

1 Identity-by-descent relatedness estimates with uncertainty characterise
2 departure from isolation-by-distance between *Plasmodium falciparum*
3 populations on the Colombian-Pacific coast

4 Aimee R. Taylor^{*,1,2}, Diego F. Echeverry^{3,4,5}, Timothy J. C. Anderson⁶, Daniel E. Neafsey^{2,7}, Caroline O. Buckee¹

5 *Corresponding author: ataylor@hsph.harvard.edu

6
7 **1** Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard T. H. Chan School of Public
8 Health, Boston, Massachusetts, USA

9 **2** Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

10 **3** Centro Internacional de Entrenamiento e Investigaciones Médicas (CIDEIM), Cali, Colombia

11 **4** Universidad Icesi, Cali, Colombia

12 **5** Departamento de Microbiología, Facultad de Salud, Universidad del Valle, Cali, Colombia

13 **6** Department of Genetics, Texas Biomedical Research Institute, San Antonio, Texas, USA

14 **7** Department of Immunology and Infectious Diseases, Harvard T. H. Chan School of Public Health, Boston, Mas-
15 sachusetts, USA

16 **Keywords:** malaria; *Plasmodium falciparum*; relatedness; identity-by-descent; molecular epidemiology;
17 spatial connectivity; isolation-by-distance

18 **Abstract**

19 Characterising connectivity between geographically separated biological populations is a common goal in
20 many fields. Recent approaches to understanding connectivity between malaria parasite populations, with
21 implications for disease control efforts, have used estimates of relatedness based on identity-by-descent (IBD).
22 However, uncertainty around estimated relatedness has not been accounted for to date. IBD-based relat-
23 edness estimates with uncertainty were computed for pairs of monoclonal *Plasmodium falciparum* samples
24 collected from five cities on the Colombian-Pacific coast where long-term clonal propagation of *P. falciparum*
25 is frequent. The cities include two official ports, Buenaventura and Tumaco, that are separated geographically
26 but connected by frequent marine traffic. The fraction of highly-related sample pairs (whose classification
27 accounts for uncertainty) was greater within cities versus between. However, based on both the fraction
28 of highly-related sample pairs and on a threshold-free approach (Wasserstein distances between parasite
29 populations) connectivity between Buenaventura and Tumaco was disproportionately high. Buenaventura-
30 Tumaco connectivity was consistent with three separate transmission events involving parasites from five
31 different clonal components (groups of statistically indistinguishable parasites identified under a graph theo-
32 retic framework). To conclude, *P. falciparum* population connectivity on the Colombian-Pacific coast abides
33 by accessibility not isolation-by-distance, potentially implicating marine traffic in malaria transmission with
34 opportunities for targeted intervention. Further investigations are required to test this and alternative hy-
35 potheses. For the first time in malaria epidemiology, we account for uncertainty around estimated relatedness
36 (an important consideration for future studies that plan to use genotype versus whole genome sequence data
37 to estimate IBD-based relatedness); we also use a threshold-free approach to compare parasite populations,
38 and identify clonal components in a statistically principled manner. The approaches we employ could be
39 adapted to other recombining organisms with mixed mating systems, thus have broad relevance.

40 Introduction

41 In many research fields genetic data are used to help characterise connectivity between geographically dis-
42 parate biological populations, with numerous applications in conservation, agriculture, and public health.
43 Patterns of genetic similarity between pathogen populations help us understand how the disease spreads.
44 Patterns of relatedness (a measure of genetic similarity) between malaria parasites in different human pop-
45 ulations, for instance, help characterise the connectivity between them, thus guide the design of targeted
46 public health interventions [1].

47 Several methods are employed to measure genetic similarity and thus characterise connectivity. Phylo-
48 genetic methods, in which genetic distances between individuals are measured in units of mutation [2], are
49 most applicable to rapidly mutating organisms that do not recombine (e.g RNA viruses) [3]. Studies of re-
50 latedness, in which relatedness is a measure of probability of identity-by-descent (IBD) between individuals,
51 are applicable to organisms that do recombine (e.g. malaria parasites). Population genetic parameters of
52 allelic variation (e.g. F_{ST}) are applicable to all organisms (those that do and do not recombine), but do not
53 generate measures of genetic distance or similarity on an inter-individual level, thus provide less granularity.
54 Moreover, among recombining organisms, inter-population allelic variation tends to accumulate more slowly
55 than inter-individual variation in IBD [4]. As such, analyses of relatedness can sometimes nearby and recent
56 connectivity where analyses of F_{ST} cannot [5].

57 Malaria parasites are protozoan parasites that undergo an obligate stage of sexual recombination in the
58 mosquito midgut. Like many organisms (e.g. many plants [6, 7]), malaria parasites have a mixed mating
59 system that encompasses both inbreeding and outcrossing. The extent to which malaria parasites outcross
60 depends on transmission intensity and is not fully understood [8]. In any event, for outcrossing to occur
61 a mosquito must ingest genetically distinct gametocytes. Humans can be infected by multiple genetically
62 distinct parasite clones that are either co-transmitted via inoculation from a single mosquito, in which case
63 they are likely recombinants so inter-related, or transmitted independently by multiple mosquitoes, in which
64 case the parasite clones are likely unrelated [9, 10]. The latter can occur in a setting where the entomological
65 inoculation rate is high; recent work suggests co-transmission is important in both low and high transmission
66 settings [10].

67 Malaria genomic epidemiology studies of connectivity are increasingly common, especially in the context
68 of public health and using genotype (versus whole genome sequence) data [5, 11–14]. Using IBD-based
69 relatedness but not F_{ST} , evidence of isolation-by-distance among *P. falciparum* populations along a 100
70 km stretch of the Thailand-Myanmar border was found [5]. This study was based, in part, on analyses of
71 monoclonal *P. falciparum* samples genotyped at 93 single nucleotide polymorphisms (SNPs). Based on F_{ST}
72 estimated using *P. falciparum* samples genotyped at 250 SNPs, a different study found evidence of departure
73 from isolation-by-distance among *P. falciparum* populations along a 500 km stretch of the Colombian-Pacific
74 coast where transmission is mixed (low and high in some regions) and outcrossing limited [11, 15]. Departure
75 from isolation-by-distance on the Colombian-Pacific coast was based on F_{ST} alone [11]. In the current study,
76 we re-explore departure from isolation-by-distance with more granularity using IBD-based relatedness. For
77 the first time in malaria epidemiology, we account for uncertainty in relatedness estimates; we also use a
78 threshold-free approach to compare parasite populations, and identify clonal components in a statistically
79 principled manner. The original study [11] and our response to it are described in more detail below.

80 Malaria epidemiology in Colombia is associated with a multitude of ecological, evolutionary and social
81 factors (Table A.1), including human migration due to deforestation, illegal crops, gold mining [16–20], and
82 the mass emigration of people fleeing the humanitarian crisis in Venezuela [21–24]. Understanding the in-
83 terplay between e.g. migration, parasite population connectivity and the spread of antimalarial resistance
84 is critical [16, 18]. In preparation for studies of resistance, Echeverry et al. genotyped *P. falciparum* sam-
85 ples from four provinces on the Colombian-Pacific coast [11]. Clonality, population structure and linkage
86 disequilibrium (LD) were characterised using a suite of population genetic analyses. The results were highly
87 informative: the vast majority of successfully genotyped *P. falciparum* samples were deemed monoclonal
88 (325 of 400) with a strong association between clonality and incidence. Among the 325 monoclonal samples,
89 136 unique haploid multilocus genotypes (MLGs) were identified using relatedness based on identity-by-state
90 (IBS), which is a correlate of IBD [25] (and has been used elsewhere to characterise connectivity between
91 nearby malaria parasite populations [12–14]). Of the 136 MLGs, 44 infected two or more patients (max.
92 28 patients), 45 persisted for two or more days (max. 8 years), and 7 of the 15 most common MLGs were

93 sampled in two or more provinces (max. all four provinces). Panmixia was rejected based on evidence of
94 four sympatric but geographically structured subpopulations; and, overall, LD decayed at a rate that was
95 faster than expected for South American *P. falciparum* populations (compare with e.g. [26]). Echeverry et al.
96 concluded that evidence of low genetic diversity, persistent MLGs and population structure is consistent with
97 low transmission and limited outcrossing, while evidence of a relatively fast rate of LD decay and of shared
98 MLGs across provinces is consistent with extensive human movement connecting *P. falciparum* populations.

99 Although the study by Echeverry et al. features analyses of IBS-based relatedness (i.e. MLGs), evidence
100 of departure from isolation-by-distance was based on F_{ST} alone. To explore in more granularity while ac-
101 counting for uncertainty, we compute IBD-based relatedness estimates and confidence intervals for all pairs of
102 325 monoclonal parasite samples. Akin to previous studies (e.g. [5]), highly-related parasites were classified
103 using a threshold; however, confidence intervals allow uncertainty to be accounted for. This is important
104 because relatedness estimated using limited genotype data can be overwhelmed by uncertainty [25]. Our ap-
105 proach includes two additional contributions. First, we complement our analysis of highly-related parasites
106 with a threshold-free approach based on optimal transport using Wasserstein distances between parasite
107 populations. Second, we identify groups of statistically indistinguishable parasites, which we call clonal com-
108 ponents, using the simple concept of components from graph theory and confidence intervals. Confidence
109 intervals circumvent reliance on an arbitrary clonal threshold (i.e. some number of differences tolerated
110 between parasites samples considered clonal). Graph components circumvent reliance on unsupervised clus-
111 tering methods that are notoriously brittle [27]. Overall, our approach could be adapted to viruses and
112 bacteria that show recombination or reshuffling of segments as well as clonal propagation [28–31], to other
113 protozoans (e.g. *Toxoplasma*, *Cryptosporidium* [32–34]), and to the many fungal pathogens [35], plants [6,7],
114 and animals with mixed mating systems. Due to our treatment of uncertainty, it is especially relevant for a
115 growing number of studies that plan estimate IBD-based relatedness using genotype (versus sequence) data.

116 Results

117 Relatedness estimates between *P. falciparum* sample pairs

118 Relatedness was estimated for all 52650 pairwise comparisons of 325 previously published monoclonal *P.*
119 *falciparum* samples with data on 250 biallelic single nucleotide polymorphisms (SNPs) [11]. The parasite
120 samples were collected between 1993 and 2007 from symptomatic patients participating in studies at five cities
121 on the Colombian-Pacific coast (Table A.2). Despite considerable uncertainty, all estimates are informative
122 (Figure 1). That is to say, there are no relatedness estimates whose 95% confidence intervals span entirely
123 from zero to one. The vast majority of relatedness estimates were classified unrelated.

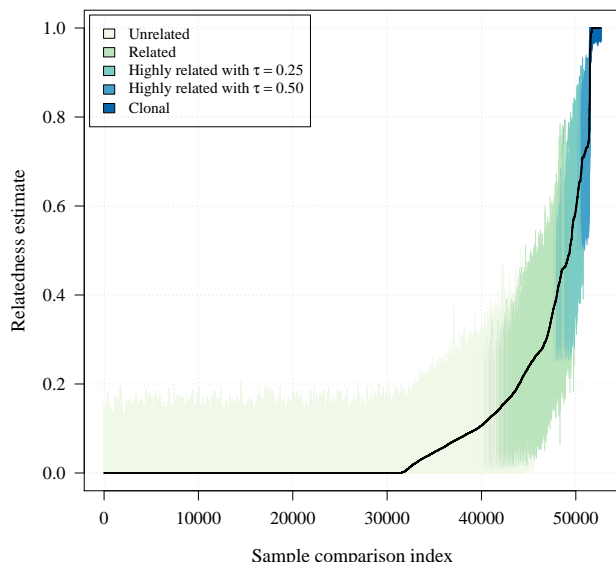


Figure 1: Estimates of relatedness with colour-coded 95% confidence intervals for all 325 choose two (52650) *P. falciparum* sample pairs in order of increasing relatedness estimate. Confidence intervals are coloured according to classifications based on lower and upper confidence interval bounds (Table 1).

Classification	Interpretation	Definition
Unrelated	\hat{r} statistically indistinguishable from zero	$\text{LCI} < \epsilon$
Related	\hat{r} statistically distinguishable from zero	$\text{LCI} > \epsilon$
highly-related	\hat{r} statistically distinguishable from a specified threshold	$\text{LCI} > \tau$
Clonal	\hat{r} statistically indistinguishable from one	$\text{UCI} > 1 - \epsilon$

Table 1: Classification of parasite sample pairs. Lower and upper confidence interval bounds (LCI and UCI, respectively) are used to classify pairs with $\epsilon = 0.01$, $\tau = 0.25$ (main analysis, Figure 2) and $\tau \in \{0.25, 0.50\}$ (sensitivity analysis, Figure A.1).

124 highly-related *P. falciparum* sample pair fractions partitioned in space and time

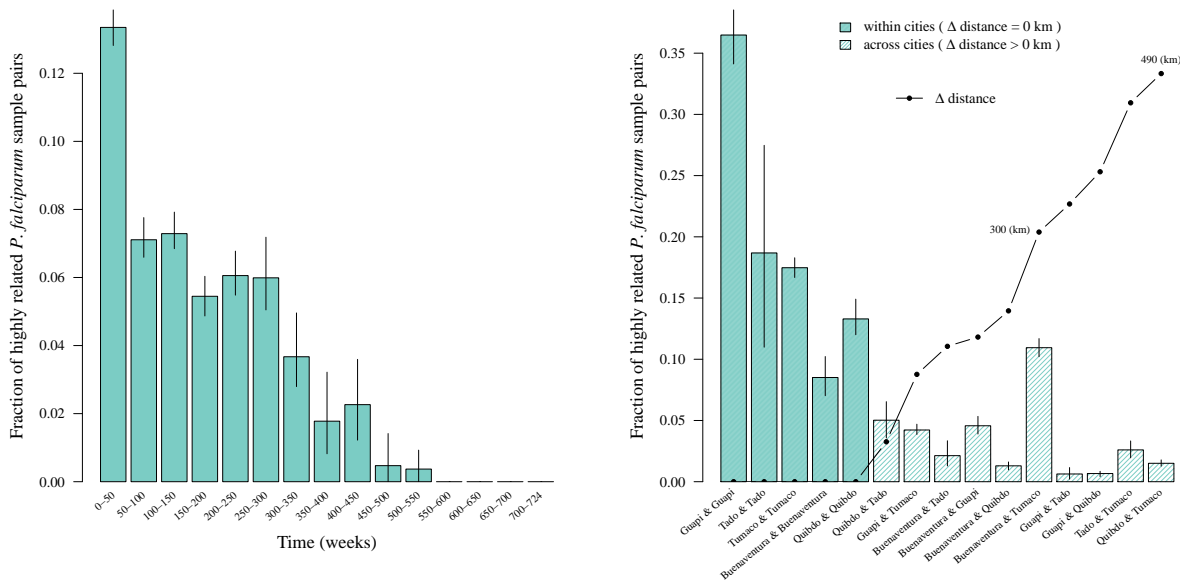
125 In our main analysis (Figure 2), highly-related parasite samples were classified using a high-relatedness
126 threshold of 0.25 (Table 1). Despite few highly-related *P. falciparum* sample pairs overall, there are three
127 notable observations regarding their fraction partitioned in space and time. First, there is a greater fraction
128 of highly-related sample pairs among those collected close together in time versus far apart (Figure 2a).
129 Second, the fraction of highly-related sample pairs is generally greater within cities versus between, with
130 Guapi having the largest fraction of highly-related pairs and Buenaventura having the lowest (Figure 2b).
131 However, third, the fraction shared between Buenaventura and Tumaco is exceptionally high considering
132 inter-city distance (Figure 2b). These observations are largely robust to different high-relatedness thresholds
133 (Figure A.1). Spatial trends evaluated using a threshold-free approach are also consistent: they show a
134 general increase in 1-Wasserstein distance with geographic distance between cities besides Buenaventura and
135 Tumaco (Figure 3). The 1-Wasserstein distance can be interpreted as the effort required to transform a
136 distribution of parasite samples from one city into a distribution of parasite samples from another [16]. The
137 small 1-Wasserstein distance between Buenaventura and Tumaco is thus consistent with elevated gene flow
138 between *P. falciparum* populations sampled from these cities.

139 Figure 4a shows the inter-city *P. falciparum* population connectivity plotted in Figure 2b projected onto
140 a map of the Colombian-Pacific coast. Buenaventura and Tumaco are the two largest official ports on the
141 Colombian-Pacific coast (Buenaventura is the largest) and are connected by frequent marine traffic (Figure
142 4b). Although Tumaco is connected to Buenaventura via the pan-american highway, which connects all sites
143 but Guapi, primary access to Tumaco is via the port due to difficult and unsafe country roads in Nariño.
144 Guapi, which is effectively unreachable by road and not an official port, is connected by marine traffic
145 but with less frequency (Figure 4b). Consistent with its isolation, the fraction of highly-related parasite
146 pairs is relatively large within Guapi (Figure 2b), and very small between Guapi and the two inland cities,
147 Quibdó and Tadó (Figures 2b and 4a). Moreover and importantly regarding the elevated fraction of highly-
148 related samples pairs within both Guapi and Tadó (Figures 2b), all samples from Guapi and Tadó were
149 collected within a single year (Table A.2). The low fraction of highly-related parasite sample pairs within
150 Buenaventura (Figure 2b) is consistent with it having contributed samples over many years (Table A.2) and
151 with it being the most important port on the Pacific coast (Figure 4b), i.e. a hub through which human
152 traffic and thus potential parasite mixing is high [11].

153 The apparent association between *P. falciparum* population connectivity and the frequency of marine
154 traffic raises questions about the latter's role in malaria transmission. However, other scenarios could lead
155 to these relationships, for example the high connectivity could result from a single travel event between
156 Buenaventura and Tumaco, followed by expansion of highly-related and clonal parasites. To further explore
157 the genetic signal that supports this association we next consider clonal components.

158 Clonal components

159 We define clonal components as groups of statistically indistinguishable parasite samples identified under a
160 graph theoretic framework; see Methods. In total, 46 distinct clonal components were detected, ranging in
161 size from 2 to 28 statistically indistinguishable parasite samples (Figure 5). They are spatially clustered.
162 Ten of the 46 contain parasite samples collected from two or more cities. Each clonal component besides
163 one (clonal component four) is on average related to at least one other (Figure 5). The unrelated clonal



(a) Partitioned by time between collection dates.

(b) Partitioned by collection city.

Figure 2: Fractions of highly-related sample pairs partitioned in time and space.

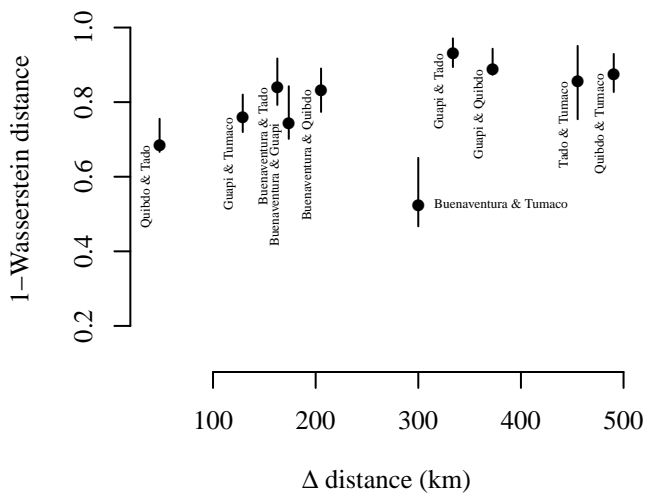
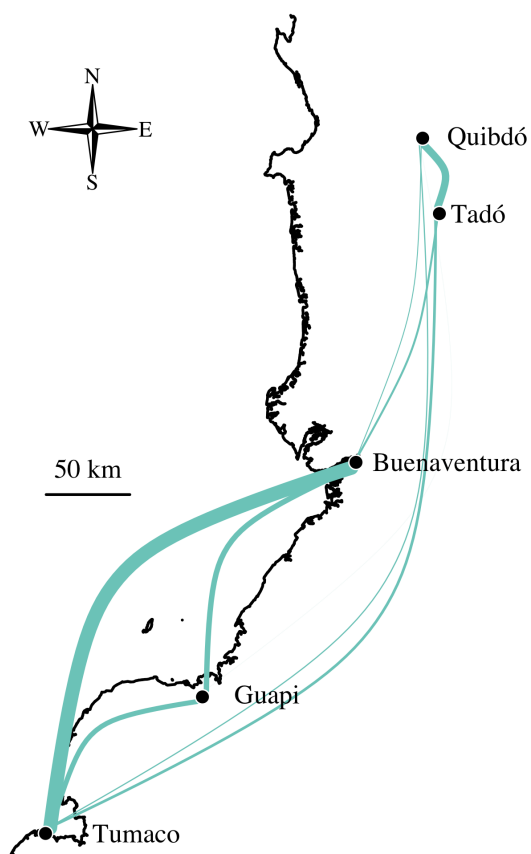


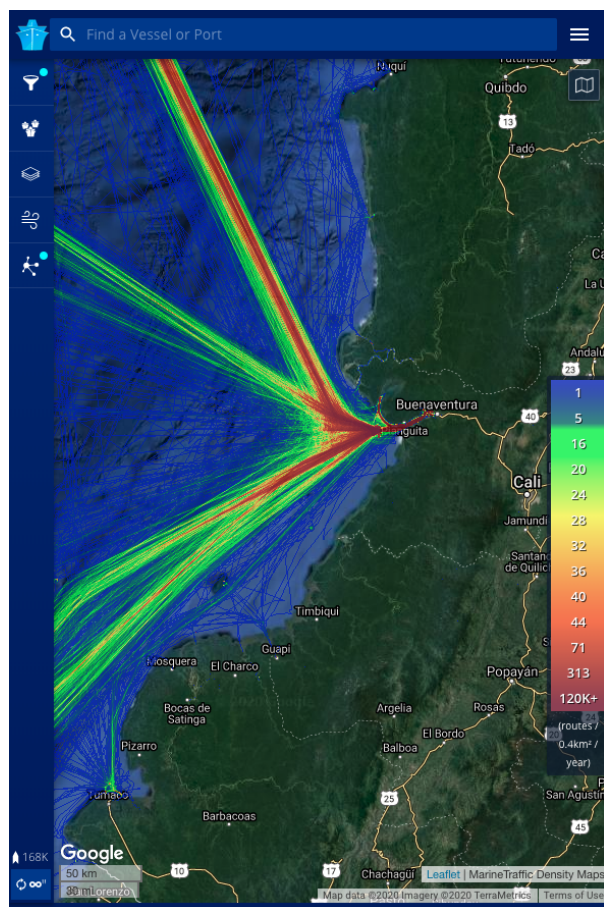
Figure 3: 1-Wasserstein distance between parasite populations across cities versus inter-city distance in kilometres (km).

164 component is likely a contaminant: it accords with MLG 036 reported in [11], where contamination during
 165 *in vitro* adaptation or DNA manipulation was suspected.

166 Clonal parasite samples detected in both Buenaventura and Tumaco belong to five distinct clonal compo-
 167 nents (one, 12, 14, 20 and 40, Figure 5). We thus dismiss a single travel event connecting Buenaventura and
 168 Tumaco involving a single parasite clone. We cannot dismiss a single travel event involving multiple parasite
 169 clones, however. Indeed, three of the five clonal components are inter-related on average (Table A.3). As
 170 such, they could derive from co-transmitted recombinant parasites transported in a single individual with a
 171 multiclonal infection. On the contrary, the remaining two clonal components have relatedness estimates that
 172 are not statistically distinguishable from zero. As such, they likely derive from different individuals with
 173 independent monoclonal infections. Given dates and cities of first detection (Table A.4), it is tempting to
 174 suggest some clonal components predate others and originate in specific locations. For example, it is possible
 175 that parasite samples from clonal components 1 and 20 in Buenaventura and Tumaco emanated from Guapi,



(a) *P. falciparum* population connectivity: the width of each inter-city edge is proportional to the fraction of highly-related sample pairs across cities plotted in Figure 2b. Note that the edges between Guapi and Quibdó and Guapi and Tadó are plotted but too thin to visually discern.



(b) Screen shot of official marine traffic frequency (routes per 0.4km^2 per year) taken from www.marinetraffic.com 2020-02-12. The road system, which includes the pan-american highway, is visible (faint orange line) on the land map (Google satellite), as are some of the province boundaries (dashed grey lines). Zoom in to see city names Tumaco, Guapi, Buenaventura, Quibdó and Tadó among others.

Figure 4: Comparison between *P. falciparum* population connectivity and the frequency of marine traffic.

176 creating a spurious link between Buenaventura and Tumaco. However, because these data are from sparsely
177 sampled symptomatic cases in setting where clonal propagation is frequent, sample collection chronology is
178 not necessarily representative of transmission chain events (Figure 6).

179 Regarding transmission chain events, we note that clonal component 20 relates to the three inter-related
180 clonal components (1, 12 and 14) via an intermediate clonal component detected in Tumaco only (clonal
181 component 15) as well as an intermediate parasite sample from Quibdó that does not belong to a clonal
182 component (Figure A.2). These intermediates likely derive from recombination between parasites related to
183 the clonal components they connect. Several connections consistent with recombinants can be found among
184 the relatedness graphs (Figures 5 and A.2). As such, it seems it may be at least theoretically possible to
185 construct approximate *P. falciparum* transmission chains given more dense sampling of malaria infections
186 on the Colombian-Pacific coast.

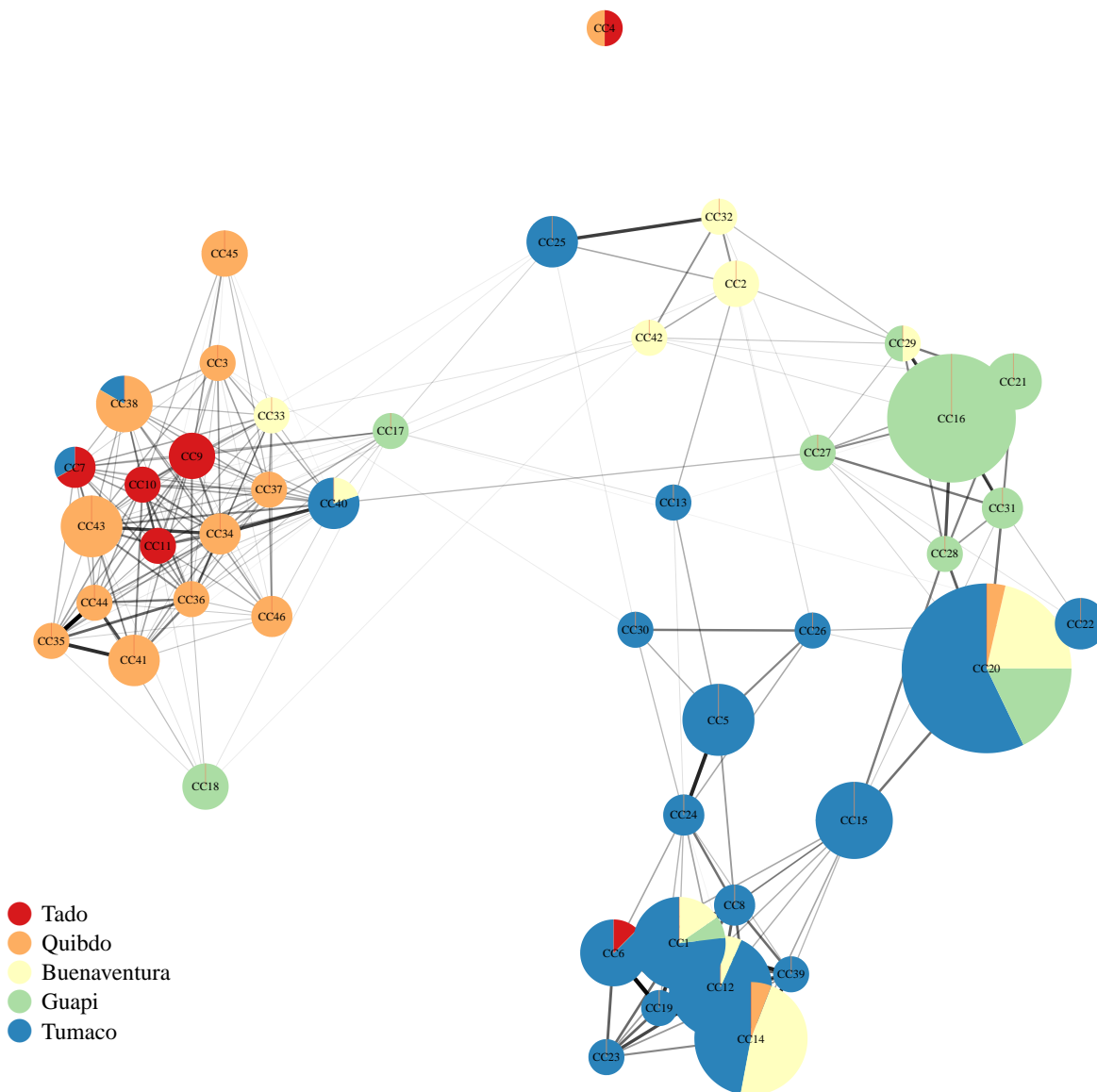


Figure 5: Clonal components (vertices) and the average relatedness between them (edges). Clonal components (CCs) are groups of two or more statistically indistinguishable parasite samples. CC vertices are plotted using the Fruchterman-Reingold layout algorithm [36], thereby clustering inter-related CCs. The size of each CC vertex is proportional to the number of parasite samples per CC, ranging from 2 to 28 statistically indistinguishable parasite samples. CCs are named in order of the collection date of the earliest parasite sample per CC (Table A.4. CCs with parasite samples collected from two or more cities are depicted as pie charts. Colour denotes the city of parasite sample collection. Edge transparency and weight is proportional to average relatedness, ranging from 0.003 to 0.840. Relatedness estimates that are indistinguishable from zero were set to zero. Edges whose average relatedness is zero are not plotted. Each CC besides CC4 is related to at least one other. CC4 contains two samples (one from Tadó, another from Quibdó). It is likely a contaminant; see main text. A plot of CCs that includes singletons (individual parasite samples that do not belong to a CC) can be found in the Appendix (Figure A.2).

187 Discussion

188 Here we show that estimates of IBD-based relatedness and their associated uncertainty can be used to
 189 uncover epidemiologically meaningful connectivity between *P. falciparum* populations on a relatively local

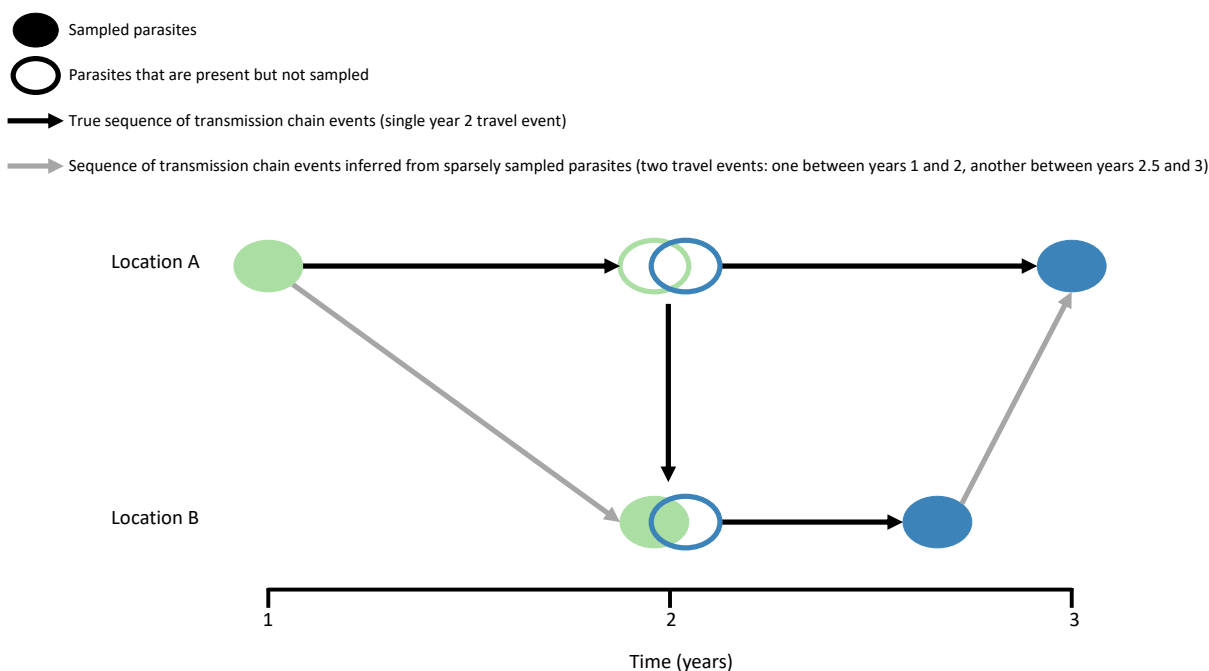


Figure 6: Schematic illustrating why sample collection chronology is not necessarily representative of the true sequence of transmission chain events when sampling is sparse and clonal propagation is frequent. The schematic shows two hypothetical locations A and B where malaria parasites have been sampled sparsely: solid ellipses represent sampled parasites, open ellipses represent parasites that were present but not sampled, different colours denote different parasite genotypes.

190 spatial scale: along the Colombian-Pacific coast where clonal propagation is frequent. While our approach
 191 largely confirms a previous report based on F_{ST} [11], estimates of relatedness provide more granularity while
 192 their confidence intervals account for uncertainty thus provide statistical rigor, e.g. when highly-related
 193 parasite sample pairs are classified. Our approach includes two additional contributions: concepts from
 194 optimal transport (1-Wasserstein distance) are used to compare parasite populations in an entirely threshold-
 195 free manner; and clonal components are identified using graph components and confidence intervals, thereby
 196 circumventing reliance on an arbitrary clonal threshold. Our overall approach could be adapted for analyses of
 197 viruses and bacteria that show recombination or reshuffling of segments as well as clonal propagation [28–31],
 198 to other protozoans (e.g. *Toxoplasma*, *Cryptosporidium* [32–34]), and to the many fungal pathogens [35],
 199 plants [6, 7], and animals with mixed mating systems.

200 IBD-based relatedness estimates recovered 1) a large fraction of highly-related parasite sample pairs
 201 within Guapi, a city on the Colombian-Pacific coast that is relatively isolated besides infrequent marine
 202 traffic; 2) a low fraction of highly-related parasite sample pairs within Buenaventura, the most important
 203 port on the Colombian-Pacific coast and thus the least isolated city in this study; and 3) a disproportionately
 204 large fraction of highly-related parasite pairs between Buenaventura and Tumaco (departure from isolation-
 205 by-distance), where Tumaco is the second largest port on the Colombian-Pacific coast. These observations
 206 accord with several published previously: 1) elevated LD in a *P. falciparum* subpopulation (identified using
 207 STRUCTURE [27,37]) predominant in Guapi; 2) rapid LD decay in a *P. falciparum* subpopulation predom-
 208 inant in Buenaventura; and 3) lowest genetic differentiation (based on F_{ST} estimates) between provinces
 209 Valle (Buenaventura) and Nariño (Tumaco) [11]. LD, STRUCTURE and F_{ST} analyses all rely on allelic
 210 variation. The concordance between results based on relatedness and allelic variation suggests that *P. falciparum*
 211 outbreeding on the Colombian-Pacific coast is infrequent enough that both types of analyses generate

212 insight on approximately the same time scale.

213 The aforementioned results generate hypotheses around the frequency of marine traffic and malaria
214 transmission on the Colombian-Pacific coast. Notwithstanding long-range windborne dispersal, which may
215 be critical for malaria transmission in Africa [38], anopheline flight range is generally small (around 3.5
216 km [39]). As such, long-range malaria parasite dispersal on the Colombian-Pacific coast is almost certainly
217 human-mediated. A recent study of *P. vivax* proposed that human movement across a “malaria corridor”
218 stretching from the northwest to the south of the Colombian-Pacific Coast likely promotes *P. vivax* gene
219 flow, and that mining activities may provide transmission “contact zones” [40], similarly proposed for *P.*
220 *falciparum* [20]. *P. falciparum* population connectivity is consistent with the human “malaria corridor”
221 hypothesis, especially since it correlates with accessibility, not isolation-by-distance. Both infected humans
222 and mosquitoes are compatible with this hypothesis, i.e. checks for infected *Anopheles spp.* on boats may
223 be merited [41, 42]. However, relatively high differentiation between populations of *An. albimanus* (one of
224 the three primary vectors of malaria in Colombia [43]) from Buenaventura and Tumaco [44], points towards
225 human carriage.

226 The Colombian-Pacific coast has long been associated with the international trade of gold and narcotics,
227 but until recently human migration in the region was largely domestic. The flow of international migrants
228 infected with *Plasmodium spp.* has increased significantly in recent years. In 2019, 2190 of 2288 (95.7%)
229 of non-domestic malaria cases reported in Colombia were from Venezuela; other sources included South
230 America (Peru, Panama, French Guyana, Ecuador, Brazil) and some African countries (Uganda, Republic
231 of the Congo, Nigeria, Ivory Coast, Cameroon, Angola) [45]. Some of the infected Venezuelan nationals are
232 migrating southward to Ecuador and Peru [22]. Other non-domestic cases may be associated with the traffic
233 of people who arrive at Colombian ports with a view towards northward travel e.g. to the USA via Central
234 America and Panama [46]. Genetic surveillance of “international parasites” may help malaria control efforts
235 in Colombia.

236 Considerable violence in the South Pacific region of Colombia between 1993 and 2007 combined with
237 historically high malaria case counts [15] could have caused fleeting connectivity between *P. falciparum*
238 populations. Based on relatedness between clonal components, we refute the hypothesis that connectivity
239 between Buenaventura and Tumaco was due to a single individual with a multiclonal infection. We cannot
240 reject a single travel event involving multiple individuals with independent infections, however. Contempo-
241 rary data on more densely sampled cases and on mosquito and human movement are required to characterise
242 extant connectivity, its reach beyond Colombia (see e.g. [47]), and to rule out alternative hypotheses.

243 Regarding alternative hypotheses, heterogeneous vectorial capacity and antimalarial drug pressure could
244 selectively enhance parasite survival in such a way that generates apparent connectivity between Buenaven-
245 tura and Tumaco, e.g. if parasites are adapted to local vectors whose distributions are more similar between
246 Buenaventura and Tumaco than elsewhere. Although adult *An. albimanus* B and *An. neivai* s.l. have
247 been detected in the vicinities of both cities [44, 48], the species distributions in the vicinities of Buenaven-
248 tura and Tumaco differ more than those in the vicinities of Tumaco and Guapi [48]. As such, heterogeneous
249 vectorial capacity seems an unlikely alternative hypothesis. Similarly, relatedness may be greater among par-
250 asites with comparable antimalarial resistance: a recent study of South East Asian *P. falciparum* parasites
251 found greater relatedness in the recent past among parasites with artemisinin resistance mutations versus
252 without [49]. This study used size-stratified IBD segments to date relatedness [49].¹ On the Colombian-
253 Pacific coast, IBD segment size inference could help identify some recently related parasites. However, it
254 requires whole genome sequence data and is hard (if not presently impossible) to interpret in the face of
255 frequent selfing that is transmission dependent [25]. The development of an ancestral recombination model
256 that incorporates transmission-dependent selfing is a research priority in malaria genetic epidemiology and
257 would aid research on other organisms that show both outbreeding and clonal propagation.

¹In malaria, IBD segments (genome segments that are descended from a common ancestor unbroken by recombination [4]) are broken down each time genetically distinct parasites outcross. As such, IBD segment size is distributed according to a number of out-crossed generations, which increase over time, albeit in a complex transmission-dependent manner, which is not yet fully understood [9, 10].

258 Methods

259 Data

260 This study relies entirely on previously published data that are publicly available [11,25]. In the original
261 study by Echeverry et al., finger-prick blood spot samples were obtained from patients with symptomatic
262 uncomplicated malaria [11]. Samples were collected between 1993 and 2007 from five cities in four provinces:
263 Tadó and Quibdó in Chocó, Buenaventura in Valle, Guapi in Cauca and Tumaco in Nariño (Table A.2) [11].
264 Informed consent was obtained from all the subjects enrolled, as approved by CIDEIM Institutional Review
265 Board (IRB) [11]. The Colombian-pacific coast is one of the rainiest regions of the world [44,50]. At that
266 time, Colombia had approximately 100,000 malaria cases per year [11,15]. Collectively Chocó, Valle, Cauca
267 and Nariño accounted for up to 75% of the *P. falciparum* cases reported, with relatively high transmission
268 in Chocó and relatively low transmission in Valle and Cauca [11].

269 The data that feature in this descriptive study also feature in a recent methodological study concerning
270 data requirements for relatedness inference [25]. As in [25], we did not post-process the data in any way
271 besides mapping SNP positions to the *P. falciparum* 3d7 v3 reference genome and recoding heteroallelic
272 calls as missing (since all samples with fewer than 10 heteroallelic SNP calls were classified monoclonal
273 previously [11]). The monoclonal data include 325 *P. falciparum* samples with data on 250 biallelic SNPs
274 whose minor allele frequency estimates (the minor allele sample count divided by 325) range from 0.006 to
275 0.495 (Figure A.3).

276 Relatedness inference and classification of parasite sample pairs and groups

277 For each pairwise parasite sample comparison, we generated a relatedness estimate and 95% confidence
278 interval using the hidden Markov model and parametric bootstrap described in [25]. Sample pairs were
279 classified as unrelated, related, highly-related and clonal using confidence interval bounds as follows and
280 summarised in Table 1. A pair was classified unrelated if its relatedness estimate, \hat{r} , was statistically
281 indistinguishable from zero with lower confidence interval bound (LCI) less than ϵ . A pair was classified
282 related if its relatedness estimate, \hat{r} , was statistically distinguishable from zero with $\text{LCI} > \epsilon$. A pair was
283 considered highly-related if its relatedness estimate, \hat{r} , was statistically distinguishable from some specified
284 threshold, τ , with $\text{LCI} > \tau$. A pair was considered clonal if its relatedness estimate, \hat{r} , was statistically
285 indistinguishable from one with upper confidence interval bound (UCI) $> 1 - \epsilon$. Note that these classifications
286 are possible because all estimates are informative, i.e. no confidence intervals span the entire zero to one
287 range (Figure 1). These classifications are neither necessarily exclusive nor conversely true: a clonal parasite
288 pair is related, but a related parasite pair is not necessarily clonal. Throughout, $\epsilon = 0.01$. In the main
289 analysis (Figure 2) $\tau = 0.25$, in the sensitivity analysis (Figure A.1a) $\tau \in \{0.25, 0.50\}$.

290 In addition to classifying parasite sample pairs, we classify groups of statistically indistinguishable parasite
291 samples, which we call clonal components because they are defined using the simple concept of components
292 from graph theory. First, we construct a super-graph whose vertices are parasite samples connected by edges
293 that are weighted by relatedness estimates. Within the super-graph, a clonal component is a sub-graph
294 within which all parasite samples are connected to one another (directly or not) via edges whose weights are
295 statistically indistinguishable from one, while being connected to parasites samples outside the sub-graph
296 via edges whose weights are not statistically indistinguishable from one. Clonal components tend to be fully
297 connected (i.e. all parasite samples within the clonal component are directly connected to one another by
298 edges whose weights are statistically indistinguishable from one). The *igraph* package [51] in R [52] was used
299 to identify clonal components and to visualise them using the Fruchterman-Reingold layout algorithm [36].

300 Spatiotemporal trends in *P. falciparum* population connectivity

301 Spatiotemporal trends in population connectivity were explored visually by partitioning parasite sample
302 pairs by their collection cities and dates, then plotting the per-partition fraction of highly-related pairs.
303 Error bars were constructed by re-sampling per-partition parasite sample pairs 100 times with replacement
304 and taking the 2.5th and 97.5th percentiles of the fraction of highly-related pairs as the lower and upper
305 limits, respectively. Sensitivity to $\tau = 0.25$ (high relatedness threshold used in Figure 2) was explored using
306 alternative $\tau \in \{0.25, 0.50\}$ (Figure A.1) and also by using a threshold-free approach (Figure 3) as follows.

307 To explore population connectivity using a threshold-free approach, we calculated 1-Wasserstein distances
308 between groups of parasite samples from different cities using the `transport` [53] package in R [52]. Specifi-
309 cally, for a pair of cities a and b , we construct a $n_a \times n_b$ genetic distance matrix, G , of $1 - \hat{r}_{ij}$ (where n_a and
310 n_b are the parasite sample counts from cities a and b , respectively, $i = 1, \dots, n_a$ and $j = 1, \dots, n_b$) and two
311 vectors $w_a = (1/n_a, \dots, 1/n_a)$ and $w_b = (1/n_b, \dots, 1/n_b)$ of length n_a and n_b , respectively. We then calculate the
312 1-Wasserstein distance, which minimises the total cost of transporting w_a to w_b , where $1 - \hat{r}_{ij}$ is the cost of
313 transporting a single unit, using `transport::transport(w_a, w_b, costm = G, method = "shortsimplex")`.
314 This amounts to treating parasite samples from different cities as draws from different distributions, where
315 the 1-Wasserstein distance can be interpreted as the effort required to transform a distribution of parasite
316 samples from one city into a distribution from another [16]. City pairs with smaller 1-Wasserstein distances
317 are interpreted as having greater connectivity between the *P. falciparum* populations collected from them.
318 Error bars were constructed by re-sampling parasite sample pairs per inter-city partition 100 times with
319 replacement and taking the 2.5th and 97.5th percentiles of the distribution 1-Wasserstein distances based
320 on the re-sampled sample pairs as the lower and upper limits, respectively.

321 Data and code availability

322 All data analyses were performed in R [52]. The data are publicly available as a `.RData` files and the code is
323 publicly available as `.R` scripts at <https://github.com/artaylor85/ColombianBarcode>.

324 Funding

325 A.R.T. and C.O.B. are supported by a Maximizing Investigators' Research Award for Early Stage Inves-
326 tigators (R35 GM-124715) (<https://www.nih.gov/>). D.F.E. received financial support from Colciencias,
327 call 656-2014 "Es Tiempo de Volver" award FP44842-503-2014 (<https://minciencias.gov.co/>). T.J.C.A.
328 is supported by funds from the National Institute of Allergy and Infectious Diseases, National Institutes
329 of Health (R37 AI048071) (<https://www.niaid.nih.gov/>). This project was funded in part with federal
330 funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Depart-
331 ment of Health and Human Services, under grant number U19 AI-110818 to the Broad Institute (D.E.N.)
332 (<https://www.niaid.nih.gov/>). The funders had no role in study design, data collection and analysis, decision
333 to publish, or preparation of the manuscript.

334 Acknowledgments

335 Thank you to Pierre Jacob for guidance on the calculation of the 1-Wasserstein distances and to James
336 Watson, Manuela Carrasquilla and Vladimir Corredor for helpful comments and discussion.

337 Competing interests statement

338 The authors have declared that no competing interests exist.

339 References

- 340 [1] Dalmat R, Naughton B, Kwan-Gett TS, Slyker J, Stuckey EM. Use cases for genetic epidemiology in
341 malaria elimination. *Malaria journal*. 2019;18(1):163.
- 342 [2] Holder M, Lewis PO. Phylogeny estimation: traditional and Bayesian approaches. *Nature reviews*
343 *genetics*. 2003;4(4):275–284.
- 344 [3] Biek R, Pybus OG, Lloyd-Smith JO, Didelot X. Measurably evolving pathogens in the genomic era.
345 *Trends in ecology & evolution*. 2015;30(6):306–313.
- 346 [4] Thompson EA. Identity by descent: variation in meiosis, across genomes, and in populations. *Genetics*.
347 2013;194(2):301–326.

- 348 [5] Taylor AR, Schaffner SF, Cerqueira GC, Nkhoma SC, Anderson TJ, Sriprawat K, et al. Quantify-
349 ing connectivity between local *Plasmodium falciparum* malaria parasite populations using identity by
350 descent. *PLoS genetics*. 2017;13(10):e1007065.
- 351 [6] Grant AG, Kalisz S. Do selfing species have greater niche breadth? Support from ecological niche
352 modeling. *Evolution*. 2019;.
- 353 [7] Mattila TM, Laenen B, Slotte T. Population genomics of transitions to selfing in Brassicaceae model
354 systems. In: *Statistical Population Genomics*. Springer; 2020. p. 269–287.
- 355 [8] Siegel SV, Rayner JC. Single cell sequencing shines a light on malaria parasite relatedness in complex
356 infections. *Trends in Parasitology*. 2020;36(2):83–85.
- 357 [9] Nkhoma SC, Nair S, Cheeseman IH, Rohr-Allegrini C, Singlam S, Nosten F, et al. Close kinship within
358 multiple-genotype malaria parasite infections. *Proceedings of the Royal Society B: Biological Sciences*.
359 2012;279(1738):2589–2598.
- 360 [10] Nkhoma SC, Trevino SG, Gorena KM, Nair S, Khoswe S, Jett C, et al. Co-transmission of Related
361 Malaria Parasite Lineages Shapes Within-Host Parasite Diversity. *Cell Host & Microbe*. 2020;27(1):93–
362 103.
- 363 [11] Echeverry DF, Nair S, Osorio L, Menon S, Murillo C, Anderson TJC. Long term persistence of clonal
364 malaria parasite *Plasmodium falciparum* lineages in the Colombian Pacific region. *BMC Genetics*.
365 2013;14(2).
- 366 [12] Omedo I, Mogeni P, Rockett K, Kamau A, Hubbart C, Jeffreys A, et al. Geographic-genetic analysis of
367 *Plasmodium falciparum* parasite populations from surveys of primary school children in Western Kenya.
368 Wellcome open research. 2017;2.
- 369 [13] Omedo I, Mogeni P, Bousema T, Rockett K, Amambua-Ngwa A, Oyier I, et al. Micro-epidemiological
370 structuring of *Plasmodium falciparum* parasite populations in regions with varying transmission inten-
371 sities in Africa. Wellcome open research. 2017;2.
- 372 [14] Tessema S, Wesolowski A, Chen A, Murphy M, Wilhelm J, Mupiri AR, et al. Using parasite genetic
373 and human mobility data to infer local and cross-border malaria connectivity in Southern Africa. *Elife*.
374 2019;8:e43510.
- 375 [15] Rodríguez JCP, Uribe GÁ, Araújo RM, Narváez PC, Valencia SH. Epidemiology and control of malaria
376 in Colombia. *Memórias do Instituto Oswaldo Cruz*. 2011;106:114–122.
- 377 [16] Feged-Rivadeneira A, Ángel A, González-Casabianca F, Rivera C. Malaria intensity in Colombia by
378 regions and populations. *PLoS ONE*. 2018;13(9):e0203673. doi:10.6084/m9.figshare.6863780.
- 379 [17] Castellanos A, Chaparro-Narváez P, Morales-Plaza CD, Alzate A, Padilla J, Arévalo M, et al.
380 Malaria in gold-mining areas in Colombia. *Memorias do Instituto Oswaldo Cruz*. 2016;111(1):59–66.
381 doi:10.1590/0074-02760150382.
- 382 [18] Recht J, Siqueira AM, Monteiro WM, Herrera SM, Herrera S, Lacerda MVG. Malaria in Brazil, Colom-
383 bia, Peru and Venezuela: current challenges in malaria control and elimination. *Malaria Journal*.
384 2017;16(273):1–18. doi:10.1186/s12936-017-1925-6.
- 385 [19] Daniels JP. Increasing malaria in Venezuela threatens regional progress. *The Lancet Infectious diseases*.
386 2018;18(3):257. doi:10.1016/S1473-3099(18)30086-0.
- 387 [20] Knudson A, González-Casabianca F, Feged-Rivadeneira A, Pedreros MF, Aponte S, Olaya A, et al.
388 Spatio-temporal dynamics of *Plasmodium falciparum* transmission within a spatial unit on the Colom-
389 bian Pacific Coast. *Scientific Reports*. 2020;10(1):1–16.
- 390 [21] Grillet ME, Leopoldo V, Oletta JF, Tami A, Conn JE. Malaria in Venezuela requires response. *Science*.
391 2018;359(6375):528.

- 392 [22] Jaramillo-ochoa R, Sippy R, Farrell DF, Cueva-aponte C, Beltrán-ayala E, Gonzaga JL, et al. Effects
393 of Political Instability in Venezuela on Malaria Resurgence at Ecuador–Peru Border, 2018. *Emerging*
394 *Infectious Diseases*. 2019;25(4):834–836.
- 395 [23] Daniels JP. Venezuela in crisis. *The Lancet Infectious Diseases*. 2019;19(1):28. doi:10.1016/s1473-
396 3099(18)30745-x.
- 397 [24] Rodríguez-Morales AJ, Suárez JA, Risquez A, Villamil-Gómez WE, Paniz-Mondolfi A. Consequences
398 of Venezuela’s massive migration crisis on imported malaria in Colombia, 2016–2018. *Travel Medicine*
399 *and Infectious Disease*. 2019;28(February):98–99. doi:10.1016/j.tmaid.2019.02.004.
- 400 [25] Taylor AR, Jacob PE, Neafsey DE, Buckee CO. Estimating relatedness between malaria parasites.
401 *Genetics*. 2019; p. genetics–302120.
- 402 [26] Neafsey DE, Schaffner SF, Volkman SK, Park D, Montgomery P, Milner DA, et al. Genome-wide SNP
403 genotyping highlights the role of natural selection in *Plasmodium falciparum* population divergence.
404 *Genome biology*. 2008;9(12):R171.
- 405 [27] Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype
406 data. *Genetics*. 2000;155(2):945–959.
- 407 [28] Wille M, Holmes EC. The Ecology and Evolution of Influenza Viruses. *Cold Spring Harbor Perspectives*
408 *in Medicine*. 2019; p. a038489.
- 409 [29] Katz EM, Esona MD, Betrapally NS, Lucia A, Neira YR, Rey GJ, et al. Whole-gene analysis of inter-
410 genogroup reassortant rotaviruses from the Dominican Republic: Emergence of equine-like G3 strains
411 and evidence of their reassortment with locally-circulating strains. *Virology*. 2019;534:114–131.
- 412 [30] Caugant DA, Brynildsrud OB. *Neisseria meningitidis*: using genomics to understand diversity, evolution
413 and pathogenesis. *Nature Reviews Microbiology*. 2019; p. 1–13.
- 414 [31] Smith JM, Feil EJ, Smith NH. Population structure and evolutionary dynamics of pathogenic bacteria.
415 *Bioessays*. 2000;22(12):1115–1122.
- 416 [32] Tibayrenc M, Ayala FJ. The clonal theory of parasitic protozoa: 12 years on. *Trends in parasitology*.
417 2002;18(9):405–410.
- 418 [33] Rajendran C, Su C, Dubey JP. Molecular genotyping of *Toxoplasma gondii* from Central and South
419 America revealed high diversity within and between populations. *Infection, Genetics and Evolution*.
420 2012;12(2):359–368.
- 421 [34] Nader JL, Mathers TC, Ward BJ, Pachebat JA, Swain MT, Robinson G, et al. Evolutionary genomics
422 of anthroponosis in *Cryptosporidium*. *Nature microbiology*. 2019;4(5):826–836.
- 423 [35] Nieuwenhuis BP, James TY. The frequency of sex in fungi. *Philosophical Transactions of the Royal*
424 *Society B: Biological Sciences*. 2016;371(1706):20150540.
- 425 [36] Fruchterman TM, Reingold EM. Graph drawing by force-directed placement. *Software: Practice and*
426 *experience*. 1991;21(11):1129–1164.
- 427 [37] Falush D, Stephens M, Pritchard JK. Inference of population structure using multilocus genotype data:
428 linked loci and correlated allele frequencies. *Genetics*. 2003;164(4):1567–1587.
- 429 [38] Huestis DL, Dao A, Diallo M, Sanogo ZL, Samake D, Yaro AS, et al. Windborne long-distance migration
430 of malaria mosquitoes in the Sahel. *Nature*. 2019;574(7778):404–408.
- 431 [39] Verdonshot PF, Besse-Lototskaya AA. Flight distance of mosquitoes (Culicidae): a metadata anal-
432 ysis to support the management of barrier zones around rewetted and newly constructed wetlands.
433 *Limnologica-Ecology and Management of Inland Waters*. 2014;45:69–79.

- 434 [40] Pacheco MA, Schneider KA, Céspedes N, Herrera S, Arévalo-Herrera M, Escalante AA. Limited
435 differentiation among *Plasmodium vivax* populations from the northwest and to the south Pacific
436 Coast of Colombia: A malaria corridor? *PLoS Neglected Tropical Diseases*. 2019;13(3):e0007310.
437 doi:10.1371/journal.pntd.0007310.
- 438 [41] Guagliardo SA, Morrison AC, Barboza JL, Requena E, Astete H, Vazquez-Prokopec G, et al. River
439 boats contribute to the regional spread of the dengue vector *Aedes aegypti* in the Peruvian Amazon.
440 *PLoS neglected tropical diseases*. 2015;9(4):e0003648.
- 441 [42] Lounibos LP. Invasions by insect vectors of human disease. *Annual review of entomology*.
442 2002;47(1):233–266.
- 443 [43] Montoya-Lerma J, Solarte YA, Giraldo-Calderón GI, Quiñones ML, Ruiz-López F, Wilkerson RC, et al.
444 Malaria vector species in Colombia: a review. *Memórias do Instituto Oswaldo Cruz*. 2011;106:223–238.
- 445 [44] Gutiérrez LA, Naranjo NJ, Cienfuegos AV, Muskus CE, Luckhart S, Conn JE, et al. Population structure
446 analyses and demographic history of the malaria vector *Anopheles albimanus* from the Caribbean and
447 the Pacific regions of Colombia. *Malaria journal*. 2009;8(1):259.
- 448 [45] Instituto Nacional de Salud Colombia, Dirección de Vigilancia y Analisis del Riesgo en Salud Pública.
449 Boletín Epidemiológico Semanal: semana epidemiológica 52. 2019;doi:10.33610/23576189.2019.52.
- 450 [46] Wabgou M, Vargas D, Carabali JA. Las migraciones internacionales en Colombia. *Investigación &*
451 *Desarrollo*. 2012;20(1):142–167.
- 452 [47] Vera-Arias CA, Castro LE, Gómez-Obando J, Sáenz FE. Diverse origin of *Plasmodium falciparum* in
453 northwest Ecuador. *Malaria journal*. 2019;18(1):251.
- 454 [48] Ahumada ML, Orjuela LI, Pareja PX, Conde M, Cabarcas DM, Cubillos EFG, et al. Spatial distributions
455 of *Anopheles* species in relation to malaria incidence at 70 localities in the highly endemic Northwest
456 and South Pacific coast regions of Colombia. *Malaria Journal*. 2016;15(407):1–16. doi:10.1186/s12936-
457 016-1421-4.
- 458 [49] Shetty AC, Jacob CG, Huang F, Li Y, Agrawal S, Saunders DL, et al. Genomic structure and di-
459 versity of *Plasmodium falciparum* in Southeast Asia reveal recent parasite migration patterns. *Nature*
460 *communications*. 2019;10(1):2665.
- 461 [50] Naranjo-Díaz N, Altamiranda M, Luckhart S, Conn JE, Correa MM. Malaria vectors in eco-
462 logically heterogeneous localities of the Colombian Pacific region. *PLoS ONE*. 2014;9(8):e103769.
463 doi:10.1371/journal.pone.0103769.
- 464 [51] Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal*.
465 2006;Complex Systems:1695.
- 466 [52] R Core Team. R: A Language and Environment for Statistical Computing; 2018. Available from:
467 <https://www.R-project.org/>.
- 468 [53] Schuhmacher D, Bähre B, Gottschlich C, Hartmann V, Heinemann F, Schmitzer B. transport: Com-
469 putation of Optimal Transport Plans and Wasserstein Distances; 2019. Available from: [https://cran.r-](https://cran.r-project.org/package=transport)
470 [project.org/package=transport](https://cran.r-project.org/package=transport).
- 471 [54] Diaz G, Lasso AM, Murillo C, Montenegro LM, Echeverry DF. Evidence of self-medication with chloro-
472 quine before consultation for malaria in the southern pacific coast region of Colombia. *American Journal*
473 *of Tropical Medicine and Hygiene*. 2019;100(1):66–71. doi:10.4269/ajtmh.18-0515.
- 474 [55] Valencia SH, Ocampo ID, Arce-Plata MI, Recht J, Arévalo-Herrera M. Glucose-6-phosphate dehydroge-
475 nase deficiency prevalence and genetic variants in malaria endemic areas of Colombia. *Malaria Journal*.
476 2016;15(291):1–9. doi:10.1186/s12936-016-1343-1.

- 477 [56] Vallejo AF, Chaparro PE, Benavides Y, Álvarez Á, Quintero JP, Padilla J, et al. High prevalence of
478 sub-microscopic infections in Colombia. *Malaria Journal*. 2015;14(201):1–7. doi:10.1186/s12936-015-
479 0711-6.
- 480 [57] Valero-Bernal MV, Tanner M, Muñoz-Navarro S, Valero-Bernal JF. Proportion of fever attributable to
481 malaria in Colombia: Potential indicators for monitoring progress towards malaria elimination. *Revista*
482 *de Salud Pública*. 2017;19:45–51.
- 483 [58] Pava Z, Echeverry DF, Díaz G, Murillo C. Large variation in detection of histidine-rich protein 2 in
484 *Plasmodium falciparum* isolates from Colombia. *The American journal of tropical medicine and hygiene*.
485 2010;83(4):834–837.
- 486 [59] Solano CM, Okoth SA, Abdallah JF, Pava Z, Dorado E, Incardona S, et al. Deletion of *Plasmodium*
487 *falciparum* histidine-rich protein 2 (pfrp2) and histidine-rich protein 3 (pfrp3) genes in Colombian
488 parasites. *PloS one*. 2015;10(7):e0131576.
- 489 [60] Dorado EJ, Okoth SA, Montenegro LM, Diaz G, Barnwell JW, Udhayakumar V, et al. Genetic character-
490 isation of *Plasmodium falciparum* isolates with deletion of the pfrp2 and/or pfrp3 genes in Colombia:
491 the Amazon region, a challenge for malaria diagnosis and control. *PLoS One*. 2016;11(9):e0163137.
- 492 [61] Padilla JC, Chaparro PE, Molina K, Arevalo-Herrera M, Herrera S. Is there malaria transmission in
493 urban settings in Colombia? *Malaria Journal*. 2015;14(453):1–9. doi:10.1186/s12936-015-0956-0.

494 Appendix A

Climate e.g. very heavy rainfall and flooding on the Pacific Coast [44, 50]
 Rich vector diversity [48, 50]†
 Parasite resistance to antimalarial drugs [18, 20]
 Persistent drug pressure due to frequent self-medication [54]
 Prevalence of glucose-6-phosphate dehydrogenase deficiency [18, 55]†
 Prevalence of duffy-negative individuals resistant to *P. vivax* invasion [18, 56]†
 Sub-microscopic and asymptomatic infections [18, 20, 56, 57]† and infections that can evade detection by some rapid diagnostic tests [18, 20, 58–60]
 Per-urban and urban transmission [18, 61]
 Human migration due to deforestation and gold mining [16–19]

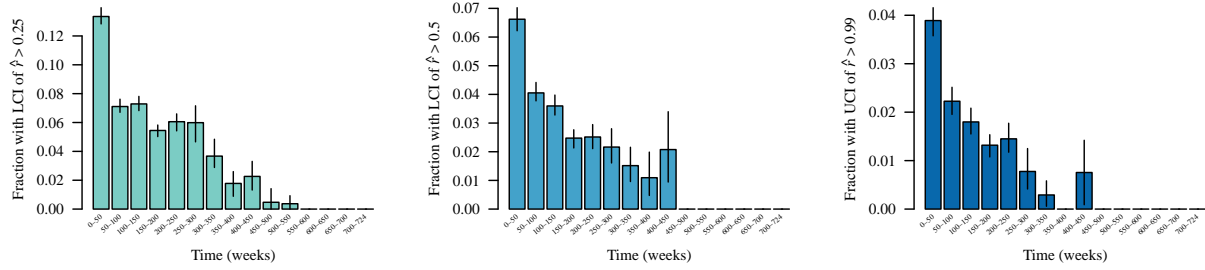
Table A.1: Factors associated with malaria epidemiology in Colombia: a non exhaustive list. †One or more citations are specific to Tierralta and the (South) Pacific coast.

City (Province)	1993	1994	1997	1999	2000	2001	2002	2003	2004	2005	2006	2007	Total
Tumaco (Nariño)	0	0	0	2	2	10	11	59	0	23	0	25	132
Guapi (Cauca)	0	0	0	1	1	0	0	66	0	0	0	0	68
Buenaventura (Valle)	4	1	0	5	0	0	0	0	12	15	10	0	47
Quibdó (Chocó)	0	0	2	0	6	1	0	0	14	6	13	22	64
Tadó (Chocó)	0	0	0	0	0	12	2	0	0	0	0	0	14
Total	4	1	2	8	9	23	13	125	26	44	23	47	325

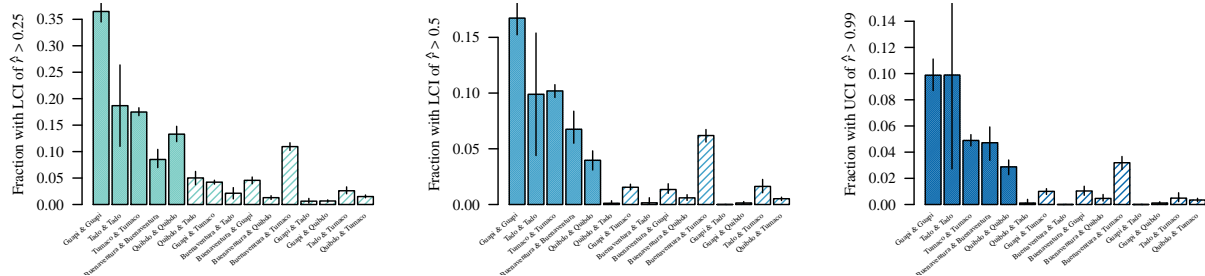
Table A.2: Yearly monoclonal *P. falciparum* sample counts per city.

	CC1	CC12	CC14	CC20
CC12	0.712 (0.635)			
CC14	0.733 (0.709)	0.648 (0.517)		
CC20	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	
CC40	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.079 (0.000)

Table A.3: Average relatedness to three decimal places between clonal components (CCs) 1, 12, 14, 20 and 40 with maximum 2.5% confidence interval bound in parentheses. The maximum 2.5% confidence interval bound indicates that relatedness between C20 and C40 is not statistically distinguishable from zero, for example.



(a) Partitioned by time between collection dates.



(b) Partitioned by collection city.

Figure A.1: Fractions of highly-related sample pairs partitioned in time and space: sensitivity to high-relatedness thresholds. highly-related samples pairs are defined as those with lower confidence interval bound (LCI) of relatedness estimate, \hat{r} , greater than thresholds 0.25 and 0.50; or with upper confidence interval bound (UCI) of $\hat{r} > 0.99$ (i.e. clonal parasite sample pairs; Table 1). Colours correspond to Figure 1.

Clonal component	Date	City
CC1	1999-03-15	Guapi
CC2	1999-04-13	Buenaventura
CC3	2000-04-13	Quibdo
CC4	2000-06-29	Quibdo
CC5	2000-11-23	Tumaco
CC6	2001-01-29	Tumaco
CC7	2001-02-08	Tumaco
CC8	2001-02-17	Tumaco
CC9	2001-06-07	Tado
CC10	2001-06-08	Tado
CC11	2001-12-03	Tado
CC12	2002-04-03	Tumaco
CC13	2002-04-03	Tumaco
CC14	2002-04-04	Tumaco
CC15	2002-04-05	Tumaco
CC16	2003-01-07	Guapi
CC17	2003-03-03	Guapi
CC18	2003-03-07	Guapi
CC19	2003-03-17	Tumaco
CC20	2003-03-22	Guapi
CC21	2003-04-11	Guapi
CC22	2003-05-15	Tumaco
CC23	2003-05-20	Tumaco
CC24	2003-05-29	Tumaco
CC25	2003-09-08	Tumaco
CC26	2003-10-01	Tumaco
CC27	2003-10-03	Guapi
CC28	2003-10-06	Guapi
CC29	2003-10-07	Guapi
CC30	2003-10-21	Tumaco
CC31	2003-11-05	Guapi
CC32	2004-05-31	Buenaventura
CC33	2004-07-23	Buenaventura
CC34	2004-08-24	Quibdo
CC35	2004-10-27	Quibdo
CC36	2004-11-10	Quibdo
CC37	2004-12-10	Quibdo
CC38	2005-02-28	Quibdo
CC39	2005-07-27	Tumaco
CC40	2005-07-28	Tumaco
CC41	2006-08-03	Quibdo
CC42	2006-08-22	Buenaventura
CC43	2006-09-05	Quibdo
CC44	2006-10-10	Quibdo
CC45	2007-03-22	Quibdo
CC46	2007-03-28	Quibdo

Table A.4: Date and city of collection of earliest parasite sample per clonal component.

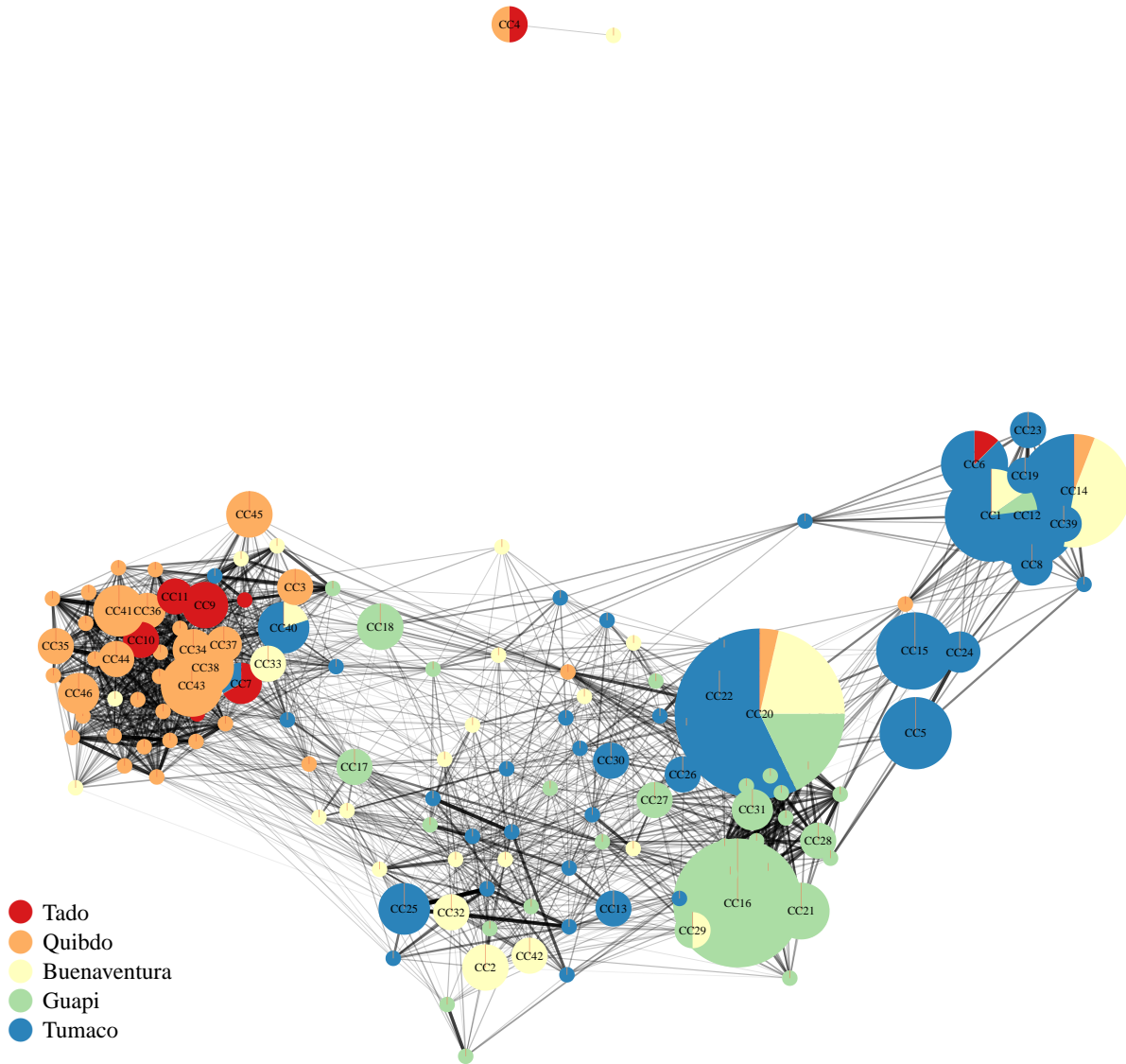


Figure A.2: Clonal components and singletons (vertices) and the average relatedness between them (edges). Clonal components (CCs) are groups of two or more statistically indistinguishable parasite samples. Singletons are individual parasite samples that do not belong to a CC. Vertices are plotted using the Fruchterman-Reingold layout algorithm [36], thereby clustering inter-related vertices. The size of each CC vertex is proportional to the number of parasite samples per CC, ranging from 2 to 28 statistically indistinguishable parasite samples. CCs are named in order of the collection date of the earliest parasite sample per CC (Table A.4). CCs with parasite samples from two or more sites are depicted as pie charts. Colour denotes the city of parasite sample collection. Edge transparency and weight is proportional to average relatedness, ranging from 0.003 to 0.912. Relatedness estimates that are indistinguishable from zero were set to zero. Edges whose average relatedness is zero are not plotted. Each CC besides CC4 is related to at least one other. CC4 is likely a contaminant; see main text. A singleton from Buenaventura, which is loosely related to CC4, may also be a contaminant.

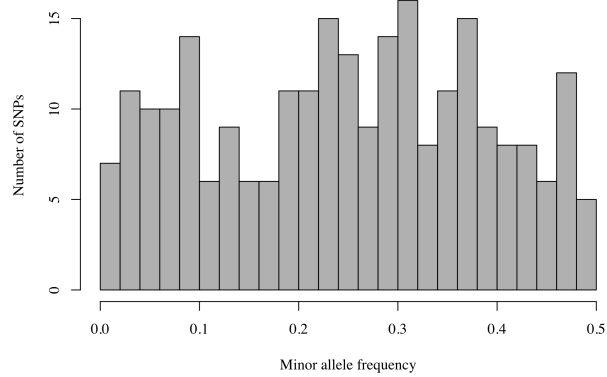


Figure A.3: Histogram of minor allele frequencies estimated using all 325 monoclonal *P. falciparum* samples genotyped at 250 biallelic SNPs.