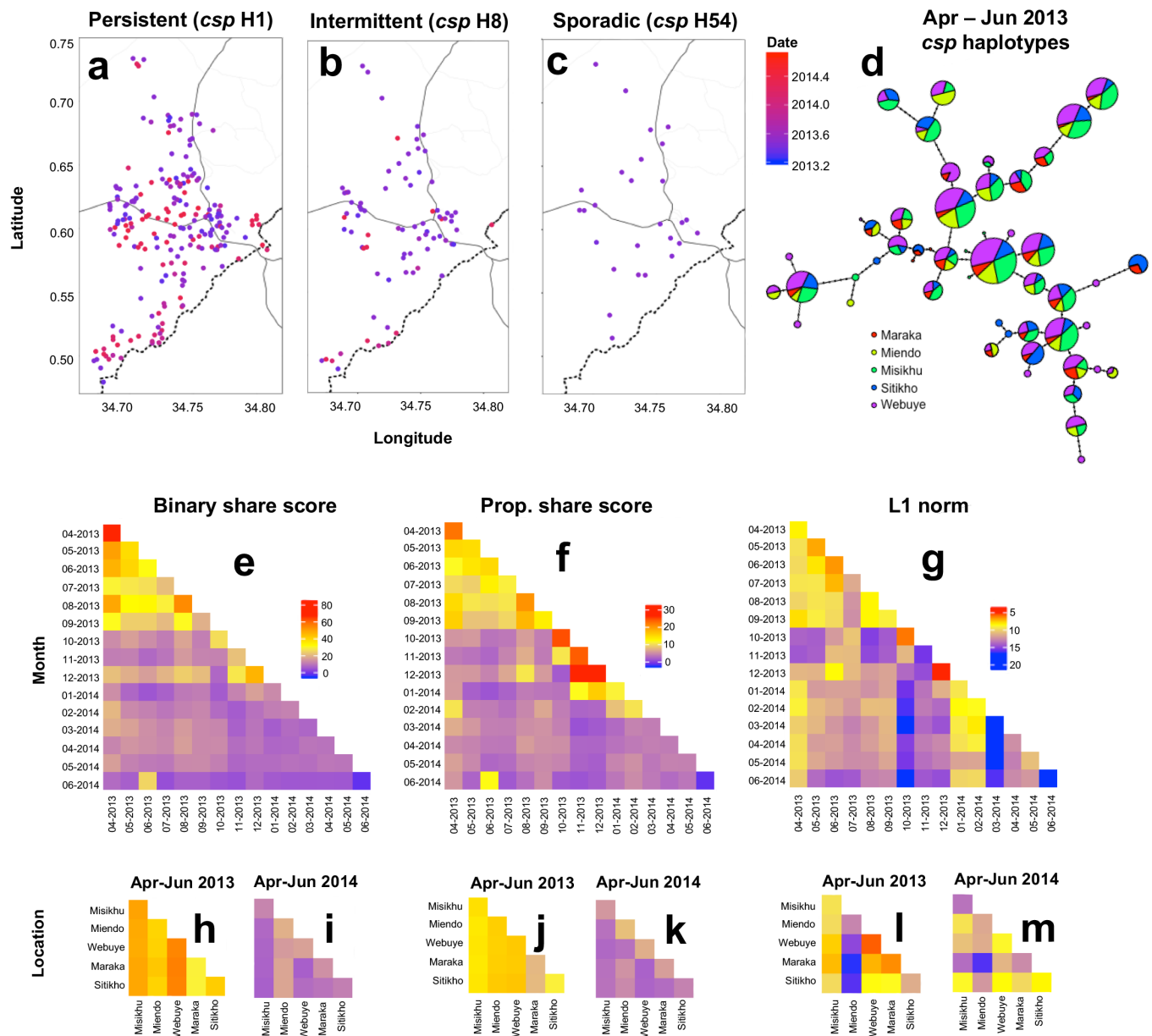
795

796    **FIGURES AND FIGURE LEGENDS**



797

798    **Fig. 1** Study setup, sequencing, and haplotype calling. All *P. falciparum* RDT+ subjects were

799    selected for PCR amplification and sequence analysis. Amplicons within both *csp* and *ama1*

800    hypervariable membrane proteins were sequenced. Approximately 30% of samples at both loci

801    failed haplotype assignment due to low number of sequencing reads. *csp* sequencing data was

802    excluded for 7 study participants and *ama1* data for 11 participants due to either data

803    inconsistencies or erroneous sample tracking/identification.
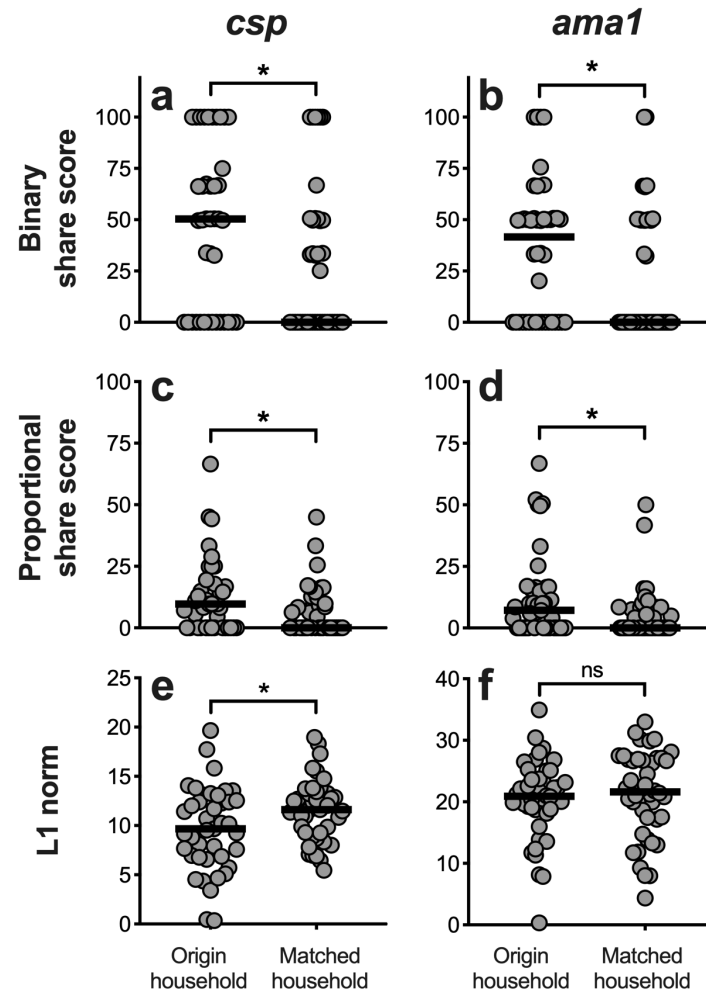
804

**Fig. 2** Heterogeneous persistence of *csp* haplotypes over time. **a** Total number of study participants with *csp* haplotypes by month. Black denotes case children, and gray indicates both control and case household members. **b** Monthly prevalence of 120 unique *csp* haplotypes, sorted by overall prevalence. Size of circle indicates number of study participants with a particular haplotype in a given month.
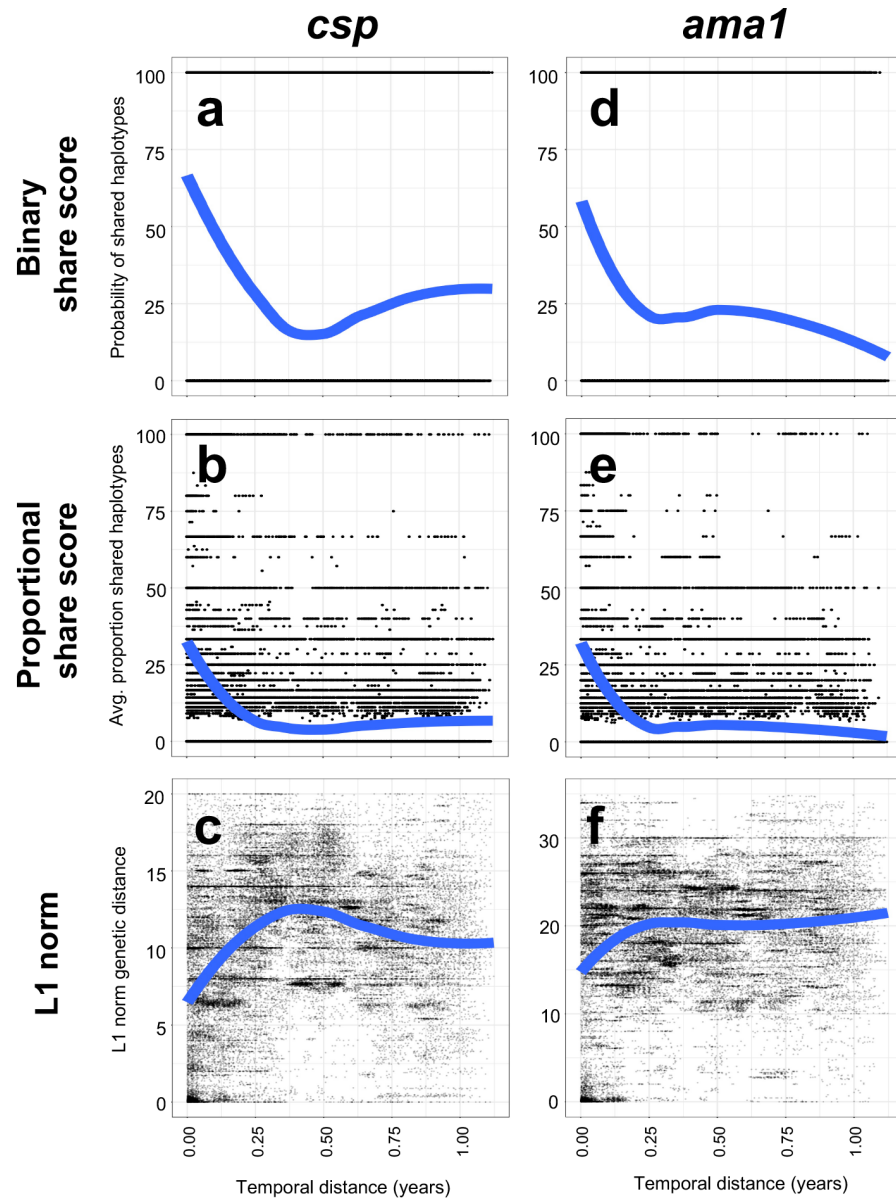
34

**Fig. 3** Genetic similarity of *csp* haplotypes is structured by time more than space. **a-c** Location of study participants with 'persistent' haplotype *csp* H1 (**a**), 'intermittent' haplotype *csp* H8 (**b**), and 'sporadic' haplotype *csp* H54 (**c**). Blue color indicates the beginning (April 2013) and red the end (June 2014) of the study period, with the date denoting fractional years in decimal notation. **d** *csp* haplotype network for the high transmission season of April through June 2013 (for 5 most

816    represented administrative locations: Maraka, Miendo, Misikhu, Sitikho, and Webuye). Each

817    circle indicates a unique haplotype with size proportional to the $\log_2$-scaled haplotype prevalence

818    and color denoting the fractional prevalence in each administrative location. Haplotype

819    connections were calculated using an infinite site model (Hamming distance) of DNA sequences,

820    with dots along connections indicating the number of base-pair differences between sequences.
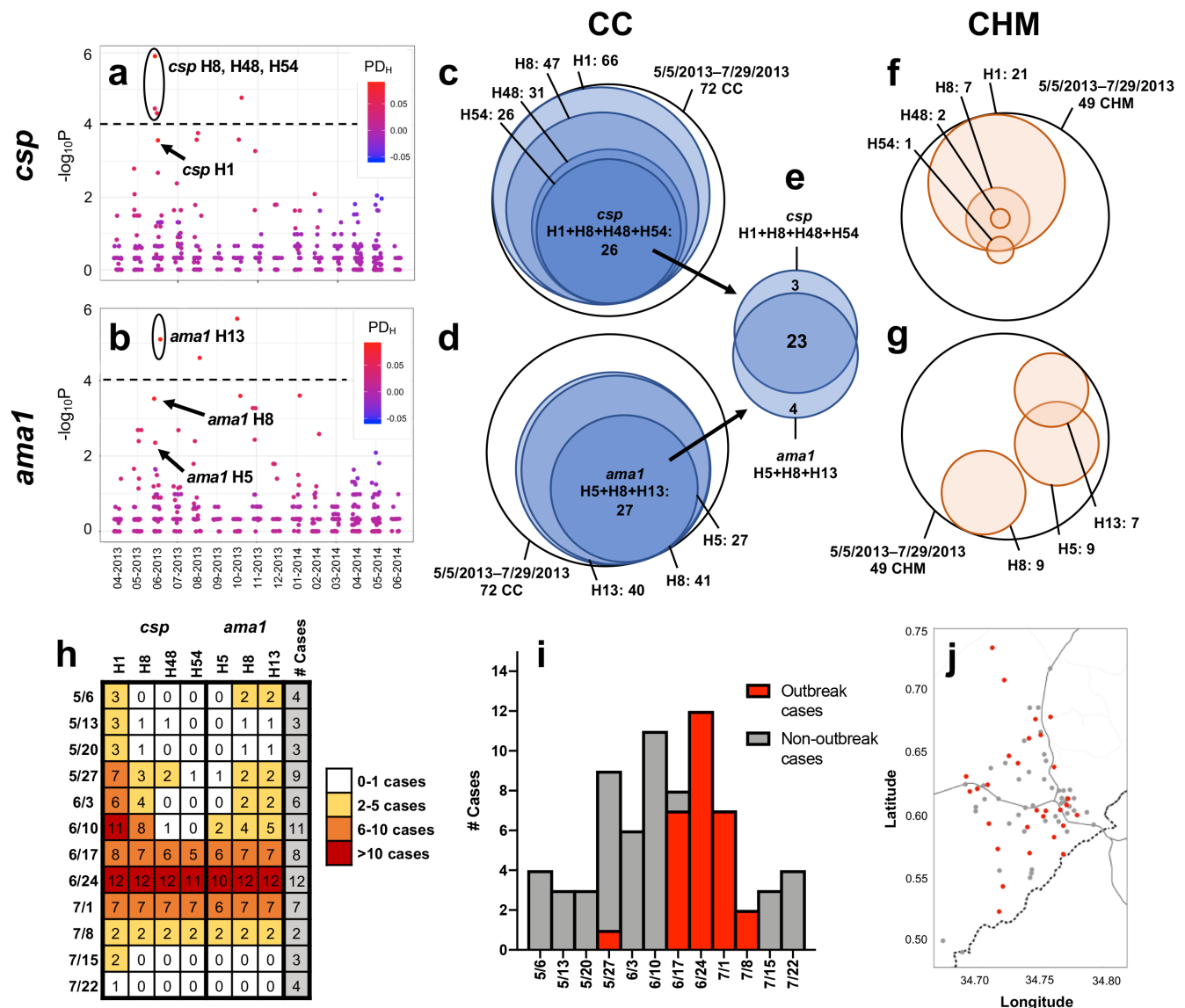
821    **e-g** Temporal comparison heat maps of mean binary haplotype sharing (**e**), proportional

822    haplotype sharing (**f**), and L1 norm genetic distance (**g**) calculated between months of study

823    enrollment. **h-m** Spatial comparison heat maps of binary haplotype sharing (**h,i**), proportional

824    haplotype sharing (**j,k**), and L1 norm genetic distance (**l,m**) calculated for a distinct temporal

825    window (**h,j,l**: April-June 2013; **l,k,m**: April-June 2014) for the 5 most represented administrative

826    locations (see above), which are arranged from north to south (see map in Figure S1). For **e-m**,

827    blue denotes the minimum for each genetic similarity index, red the maximum, and yellow the

828    midpoint.

**Fig. 4** Enhanced parasite genetic similarity between case children and their own household members. Genetic similarity metrics (binary share score, proportional share score, and L1 norm) were computed for case children (CC) with both their household of origin as well as with an unrelated household matched on time. Each metric was computed independently for *csp* and for *ama1*. **a,b** CC had significantly higher *csp* (**a**) and *ama1* (**b**) binary share scores with infected members of their origin households compared to members of matched households. **c,d** Furthermore, a higher proportional share score was observed for CC with their household of origin at both *csp* (**c**) and *ama1* (**d**) loci. **e,f** A reduced L1 norm for origin household (i.e. greater sequence similarity) was observed for *csp* (**e**) though not *ama1* (**f**) haplotypes. *p<0.05, Wilcoxon Signed-Rank test.

**Fig. 5** Case children with temporally-proximal infections have enhanced genetic similarity of *csp/ama1* haplotypes. Genetic similarity metrics were computed for all possible pairings of case children for *csp* (n=273) and *ama1* (n=288) haplotypes. **a-c** *csp* haplotype binary sharing (**a**), proportional sharing (**b**) and L1 norm (**c**) metrics for all CC pairwise comparisons is plotted against temporal distance for CC. **d-f** *ama1* haplotype binary sharing (**d**), proportional sharing (**e**), and L1 norm (**f**) metrics for all CC pairings is plotted against temporal distance for CC. Blue lines indicate the locally-estimated scatterplot smoothing (LOESS) regression fit of data.

**Fig. 6** Discovery of unique haplotype combination in case children indicating malaria outbreak.

**a,b** The haplotype prevalence difference ($PD_H$) between case children (CC) and case household members (CHM) during each month was calculated for *csp* (**a**) and *ama1* (**b**). Color indicates $PD_H$, with red identifying haplotypes more common in CC and blue more common in CHM. The y-axis indicates the Fisher's Exact test $-\log_{10}$(p-value) for haplotype prevalence in CC vs. CHM, while the dotted line denotes the Bonferroni-corrected threshold for statistical significance. **c-g** Venn diagrams of *csp* and *ama1* haplotypes with a high incidence in case children for 72 CC and

857    49 CHM between 5/6/2013 and 7/28/2013 (only CC with both *csp* and *ama1* haplotypes included).

858    **c,f** 26 CC had evidence of co-infection with *csp* haplotypes 1, 8, 48, and 54, compared to no

859    CHM. **d,g** 27 CC had evidence of co-infection with *ama1* haplotypes 5, 8, and 13, compared to

860    no CHM. **e** 23 CC were infected with CSP haplotypes 1, 8, 48, and 54 in addition to *ama1* 5, 8,

861    and 13. **h** Number of CC each week from 5/6/2013 – 7/29/2013 infected with each haplotype

862    associated with outbreak infection. White indicates 0-1 CC, yellow 2-5 CC, orange 6-10 CC, and

863    red >10 CC infected with a given haplotype in a particular week. **i** Outbreak epidemiology curve.

864    Weekly case incidence from 5/6/2013 – 7/28/2013 is shown by overall bar height. Weekly

865    incidence of outbreak cases (defined as presence of >=5/7 outbreak haplotypes CSP

866    H1/H8/H48/H54 and AMA H5/H8/H13, with these haplotypes comprising >98% reads detected)

867    is indicated in red, and non-outbreak cases in grey. **j** Household location of case children

868    (5/6/2013–7/28/2013) with outbreak cases (red) and non-outbreak cases (grey).

869 **TABLES**

870

871 **Table 1** Descriptive statistics for study participants who were RDT+ and with successful

872 *csp*/*ama1* haplotype assignment

873

| | | RDT+ (n=1,050) | Assigned *csp* haplotypes (n=662) | Assigned *ama1* haplotypes (n=678) |
|---|---|---|---|---|
| Person type | Case child | 43.1% | 43.6% | 43.2% |
| | Case household member | 39.0% | 38.4% | 39.3% |
| | Control household member | 17.0% | 17.9% | 17.4% |
| Median age (Range) | | 6 (0.08 – 82) | 6 (0.08 – 82) | 6 (0.08 – 82) |
| Median $\log_{10}$PD (Range) | | 2.16 (-0.86 – 6.67) | 2.36 (-0.51 – 6.67) | 2.31 (-0.86 – 6.67) |
| Median # haplotypes (Range) | | N/A | 2 (1 – 16) | 2 (1 – 14) |

874 **Table 2** Distribution of number of months unique *csp*/*ama1* haplotypes were detected

| Number of months present | Designation* | *csp* | | *ama1* | |
|:---:|:---:|:---:|:---:|:---:|:---:|
| | | Number haplotypes | % | Number haplotypes | % |
| 15 | | 2 | 1.67 | 0 | 0 |
| 14 | Persistent | 4 | 3.33 | 2 | 1.11 |
| 13 | | 2 | 1.67 | 3 | 1.67 |
| 12 | | 4 | 3.33 | 5 | 2.78 |
| 11 | | 5 | 4.17 | 8 | 4.44 |
| 10 | | 3 | 2.50 | 3 | 1.67 |
| 9 | | 1 | 0.83 | 6 | 3.33 |
| 8 | Intermittent | 6 | 5.00 | 4 | 2.22 |
| 7 | | 3 | 2.50 | 4 | 2.22 |
| 6 | | 3 | 2.50 | 6 | 3.33 |
| 5 | | 5 | 4.17 | 3 | 1.67 |
| 4 | | 10 | 8.33 | 10 | 5.56 |
| 3 | | 9 | 7.50 | 12 | 6.67 |
| 2 | Sporadic | 5 | 4.17 | 12 | 6.67 |
| 1 | | 58 | 48.33 | 102 | 56.67 |

875

876 *Designation based upon number of study months each haplotype was detected: persistent

877 >80% study months, intermittent 20% ≤ study months ≤ 80%, and sporadic <20% study months

42

**Table 3** Maximal genetic similarity metrics are highly predictive of case child temporal distance from origin household

| | | Binary share score | | Proportional share score | | L1 norm | |
|---|---|---|---|---|---|---|---|
| | | %* | p-value$^\dagger$ | %* | p-value$^\dagger$ | %* | p-value$^\dagger$ |
| **Household of origin** | | 18.4 | 0.06 | 15.7 | 0.11 | 10.5 | 0.20 |
| **Temporal distance from origin household (% maximal temporal distance)** | ± 10 days (5%) | 42.1 | **<0.01** | 42.1 | **<0.01** | 34.2 | **<0.01** |
| | ± 30 days (15%) | 50.0 | **0.03** | 50.0 | **0.03** | 55.3 | **<0.01** |
| | ± 60 days (30%) | 60.5 | **0.02** | 60.5 | **0.02** | 65.8 | **<0.01** |
| **Geographic distance from origin household (% maximal geographic distance)** | within 0.75 km (5%) | 18.4 | 0.15 | 21.1 | 0.09 | 13.2 | 0.43 |
| | within 2.25 km (15%) | 26.3 | 0.40 | 26.3 | 0.40 | 23.7 | 0.57 |
| | within 4.5 km (30%) | 36.8 | 0.81 | 36.8 | 0.81 | 39.4 | 0.63 |

*Percentage of case children (n=38) for which maximal binary share score, proportional share score, and L1 norm correctly predict case child household of origin OR temporal distance from origin household (within 3 temporal windows: ± 10, 30, or 60 days representing 5, 15, or 30% of the maximum time between cases) OR geographic distance from origin household (within 3 ranges – 5, 15, and 30% maximal geographic distance)

$^\dagger$ Fisher's Exact Test p-value based on assumption that cases are evenly distributed throughout time and across geographic space.