

1 Resolving within-host malaria parasite diversity using single-cell sequencing

2
3 Standwell C. Nkhoma^{1,2,3,4}, Simon G. Trevino⁴, Karla M. Gorena⁵, Shalini Nair^{4,1}, Stanley
4 Khoswe¹, Catherine Jett⁴, Roy Garcia⁴, Benjamin Daniel⁵, Aliou Dia⁴, Dianne J. Terlouw^{1,2},
5 Stephen A. Ward², Timothy J.C. Anderson⁴, Ian H. Cheeseman⁴

6
7 ¹Malawi-Liverpool-Wellcome Trust Clinical Research Programme, Chichri, Blantyre,
8 Malawi

9 ²Liverpool School of Tropical Medicine, Liverpool, United Kingdom

10 ³Wellcome Trust Liverpool Glasgow Centre for Global Health Research, Liverpool, United
11 Kingdom

12 ⁴Texas Biomedical Research Institute, San Antonio, Texas, United States of America

13 ⁵University of Texas Health Science Center San Antonio, San Antonio, Texas, United
14 States of America

15
16 Corresponding Authors: Ian Cheeseman, ianc@txbiomed.org, Standwell Nkhoma,
17 snkhoma@atcc.org

18
19
20 **Malaria patients can carry one or more clonal lineage of the parasite, *Plasmodium***
21 ***falciparum*, but the composition of these infections cannot be directly inferred from bulk**
22 **sequence data. Well-defined, complete haplotypes at single-cell resolution are ideal for**
23 **describing within-host population structure and unambiguously determining parasite**
24 **diversity, transmission dynamics and recent ancestry but have not been analyzed on a**
25 **large scale. We generated 485 near-complete single-cell genome sequences isolated**
26 **from fifteen *P. falciparum* patients from Chikhwawa, Malawi, an area of intense malaria**
27 **transmission. Matched single-cell and bulk genomic analyses revealed patients harbored**
28 **up to seventeen unique lineages. Estimation of parasite relatedness within patients**
29 **suggests superinfection by repeated mosquito bites is rarer than co-transmission of**
30 **parasites from a single mosquito. Our single-cell analysis indicates strong barriers to**
31 **establishment of new infections in malaria-infected patients and allows high resolution**
32 **dissection of intra-host variation in malaria parasites.**

33
34
35 Within a single host interactions between genetically distinct malaria parasites influence the
36 evolution of parasite virulence, antimalarial drug resistance, immunity, gametocyte sex ratios,
37 and malaria transmission in mouse malaria models¹⁻⁵. Complex infections (that contain more
38 than one unique parasite genetic background) confound most traditional genetic analysis,
39 preventing the accurate inference of allele frequencies and even simple phenotype
40 associations^{6,7}. Complex infections play a key role in the structure of populations as they
41 provide the substrate for sexual recombination to occur, in turn shaping local decay of linkage
42 disequilibrium and haplotype variation⁸⁻¹⁰. The impact of within-host interactions remain largely
43 unknown in human malaria. This is due to the paucity of appropriate tools for resolving infection
44 complexity on a large-scale at the level of single parasitized cells: we cannot directly infer the
45 composition of malaria infections by bulk sequencing of infected blood samples. Important
46 insights into the genetic architecture of individual malaria infections are emerging, aided by
47 recent advances in targeted capture of singly-infected erythrocytes from complex mixtures and
48 improved methods for single-cell sequencing^{7,11}, and computational approaches for interpreting
49 infection complexity^{12,13}.

50

51 Genetically distinct malaria parasites can
52 infect an individual through two routes
53 (Fig. 1). A single individual may be bitten
54 by two (or more) infected mosquitos, each
55 bearing a unique parasite genotype, or an
56 individual may be bitten by a single
57 mosquito bearing more than one parasite
58 genotype. Throughout, we refer to these
59 two processes as superinfection and co-
60 transmission respectively. Following a
61 bloodmeal, gametocyte stage parasites
62 fuse in the mosquito midgut, and an
63 obligate round of sexual recombination
64 occurs. If only a single parasite genotype
65 is present all offspring will be identical
66 (Fig. 1, top panel), with the presence of
67 multiple genotypes allowing recombinant
68 progeny to arise^{6,8,14-16} (Fig. 1, bottom
69 panel). Using single cell sequencing and
70 cloning by limiting dilution of parasites
71 from a single individual we have
72 previously seen a range of inferred
73 relationships, including identical clones,
74 siblings and unrelated individuals^{7,11,16}.
75 However, it is unknown to what extent
76 these findings can be generalized across
77 a population. The degree to which the
78 genetic diversity of individual infections is
79 driven by superinfection of unrelated
80 strains, or co-transmission of related ones
81 is needed to model how genetic diversity
82 could be maintained in the face of malaria
83 control measures.

85 Infection complexity in bulk sequenced 86 samples

87 To resolve the within-host structure of
88 malaria infections, we performed a cross-
89 sectional survey of individuals infected
90 with uncomplicated *P. falciparum* malaria
91 in Chikhwawa, Malawi, an area of high malaria transmission (entomological inoculation rate 183
92 infectious bites per person per year¹⁷). We performed bulk parasite genome sequencing of 49
93 infections to a median read depth of 31 (interquartile range 20.93-48.37). We estimated the
94 complexity of infection of bulk sequence data using 10,997 unfixed SNP positions with a minor
95 allele frequency (MAF) >0.05 using the F_{WS} statistic^{18,19} and DEploid¹³ (Fig. 2a,b, Supplementary
96 Table 1). F_{WS} grades infections on a continuous scale of complexity where infections with an
97 F_{WS} >0.95 are considered clonal and DEploid estimates the number of haplotypes (K) present in
98 sequence data by jointly estimating haplotypes and their abundances. In close agreement with
99 contemporary estimates of within host diversity²⁰, 22 of 49 infections (44.9%) were considered
100 clonal by F_{WS} . The within-host allele frequency (WHAF) captured from deep sequencing can be
101 used to infer the presence of related parasites²¹. The patterns of unfixed mutations in the

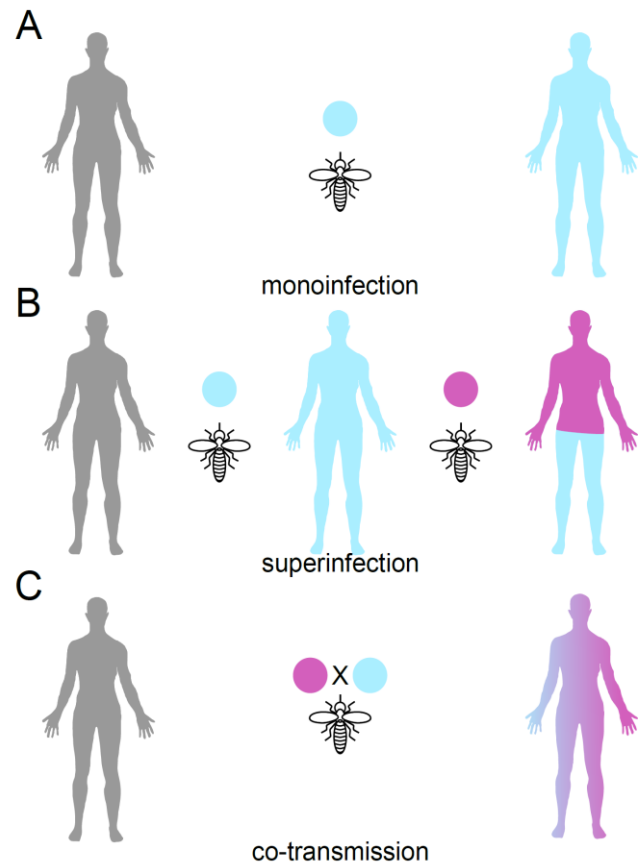


Figure 1. The within host genetic diversity of malaria parasites is shaped by transmission strategy.

(A) A simple monoinfection is generated when an uninfected individual is bitten by a mosquito bearing a single parasite genotype. (B) A superinfection occurs when an individual is bitten by two mosquitos, each bearing a single parasite genotype. (C) Co- transmission of parasites occurs when a single mosquito bearing multiple genetically distinct parasites bites an uninfected individual. As genetic recombination is an obligate stage of mosquito transmission multiple related parasites may be infected through this route.

102 remaining 27 infections suggest a simple model of superinfection, where two unrelated parasite
103 genetic backgrounds colonize an individual, are insufficient to universally capture all patterns of
104 within-host relatedness (Supplementary File 1). We selected 15 infections across the range of
105 F_{WS} and inferred K for single-cell sequencing, using a recently optimized method to generate
106 near-complete genome capture¹¹. The malaria parasite undergoes 4-5 rounds of DNA
107 replication within a single cell producing segmented schizont stage parasites with an average of
108 16 genome copies²². We isolate individual schizonts by fluorescence activated cell sorting,
109 followed by whole genome amplification (WGA) under highly sterile conditions before
110 sequencing the amplified product.

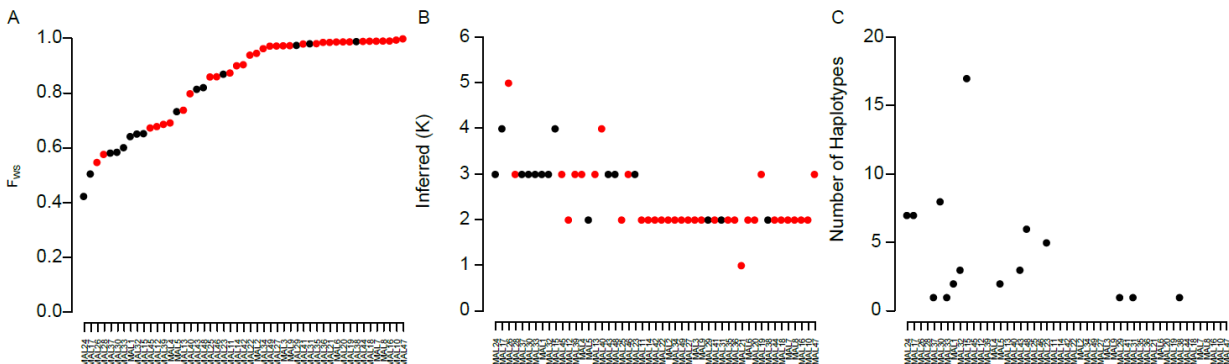


Figure 2. Complexity of infection inferred from bulk and single-cell sequencing. (A) F_{WS} scores for 49 bulk sequenced infections. Infections above the dashed line ($F_{WS}=0.95$) are assumed to be clonal. (B) Inferred number of haplotypes (K) inferred by DEploid, infections are ordered by the F_{WS} score. Black dots in (A) and (B) denote infections also deconvoluted by single-cell sequencing. (C) Number of unique haplotypes inferred by single-cell sequencing.

111

112 Single cell sequencing of malaria parasites

113 In total we sequenced the genomes of 485 single-cells subjected to WGA (437 unique to this
114 study), 49 bulk infections and 24 clones isolated from a single patient by limiting dilution^{16,23}.
115 Prior to genotype filtering we scored 175,543 biallelic SNPs with a VQSLOD>0 across the 558
116 genome sequences. The highly repetitive and AT-rich *P. falciparum* genome²⁴ presents unique
117 challenges with generating an accurate picture of the variation present in a single-cell. We were
118 particularly concerned with capture of DNA from more than one genetic background during the
119 single-cell sequencing protocol and implemented stringent quality checks. Using sequencing
120 data from the 24 clones we estimated the threshold for identifying single cell sequences where
121 there was potential contamination from exogenous DNA at 1% of mixed base calls. The
122 sequences from the cloned lines were integrated into the single cell dataset for downstream
123 analysis. After excluding low coverage libraries (<75,000 calls, n=23) and sequences with >1%
124 mixed base calls (n=38) 424 single-cell sequences remained. After including 23 of the
125 sequences from *ex vivo* expanded clones there were 13-45 sequences per infection (mean 29.9
126 sequences; Supplementary Data Fig. 1). The number of sequences per sample attempted was
127 determined by rarefaction analysis (described below).

128

129 After quality control of the dataset we retained 60,002 SNPs scored in at least 90% of the 496
130 sequences, 10,997 of which had a MAF>0.05 across the 49 bulk sequenced infections. As an
131 initial characterization of our data we estimated the genetic diversity in each infection from the
132 number of unfixed sites from read pileups in bulk sequencing or across called genotypes in
133 single-cell sequencing. For paired bulk/single-cell data from the same infection a mean of 1.6
134 fold (range 0.7-9.1 fold) more polymorphic sites were discovered by single-cell sequencing than
135 by bulk sequencing (Supplementary Data Fig. 2). This is likely due to the limits in discovery of

136 very low frequency SNPs by bulk sequencing. By subsampling our single-cell data we saw
137 diminishing returns from sequencing additional cells, with 90% of the observed polymorphic
138 sites captured by sampling a mean of 21.6 cells (range 7-43, Supplementary Data Fig. 2).

139

140 **Haplotypic diversity of malaria infections**

141 A major goal in malaria genomics has been estimating the number of unique haplotypes (or
142 complexity of infection) within an infection²⁵. We estimated the number of unique haplotypes
143 directly from the single-cell data. To exclude potential confounding of *de novo* mutation and
144 sequencing error we restricted analysis to 10,997 conservatively called sites with a MAF >0.05
145 in the 49 bulk sequenced infections. We estimated the number of unique haplotypes per
146 infection by collapsing haplotypes from the same infection that were different at <1% of sites.
147 For each infection we applied individual-based rarefaction to the haplotype abundances and
148 sequenced additional single genomes until a plateau in the rarefaction curve was reached
149 (Supplementary Data Fig. 3). Using this approach between 1 and 17 haplotypes were observed
150 in each infection (Fig. 1c, Supplementary Data Table 1). There was strong correlation between
151 the effective number of strains¹³ inferred by single cell sequencing and the effective K from
152 DEploid (Pearson's $r^2=0.61$) and F_{WS} (Pearson's $r^2=-0.51$, Supplementary Data Fig. 4).
153 Rarefaction of haplotype abundance suggested exhaustive capture of haplotypes in 12/15
154 infections. In two infections (MAL23 and MAL30) we sampled two fewer haplotypes than
155 suggested by rarefaction (Chao I estimator- MAL23=6.94, MAL30=10.18, observed haplotypes-
156 MAL23=5, MAL30=8). In both cases the observed number of haplotypes were within the 95%
157 confidence intervals of the estimation. One infection (MAL15) showed exceptionally high
158 diversity with 17 of an estimated 30.21 (95% CI=19.7-81.7) haplotypes detected. Two infections
159 (MAL37 and MAL33) show a single haplotype from single-cell sequencing, although F_{WS} scores
160 <0.95 and patterns of segregating sites suggest we have incompletely captured all haplotypes
161 (Supplementary File 1). Sequencing more cells did not capture additional haplotypes.

162

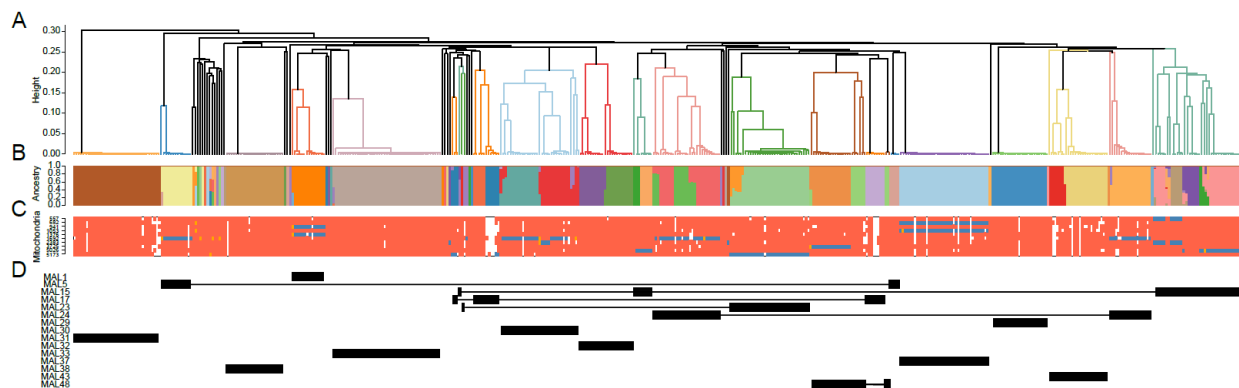


Figure 3. Clustering of single-cell and bulk genome sequences. (A) A UPGMA tree of 1-pairwise allele sharing across all samples passing quality control. Shading behind the tree and labels is specific to a particular infection. (B) Unsupervised clustering of parasites by ADMIXTURE, each bar denotes the proportion of ancestry a parasite derives from 31 latent populations. (C) Mitochondrial haplotypes for each parasite. Red blocks denote the reference allele, blue the alternative and orange where mixed base calls were observed. The ADMIXTURE proportions and mitochondrial haplotypes are oriented to be below the branch tip of the same parasite. (D) The location of cells from each infection subject to single-cell sequencing in the upper panel.

163 **Population structure of individual infections**

164 The number of unique haplotypes alone captures only a single aspect of the recent history of
165 genetically distinct parasites from the same infection. For instance it does not distinguish
166 whether diversity is due to superinfection of unrelated strains from multiple mosquito

167 inoculations, from co-transmission of related strains from a single mosquito inoculation or from a
168 combination of the two. To jointly characterize the between- and within-host genetic structure of
169 malaria infections we clustered infections based upon either a UPGMA tree of pairwise allele
170 sharing (Fig. 3a) or by unsupervised clustering in ADMIXTURE²⁶ (K=31; Fig. 3b, Supplementary
171 Data Figure 5). In both cases sequences from the same infection predominantly cluster
172 together, suggesting they were more related to each other than parasites drawn from the
173 population and diversity is likely the result of co-transmission by related parasites from the same
174 mosquito.

175
176 We used clustering estimated by the UPGMA tree, ADMIXTURE and mitochondrial genotypes
177 (Fig. 3c) to distinguish between infections where diversity results from superinfection or
178 coinfection. Mitochondrial genome sequences are useful markers in this context as they are
179 uniparentally inherited in malaria infections and do not undergo recombination. However, there
180 are few confidently scored mutations in our mitochondrial genome data (n=10) limiting the
181 resolution of our inference. To classify potential superinfections we first identified infections
182 which contained putatively unrelated parasites. Based upon the UPGMA tree 6 infections
183 (MAL5, MAL15, MAL17, MAL23, MAL24, MAL48, Fig. 3d) contained sequences that cluster
184 more closely to other infections than to sequences from the same infection. For 3 of these
185 infections (MAL5, MAL23, MAL24) ADMIXTURE results showed clustering which was congruent
186 with the UPGMA tree. For example, in MAL5 there were 2 distinct clusters by both ADMIXTURE
187 and the UPGMA tree which were in agreement and both of these clusters had a unique
188 mitochondrial genotype. The remaining 3 infections (MAL15, MAL17, MAL48) each showed
189 discordant clustering between the methods. For instance MAL15 shares both ADMIXTURE
190 clusters and mitochondrial genotypes between parasites which were separated by the tree.
191 Based upon this analysis MAL5, MAL23 and MAL24 show evidence of superinfection, while the
192 diversity present MAL15, MAL17 and
193 MAL48 can be explained by co-
194 transmission alone.

195 196 **Recent ancestry of individual 197 infections**

198 To better characterize levels of
199 relatedness within infections we
200 identified blocks of chromosomes
201 shared identical-by-descent (IBD)
202 between all paired sequences using a
203 hidden Markov model²⁷. IBD sharing
204 between clonal bulk sequenced
205 infections was rare, with a mean of
206 0.73 blocks shared between infections
207 (range 0-5), encompassing a mean of
208 88.5kb (range 3.8-342.7kb) of each
209 genome, with a mean block length of
210 50.8kb (range 3.8-142.4kb). In
211 contrast, within infections parasites
212 shared a mean of 13.0 (range 0-30)
213 IBD blocks between parasite genomes,
214 encompassing a mean of 16,334.2kb
215 (range 3.1-20,577.0kb) of each
216 genome, with a mean shared block
217 length of 1,143.6kb (range 3.1-

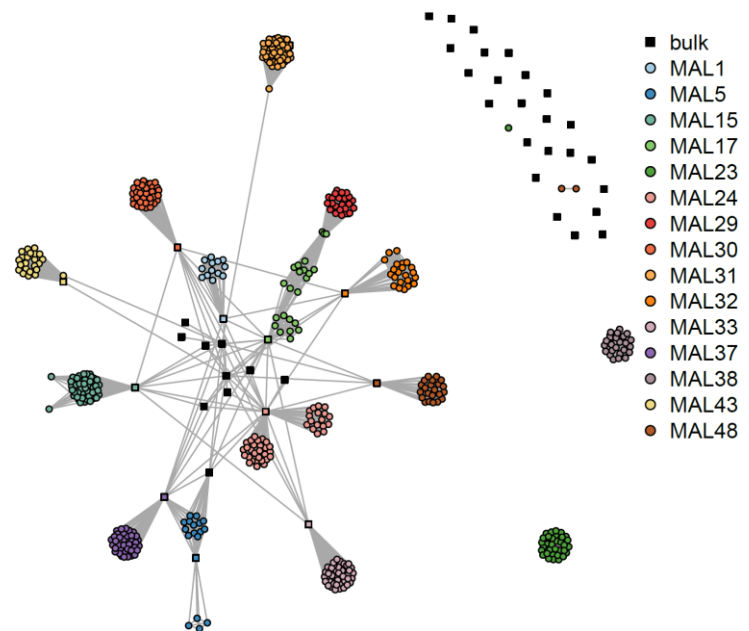


Figure 4. A network representation of pairwise IBD sharing across the genomes. Each node represents a single parasite colored by the infections of origin. Nodes are joined parasites if >15% of the genomes were shared IBD. Each node is colored by the infection it was derived from, with bulk sequences denoted by a square and single cell sequences by a circle.

218 1469.8kb). As we limit inference of IBD to the ‘core’ genome²⁸ identical parasites share
 219 20,577kb of their genomes IBD in 14 blocks (one per chromosome). The presence of IBD
 220 sharing between individuals supports recent shared ancestry. For 10/15 of the infections there
 221 was at least one block of IBD shared in all pairwise comparisons. As our filtering of IBD blocks
 222 was limited to >2.5cM we are limited to inference of relatedness over the last 25 generations (~6
 223 years)²⁹.

224
 225 Recent studies have highlighted the power of IBD networks to capture the structure of a parasite
 226 population²⁹. We built a network of pairwise shared IBD, creating links between parasites with
 227 >15% of their genomes shared IBD (Fig. 4, Supplementary File 3). This revealed close
 228 connectivity between parasites from the same infection, with much sparser connectivity between
 229 parasites from different infections. We observed subdivision within individual infections,
 230 supporting many of the observations in Fig. 2. Varying the minimum IBD required to connect
 231 genomes allowed us to visualize how relatedness subdivides individual infections
 232 (Supplementary File 3). For instance MAL5 and MAL24 where two clusters of parasites were
 233 only connected by the sequence derived from bulk sequencing to one another.

234
 235 The distribution of total pairwise shared IBD and the average shared block lengths
 236 can be used to infer the relationships
 237 between individual genomes^{30,31}. We
 238 inferred the degree of relatedness from
 239 our data using the Estimation of Recent
 240 Shared Ancestry (ERSA) algorithm. ERSA
 241 estimates relatedness between individuals
 242 from distribution of IBD tract lengths (Fig.
 243 5a,b) using assumed unrelated individuals
 244 from the same population as a reference.
 245 We see a spectrum of relationships within
 246 each infection (Fig. 5c, Supplementary
 247 File 2). In MAL5 this confirmed the lack of
 248 relatedness between the two clusters of
 249 parasites, suggesting this infection was
 250 the result of a genuine superinfection.
 251 However, no other infections can be
 252 classified so simply, commonly showing
 253 relationships as distant as 4th degree
 254 (equivalent to ‘first cousins’). Within our
 255 data this suggests that it is not uncommon
 256 for parasites to be transmitted through two
 257 generations (human-mosquito-human-
 258 mosquito-human), with up to four
 259 generations of co-transmission seen in our
 260 data in infection MAL24 and MAL17. In
 261 our data we see only a single
 262 unambiguous instance of superinfection of
 263 two unrelated parasites with no concurrent
 264 co-transmission (MAL5). Across the
 265 analysis, the genetic diversity of three infections (MAL17, MAL24 and MAL48) appears to be
 266 driven by both superinfection of unrelated parasites and co-transmission of related parasites (in
 267 addition to MAL5 where only superinfection is suspected). Surprisingly, we see substantial
 268

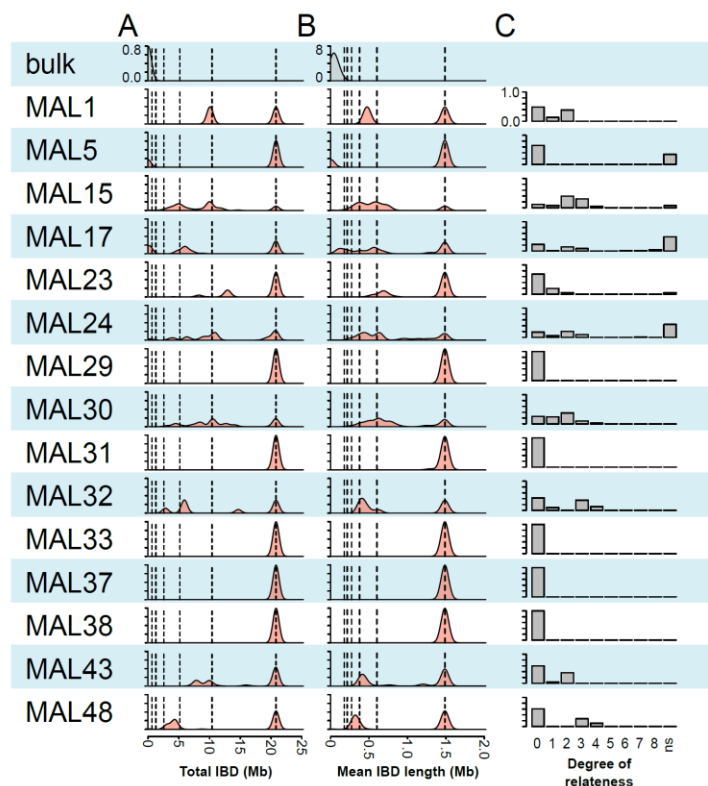


Figure 5. Recent ancestry inferred from IBD sharing.

(A) Density plot of the total IBD shared between parasites from a single infection (labeled to the left of the plot). (B) Density plot of the mean IBD block length between parasites from a single infection. (C) Relative frequency of different degrees of relatedness inferred between parasites from the same infection (ns - no significant relatedness observed).

across the analysis, the genetic diversity of three infections (MAL17, MAL24 and MAL48) appears to be driven by both superinfection of unrelated parasites and co-transmission of related parasites (in addition to MAL5 where only superinfection is suspected). Surprisingly, we see substantial

269 genetic variation amongst co-transmission parasites. In each of the three infections
270 there were more polymorphic sites segregating between co-infecting parasites than separating
271 superinfecting lineages (Supplementary Data Fig. 6).

272

273 Discussion

274 There has been a concerted effort to understand the complexity of malaria infections from either
275 deep sequencing data^{13,32,33}, or from genotyping a limited number of markers^{12,34}. We show here
276 there is considerable depth to complex infections which may be challenging to infer from bulk
277 analysis alone. Through a combination of deep sequencing of bulk infections and single-cell
278 sequencing we have generated the most comprehensive study of the within-host diversity of
279 malaria infections to date. This provides a much needed standard for developing novel tools for
280 probing the complexity of infections from deep sequencing data. By using multiple estimates of
281 relatedness targeting distinct features of the data we argue that most complex infections result
282 from parasites co-transmitted from single mosquito bites in our dataset. Strikingly, our analysis
283 supports only a single infection where simple superinfection of two unrelated strains has
284 occurred (MAL5), and a further three infections where both superinfection and co-transmission
285 have concurrently contributed to diversity (MAL17, MAL24, MAL48). The remaining infections
286 were either monomorphic or showed strong support for co-transmission of related strains only.
287 In the two infections where we were unable to capture the minor strains (MAL33 and MAL37)
288 patterns of unfixed SNPs within the infection suggest the uncaptured strain was related to the
289 captured strain (Supplementary File 1).

290

291 Only parasites which transmit gametes to the same mosquito can produce recombinant
292 offspring. Estimates of the parasite diversity and relatedness within individual mosquitos³⁵
293 (albeit in a distinct population) are in general agreement with our data – most mating is between
294 related parasites. The mechanisms underlying why inbreeding is common, even in high
295 transmission settings, is less clear. Malaria transmission is intense in Chikhwawa¹⁷ and we
296 expected superinfection to be more prevalent than we observed. A mechanism controlling the
297 outcome of superinfection, perhaps by hepcidin based inhibition of liver development in
298 superinfecting sporozoites³⁶, could explain why we do not see more superinfection. Alternatively,
299 the low numbers of superinfecting parasites emerging from the liver relative to those present in
300 established infections (which may contain 10^{11-12} blood stage parasites) may limit establishment
301 of superinfections. Analysis of parasite diversity is generally limited to single blood draws due to
302 a need to treat symptomatic patients expediently. As this sampling strategy may overlook sub-
303 clones circulating at lower frequencies there may be additional genetic variation which escapes
304 routine analysis.

305

306 The depletion of genetic variation during repeated rounds of co-transmission has been
307 previously modelled⁶, suggesting a substantial decline in the number of clones and an increase
308 in average relatedness can arise through a single transmission cycle. Our data suggest that few
309 complex infections have parasites which have been co-transmitted longer than two transmission
310 cycles. We observe substantial genetic variation is maintained despite the bottleneck of
311 mosquito transmission (Supplementary Data Fig. 6) with up to 17 unique haplotypes likely
312 inoculated by a single mosquito. Understanding how patterns of transmission and within host
313 dynamics contribute to the diversity and relatedness structure within malaria infections will be
314 critical to ongoing elimination and control efforts.

315

316 Acknowledgements

317 We thank all children who participated in this study in Chikhwawa, Malawi. We also thank
318 Andrew Mtande, Ruth Daiman and Miriam Phiri for their help with participant recruitment and
319 clinical management. We are also grateful to Clement Masesa and Lumbani Makhaza for

320 designing our data collection tool and for managing the study database. This study was
321 supported by a Wellcome Trust Intermediate Fellowship in Tropical Medicine and Public Health
322 (Grant # 099992/Z/12/Z to SCN) and an NIH grant (NIAID AI110941-01A1 to IHC). FACS data
323 were generated in the Flow Cytometry Shared Resource Facility (supported by UTHSCSA, NIH-
324 NCI P30 CA054174, UL1RR025767 (CTSA)).

325

326 **Author Contributions**

327 S.C.N., S.G.T., R.S.H, S.A.W., D.J.T, T.J.C.A and I.H.C designed the study. S.G.T. and I.H.C.
328 developed tools. S.C.N., S.G.T., K.G., S.N., A.D., C.J., R.G., B.D., and I.H.C. performed
329 experiments. S.C.N., S.K., and D.J.T. collected samples. S.C.N, S.G.T., T.J.C.A and I.H.C.
330 wrote the paper.

331

332 **Author Information**

333 The authors declare no competing interests. Correspondence and requests for materials should
334 be addressed to ianc@txbiomed.org. Raw sequence data has been deposited at the sequence
335 read archive (<https://www.ncbi.nlm.nih.gov/sra>) under study number SRP155167.

336

337 **Methods**

338

339 **Sample Collection**

340 Malaria-infected blood samples (5 ml; thin smear parasitaemia: 0.2 to 21.8%) were obtained
341 prior to treatment from children aged 19 to 116 months old presenting to Chikhwawa District
342 Hospital in Malawi with uncomplicated *P. falciparum* malaria from February to July 2016. Blood
343 samples were collected in Acid Citrate Dextrose tubes (BD, UK) following consent from parents
344 or guardians, and transported in an ice-cold container to our laboratory in Blantyre, Malawi, for
345 processing. Half of each blood sample was washed using incomplete RPMI 1640 media
346 (Sigma-Aldrich, UK) and cryopreserved in glycerolyte 57 solution (Fenwal, Lake Zurich, IL,
347 USA). Parasites used in fluorescence-activated single-cell sorting were cultured from this
348 sample. The second half of the sample was filtered using CF11 columns to deplete human
349 leucocytes³⁷, and was stored at -80°C until needed. Parasite DNA was extracted from this
350 sample using a DNA Mini Kit (QIAGEN, USA) and directly sequenced on an Illumina HiSeq
351 instrument. Ethical approval for this study was obtained from the University of Malawi College of
352 Medicine and Ethics Committee (Protocol number P.02/13/1528) and the Liverpool School of
353 Tropical Medicine Research Ethics Committee (Protocol number 14.035).

354

355 **Cell culture and FACS sorting**

356 Approximately 1 mL of frozen sample was thawed at 37°C and parasites were revived (~200ul
357 recovered pellet, ~1% parasitemia). Half of the recovered sample was frozen for bulk DNA
358 extraction and analysis. The other half was grown in 8 mL complete media for 40 hours to allow
359 for parasite progression to late stages, which generates higher quality genomic data after MDA
360 and library preparation¹¹. ~8 ul of infected red blood cell pellet was stained in 10 mL PBS which
361 included 5 ul of Vibrant DyeCycle Green at 37C with intermittent mixing for 30 minutes. Cells
362 were washed once in PBS and individually sorted by FACS, gating for trophozoite and schizont-
363 stage parasites.

364

365 **Single-cell Sequencing**

366 Library preparation for individually sorted late-stage parasites was carried out using the Qiagen
367 Single-Cell FX DNA kit without library amplification according to manufacturer's instructions.
368 Library products were analyzed by TapeStation and included off-target peaks typical of MDA
369 DNA inputs. Adapter-ligated DNA products were quantified by KAPA Hyperplus Kits. All
370 sequencing was performed on an Illumina HiSeq 2500.

371

372 **Sequence analysis**

373 Median read depth of WGA single-cells was 28.3 (interquartile range (IQR) 12.5-46.4) with
374 median of 90.5% (IQR 78.1-96.0%) of the genome covered by at least one read. In contrast the
375 non-WGA samples had a median read depth of 31.11 (IQR 20.93-48.37) and a median of
376 95.8% (IQR 93.1-97.4%) of the genome covered by at least one read. A potential source of
377 error in single-cell genomics is the inclusion of exogenous DNA amplified alongside the target
378 genome in downstream analysis. As an initial indication of the potential of non-target DNA being
379 introduced to our analysis we first examined the proportion of reads mapping to the *P.*
380 *falciparum* genome²⁴ in each sequence. We observed a median of 93.3% (IQR 87.0-95.4%) of
381 reads map to the parasite genome for single-cell sequences, compared to 35.7% (IQR 19.7-
382 48.5%) for bulk patient samples and 79.4% (IQR 74.5-86.9%) for clonally expanded samples
383 suggesting our stringent handling protocols were effective at eliminating environmental DNA.
384 For a more rigorous test we identified lines with potential cross contamination based on unfixed
385 basecall frequency. As the parasite genome is haploid during blood stages all variants are
386 expected to be fixed in genome sequencing data. The highly AT-rich and repetitive nature of the
387 parasite genome makes alignment challenging, generating false positive unfixed variants in
388 clonal lines. After excluding highly error-prone genomic regions (calls outside of the “core
389 genome”²⁸ or within microsatellites) we measured the proportion of mixed base calls (>5% of
390 reads at a locus mapping to the minority allele) at high confidence biallelic SNPs (>10 reads
391 mapped, VQSLOD>0, GQ>70). Using the cloned lines and bulk population samples as a guide
392 we estimated 1% as an appropriate threshold for excluding putatively mixed lines
393 (Supplementary Data Fig. 1).

394

395 **Estimating the complexity and diversity of bulk sequenced samples**

396 F_{WS} was calculated in moimix (<https://github.com/bahlolab/moimix>) for all bulk patient samples.
397 We estimated the number of unique haplotypes and their sequence from deep sequence of bulk
398 infections using DEploid¹³ v0.5 (<https://github.com/mcveanlab/DEploid>). We used 10,997 HQ
399 SNPs with a MAF >5%. For a reference panel we used 10 bulk Malawian samples presumed to
400 be clonal ($F_{WS}>0.95$) and population level allele frequencies from across the complete bulk
401 sequencing data. We inferred the most likely number of haplotypes (K) using the command:

402

```
403 ./dEplod -ref sample_reference_allele_counts.txt -alt sample_alternative_allele_counts.txt -plaf  
404 population_allele_freq.txt -o sample_out -ibd -noPanel -exclude highly_variable_sites.txt -sigma  
405 7 -seed 2
```

406

407 **Estimating relatedness between sequences**

408 SNP data were imported into R using SeqArray³⁸. Between all samples passing quality control
409 we calculated the proportion of shared alleles and using SNPs which were at >5% MAF in the
410 bulk sequenced samples. We used a distance matrix generated from this data (1-pairwise allele
411 sharing) to build a UPGMA tree (Fig. 2). We also used this statistic to estimate the number of
412 unique haplotypes in each infection by collapsing together sequences which differed at <1% of
413 sites. Rarefaction of haplotype abundance was performed using the rareNMtests package³⁹ in
414 R. We performed unsupervised clustering of the sequence data using ADMIXTURE v1.3²⁶
415 (<https://www.genetics.ucla.edu/software/admixture/>). This again used sites with a MAF of >5%
416 across the bulk sequenced data. We clustered data using K of 2-40 seeing a minima of CV error
417 at K=31 (Supplementary Data Fig. 5). We called regions of IBD between all samples passing
418 quality control using hmIBD v2.0.0²⁷ (<https://github.com/glipsnort/hmIBD>). We performed
419 maximum-likelihood estimation of recent shared ancestry using ERSa 2.0^{30,31}
420 (<http://www.hufflab.org/software/ersa/>) using the output from hmIBD using the flags --
421 min_cm=1.5 --adjust_pop_dist=true --number_of_chromosomes=14 --rec_per_meioses=19. We

422 converted the basepair positions to a uniform genetic map using the scaling factor $1\text{cM}=9.6\text{kb}^{40}$
423 and excluded IBD chunks $<1\text{cM}$ in length. As identical clones are not specifically modelled in
424 ERSA we excluded these from analysis, though their abundance is shown in the '0' bar in Fig.
425 3c. All other statistical analysis and visualization was performed in R v3.4.0⁴¹.

426

427 References

- 428 1 Bell, A. S., de Roode, J. C., Sim, D. & Read, A. F. Within-host competition in genetically diverse
429 malaria infections: parasite virulence and competitive success. *Evolution* **60**, 1358-1371 (2006).
- 430 2 de Roode, J. C. *et al.* Virulence and competitive ability in genetically diverse malaria infections.
431 *Proc Natl Acad Sci U S A* **102**, 7624-7628, doi:10.1073/pnas.0500078102 (2005).
- 432 3 Wargo, A. R., de Roode, J. C., Huijben, S., Drew, D. R. & Read, A. F. Transmission stage
433 investment of malaria parasites in response to in-host competition. *Proc Biol Sci* **274**, 2629-2638,
434 doi:10.1098/rspb.2007.0873 (2007).
- 435 4 Wargo, A. R., Huijben, S., de Roode, J. C., Shepherd, J. & Read, A. F. Competitive release and
436 facilitation of drug-resistant parasites after therapeutic chemotherapy in a rodent malaria
437 model. *Proc Natl Acad Sci U S A* **104**, 19914-19919, doi:10.1073/pnas.0707766104 (2007).
- 438 5 Reece, S. E., Drew, D. R. & Gardner, A. Sex ratio adjustment and kin discrimination in malaria
439 parasites. *Nature* **453**, 609-614, doi:10.1038/nature06954 (2008).
- 440 6 Wong, W., Wenger, E. A., Hartl, D. L. & Wirth, D. F. Modeling the genetic relatedness of
441 *Plasmodium falciparum* parasites following meiotic recombination and cotransmission. *PLoS*
442 *Comput Biol* **14**, e1005923, doi:10.1371/journal.pcbi.1005923 (2018).
- 443 7 Nair, S. *et al.* Single-cell genomics for dissection of complex malaria infections. *Genome Res* **24**,
444 1028-1038, doi:10.1101/gr.168286.113 (2014).
- 445 8 Mu, J. *et al.* Recombination hotspots and population structure in *Plasmodium falciparum*. *PLoS*
446 *Biol* **3**, e335, doi:10.1371/journal.pbio.0030335 (2005).
- 447 9 Mu, J. *et al.* Genome-wide variation and identification of vaccine targets in the *Plasmodium*
448 *falciparum* genome. *Nat Genet* **39**, 126-130, doi:10.1038/ng1924 (2007).
- 449 10 Neafsey, D. E. *et al.* Genome-wide SNP genotyping highlights the role of natural selection in
450 *Plasmodium falciparum* population divergence. *Genome Biol* **9**, R171, doi:10.1186/gb-2008-9-
451 12-r171 (2008).
- 452 11 Trevino, S. G. *et al.* High-Resolution Single-Cell Sequencing of Malaria Parasites. *Genome Biol*
453 *Evol* **9**, 3373-3383, doi:10.1093/gbe/evx256 (2017).
- 454 12 Chang, H. H. *et al.* THE REAL McCOIL: A method for the concurrent estimation of the complexity
455 of infection and SNP allele frequency for malaria parasites. *PLoS Comput Biol* **13**, e1005348,
456 doi:10.1371/journal.pcbi.1005348 (2017).
- 457 13 Zhu, S. J., Almagro-Garcia, J. & McVean, G. Deconvolution of multiple infections in *Plasmodium*
458 *falciparum* from high throughput sequencing data. *Bioinformatics* **34**, 9-15,
459 doi:10.1093/bioinformatics/btx530 (2018).
- 460 14 Wong, W. *et al.* Genetic relatedness analysis reveals the cotransmission of genetically related
461 *Plasmodium falciparum* parasites in Thies, Senegal. *Genome Med* **9**, 5, doi:10.1186/s13073-017-
462 0398-0 (2017).
- 463 15 Conway, D. J., Greenwood, B. M. & McBride, J. S. The epidemiology of multiple-clone
464 *Plasmodium falciparum* infections in Gambian patients. *Parasitology* **103 Pt 1**, 1-6 (1991).
- 465 16 Nkhoma, S. C. *et al.* Close kinship within multiple-genotype malaria parasite infections. *Proc Biol*
466 *Sci* **279**, 2589-2598, doi:10.1098/rspb.2012.0113 (2012).
- 467 17 Mzilahowa, T., Hastings, I. M., Molyneux, M. E. & McCall, P. J. Entomological indices of malaria
468 transmission in Chikhwawa district, Southern Malawi. *Malar J* **11**, 380, doi:10.1186/1475-2875-
469 11-380 (2012).

- 470 18 Auburn, S. *et al.* Characterization of within-host *Plasmodium falciparum* diversity using next-
471 generation sequence data. *PLoS One* **7**, e32891, doi:10.1371/journal.pone.0032891 (2012).
- 472 19 Manske, M. *et al.* Analysis of *Plasmodium falciparum* diversity in natural infections by deep
473 sequencing. *Nature* **487**, 375-379, doi:10.1038/nature11174 (2012).
- 474 20 Early, A. M. *et al.* Host-mediated selection impacts the diversity of *Plasmodium falciparum*
475 antigens within infections. *Nat Commun* **9**, 1381, doi:10.1038/s41467-018-03807-7 (2018).
- 476 21 Pearson, R. D. *et al.* Genomic analysis of local variation and recent evolution in *Plasmodium*
477 *vivax*. *Nat Genet* **48**, 959-964, doi:10.1038/ng.3599 (2016).
- 478 22 Reilly, H. B., Wang, H., Steuter, J. A., Marx, A. M. & Ferdig, M. T. Quantitative dissection of clone-
479 specific growth rates in cultured malaria parasites. *Int J Parasitol* **37**, 1599-1607,
480 doi:10.1016/j.ijpara.2007.05.003 (2007).
- 481 23 Rosario, V. Cloning of naturally occurring mixed infections of malaria parasites. *Science* **212**,
482 1037-1038 (1981).
- 483 24 Gardner, M. J. *et al.* Genome sequence of the human malaria parasite *Plasmodium falciparum*.
484 *Nature* **419**, 498-511, doi:10.1038/nature01097 (2002).
- 485 25 Volkman, S. K., Neafsey, D. E., Schaffner, S. F., Park, D. J. & Wirth, D. F. Harnessing genomics and
486 genome biology to understand malaria biology. *Nat Rev Genet* **13**, 315-328,
487 doi:10.1038/nrg3187 (2012).
- 488 26 Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated
489 individuals. *Genome Res* **19**, 1655-1664, doi:10.1101/gr.094052.109 (2009).
- 490 27 Schaffner, S. F., Taylor, A. R., Wong, W., Wirth, D. F. & Neafsey, D. E. hmIBD: software to infer
491 pairwise identity by descent between haploid genotypes. *Malar J* **17**, 196, doi:10.1186/s12936-
492 018-2349-7 (2018).
- 493 28 Miles, A. *et al.* Indels, structural variation, and recombination drive genomic diversity in
494 *Plasmodium falciparum*. *Genome Res* **26**, 1288-1299, doi:10.1101/gr.203711.115 (2016).
- 495 29 Henden, L., Lee, S., Mueller, I., Barry, A. & Bahlo, M. Identity-by-descent analyses for measuring
496 population dynamics and selection in recombining pathogens. *PLoS Genet* **14**, e1007279,
497 doi:10.1371/journal.pgen.1007279 (2018).
- 498 30 Huff, C. D. *et al.* Maximum-likelihood estimation of recent shared ancestry (ERSA). *Genome Res*
499 **21**, 768-774, doi:10.1101/gr.115972.110 (2011).
- 500 31 Li, H. *et al.* Relationship estimation from whole-genome sequence data. *PLoS Genet* **10**,
501 e1004144, doi:10.1371/journal.pgen.1004144 (2014).
- 502 32 O'Brien, J. D., Iqbal, Z., Wendler, J. & Amenga-Etego, L. Inferring Strain Mixture within Clinical
503 *Plasmodium falciparum* Isolates from Genomic Sequence Data. *PLoS Comput Biol* **12**, e1004824,
504 doi:10.1371/journal.pcbi.1004824 (2016).
- 505 33 Assefa, S. A. *et al.* estMOI: estimating multiplicity of infection using parasite deep sequencing
506 data. *Bioinformatics* **30**, 1292-1294, doi:10.1093/bioinformatics/btu005 (2014).
- 507 34 Galinsky, K. *et al.* COIL: a methodology for evaluating malarial complexity of infection using
508 likelihood from single nucleotide polymorphism data. *Malar J* **14**, 4, doi:10.1186/1475-2875-14-
509 4 (2015).
- 510 35 Annan, Z. *et al.* Population genetic structure of *Plasmodium falciparum* in the two main African
511 vectors, *Anopheles gambiae* and *Anopheles funestus*. *Proc Natl Acad Sci U S A* **104**, 7987-7992,
512 doi:10.1073/pnas.0702715104 (2007).
- 513 36 Portugal, S. *et al.* Host-mediated regulation of superinfection in malaria. *Nat Med* **17**, 732-737,
514 doi:10.1038/nm.2368 (2011).
- 515 37 Venkatesan, M. *et al.* Using CF11 cellulose columns to inexpensively and effectively remove
516 human DNA from *Plasmodium falciparum*-infected whole blood samples. *Malar J* **11**, 41,
517 doi:10.1186/1475-2875-11-41 (2012).

- 518 38 Zheng, X. *et al.* SeqArray-a storage-efficient high-performance data format for WGS variant calls.
519 *Bioinformatics* **33**, 2251-2257, doi:10.1093/bioinformatics/btx145 (2017).
520 39 Cayuela, L. & Gotelli, N. J. (2014).
521 40 Jiang, H. *et al.* High recombination rates and hotspots in a *Plasmodium falciparum* genetic cross.
522 *Genome Biol* **12**, R33, doi:10.1186/gb-2011-12-4-r33 (2011).
523 41 Team, R. C. (2017).