1    The impact of host metapopulation structure on the population genetics of colonizing bacteria

2    Elina Numminen[1], Michael Gutmann[1,2], Mikhail Shubin[1], Pekka Marttinen[2], Guillaume Méric[3],

3    Willem van Schaik[4], Teresa M. Coque[5], Fernando Baquero[5], Rob J. L. Willems[4], Samuel K.

4    Sheppard[3], Edward J. Feil[6], William P. Hanage[7], Jukka Corander[1]

5

6    [1]Department of Mathematics and Statistics, University of Helsinki, Helsinki, Finland; [2]Helsinki

7    Institute for Information Technology HIIT, Department of Information and Computer Science,

8    Aalto University, Finland; [3]College of Medicine, Swansea University, Institute of Life Science,

9    Swansea, UK; [4]Department of Medical Microbiology, University Medical Center Utrecht, Utrecht,

10   The Netherlands; [5]Department of Microbiology, Ramón y Cajal University Hospital, Madrid,

11   Spain; [6]Department of Biology and Biochemistry, University of Bath, Claverton Down, Bath, UK;

12   [7]Center for Communicable Disease Dynamics, Harvard School of Public Health, Boston,

13   Massachusetts, USA;

14

15   Abstract

16   Many key bacterial pathogens are frequently carried asymptomatically, and the emergence and
17   spread of these opportunistic pathogens can be driven, or mitigated, via demographic changes
18   within the host population. These inter-host transmission dynamics combine with basic
19   evolutionary parameters such as rates of mutation and recombination, population size and selection,
20   to shape the genetic diversity within bacterial populations. Whilst many studies have focused on
21   how molecular processes underpin bacterial population structure, the impact of host migration and
22   the connectivity of the local populations has received far less attention. A stochastic neutral model
23   incorporating heightened local transmission has been previously shown to fit closely with genetic
24   data for several bacterial species. However, this model did not incorporate transmission limiting
25   population stratification, nor the possibility of migration of strains between subpopulations, which
26   we address here by presenting an extended model. The model captures the observed population
27   patterns for the common nosocomial pathogens *Staphylococcus epidermidis* and *Enterococcus*
28   *faecalis*, while *Staphylococcus aureus* and *Enterococcus faecium* display deviations attributable to
29   adaptation. It is demonstrated analytically and numerically that expected strain relatedness may
30   either increase or decrease as a function of increasing migration rate between subpopulations, being
31   a complex function of the rate at which microepidemics occur in the metapopulation. Moreover, it
32   is shown that in a structured population markedly different rates of evolution may lead to
33   indistinguishable patterns of relatedness among bacterial strains; caution is thus required when
34   drawing evolution inference in these cases.

35   Key words: Bacterial evolution, genetic structure, migration, population dynamics

36

37     Introduction

38     Bacteria colonizing multicellular hosts are organized in a hierarchy of local interconnected
39     subpopulations forming a complex metapopulation as a whole. The subpopulations can range in
40     scale from discrete intracellular colonies residing within a single host cell to pervasive strains
41     circulating among hosts across cities, countries and continents(Fraser et al., 2009). Although most
42     bacteria are harmless or even advantageous to their host organisms, some cause infectious disease,
43     and understanding the evolutionary dynamics and the factors producing the genetic variation of
44     pathogen populations is important for combatting disease emergence and spread.

45     Previous work has demonstrated that a simple model of stochastic microepidemics arising from
46     repeated sampling of localized transmission chains, can explain genotypic variation in local
47     surveillance data from several common human pathogens(Fraser et al., 2005; Hanage et al., 2006),
48     under an assumption that all isolates are equally fit (neutrality). In these studies, populations were
49     characterized by a simple measure of the level of genotype relatedness known as the allelic
50     mismatch distribution, where isolates with more shared alleles are considered to be more closely
51     related. These comparisons have been widely used in classical ecology and population genetics and
52     different patterns in the mismatch distribution can be associated with various factors contributing to
53     the population structure, including: population growth(Harpending, 1994; Rogers and Harpending,
54     1992), selection(Bamshad et al., 2002), and host contact network structure(Plucinski et al., 2011).
55     The mismatch distribution has also been used to detect deviations from neutrality or constant
56     population size(Mousset et al., 2004) and for inference about bacterial recombination rates(Hudson,
57     1987).

58     Population structure is one of the most studied phenomena in population genetics, both from the
59     theoretical and applied perspective(Ewens, 2004; Hartl and Clark, 2007). Nevertheless in the case
60     of bacteria limited knowledge exists about the effects of population structure arising from multiple
61     host organisms such as human and different animal species or other, often poorly defined and
62     understood, ecological patches. The main reason for this is simultaneously accounting for the major
63     phenomena known to impact evolution of bacterial pathogen populations, such as recombination,
64     clonal expansion, as well as migration, which for example may be caused by anthroponosis and
65     zoonosis when multiple different host organisms are colonized by the same bacterial species. This
66     hampers both theoretical derivation of limit results for such models and empirical fitting due to
67     likelihood equations not being available in closed form. Fraser et al. solved the likelihood
68     intractability arising from microepidemics by using a stochastic mixture distribution to account for
69     the increase in the probability of sampling identical strains from the same transmission chain
70     (Fraser et al., 2005). An analogous approximation technique has later been independently
71     introduced in a more general ecological setting and it is known as the *synthetic likelihood* (Wood,
72     2010*)*.

73     To improve understanding of the evolutionary dynamics of structured bacterial populations, we
74     employ a simulation-based approach to neutral models that can account for the multiple stochastic
75     forces impacting the genetic diversity that persists over time. By capturing both a heterogeneous
76     span of microepidemics and migration events across the boundaries limiting transmission between
77     subpopulations, we characterize the expected behavior of the metapopulations as a whole. This

78  provides an opportunity to explore the limits of inferring the vital model parameters from genetic
79  surveillance data, and gives novel insight into the emergence of important human pathogens.

80  Materials and Methods

81  Model

82  We consider an infinite alleles model for a finite haploid population with $N$ individuals and discrete
83  generations, where the reproduction takes place by random sampling of $N$ individuals from the
84  current generation to the next generation(Ewens, 2004). When the population is assumed structured,
85  the subpopulation sizes are indexed by $N_1$, $N_2$. The parameters which may vary across
86  subpopulations are indexed accordingly. Mutations are introduced per generation by a Poisson
87  process with the rate $\theta = \mu N \tau$, where $\mu$ is the per locus mutation rate and $\tau$ is a scaling factor
88  representing the generation time in calendar time. In all subsequent work we set $\tau = 1$, unless
89  otherwise mentioned. We assume that each individual is characterized by a genotype comprising
90  alleles at $L$ unlinked loci, where a mutation event at any locus always introduces a novel allele.
91  Recombination between randomly chosen genotypes occurs at any locus according to a Poisson
92  process with the rate defined as $\rho = r N \tau$, where $r$ is the rate per locus in relation to the mutation rate.
93  In our simulations we simulated the population until allelic diversity reached equilibrium.

94  Microepidemics are modeled as doubly stochastic events, with the frequency of new
95  microepidemics per generation following a Poisson distribution with mean $\omega N \tau$. The size of each
96  microepidemic has a Poisson distribution with mean $\gamma$. Each micropidemic is generated
97  independently similar to the assumptions in Fraser et al. such that first a single individual is
98  randomly chosen, after which its genotype is propagated to $Y$ randomly chosen other individuals
99  such that $Y$ has Poisson distribution with mean $\gamma$. When the population is stratified, the
100 microepidemic rates of the subpopulations are denoted by $\omega_1$, $\gamma_1$ and $\omega_2$, $\gamma_2$, respectively. Migration
101 between subpopulations is a Poisson process with the rates $\tau N_1 m_{12}$, $\tau N_2 m_{21}$ per generation, where the
102 first subindex of the parameters $m_{12}$, $m_{21}$ defines the source and the second subindex the target
103 subpopulation. In migration events genotypes of a Poisson distributed number of randomly chosen
104 individuals from the source population replace the genotypes of randomly chosen individuals in the
105 target population. In our simulations the events were generated in the following order: reproduction
106 mutation, recombination, microepidemics and migration at each generation. In all the reported
107 results each subpopulation size was $N = 2000$, unless otherwise indicated. Medians and 95%
108 confidence intervals for the allelic mismatch distributions were obtained by recording the
109 population state every 100th generation after initial 500 generations until 20000 generations, and
110 using these values to calculate the corresponding quantiles of the mismatch probabilities.

111 Data and processing of genotype networks

112 eBURST networks for the populations were produced using default settings(Feil et al., 2004).
113 Turner et al. demonstrated that eBURST provides a robust recapitulation of the genetic relatedness
114 of strains in a bacterial population based on the MLST resolution(Turner et al., 2007). To quantify
115 details of the networks we calculated genotype degree distributions and distributions of geodesic
116 distances between pairs of genotypes, which are standard measures of network topology(Goh et al.,
117 2002).

118    MLST isolate data were accessed (September 15, 2014) from the following databases:
119    http://efaecalis.mlst.net/ (*E. faecalis*), http://efaecium.mlst.net/ (*E. faecium*), http://saureus.mlst.net/
120    (*S. aureus*), and (May 10, 2015) from: http://sepidermidis.mlst.net/ (*S. epidermidis*).

121

122

123    Results

124    We extended the microepidemic infinite alleles model with mutation and recombination rates
125    previously proposed by Fraser et al. (Fraser et al., 2005) to incorporate population stratification,
126    whereby genotypes are free to move between subpopulations at a defined rate. In addition, rather
127    than using a single microepidemic parameter to describe localized transmission (Fraser et al., 2005),
128    we introduced two parameters modulating the distributions of both the frequency and sizes of the
129    transmission clusters in stochastic fashion. Our microepidemic infinite alleles migration model
130    (MIAMI) can thereby encompass a wide variety of evolutionary and ecological parameter space.
131    Since the resulting patterns of genetic variation reflect a complex function of several factors, we
132    consider first a model without population stratification to delineate the influence of each of the
133    model components.

134    The frequency distribution of the number of allelic mismatches between pairs of genotypes is a
135    classical approach to describe the distribution of genetic variation within a population (Fraser et al.,
136    2005). Depending on the interplay of several factors, a population may either have a peaked or flat
137    equilibrium distribution over the space of summary statistics, such as the allelic mismatch
138    distribution (Fig. 1). For lower mutation rates, high $r/m$ will lead to bell-shaped mismatch
139    distributions, since recombination acts as a cohesive force keeping genetic variation together as a
140    cloud in the space of possible genotypes(Fraser et al., 2007). The mismatch distribution becomes
141    less sensitive to changes in recombination rate and the equilibrium distribution becomes more
142    peaked when the mutation rate increases (Fig. 1).

143    Fig. 2 shows the impact of heightened localized transmission (microepidemics) on genetic
144    relatedness visualized using eBURST (Feil et al., 2004; Francisco et al., 2009) and the allele
145    mismatch distribution. The rate of mutation and homologous recombination varies among bacterial
146    pathgoens and this can have a marked effect on the population structure. To model the interplay of
147    these two important factors at different levels, four evolutionary scenarios were considered: low
148    mutation and recombination rate (A), mutation dominates (B), recombination dominates (C), both
149    mutation and recombination effects are sizeable (D). If mutation dominates over recombination
150    (Fig. 2,B), microepidemics do not lead to as pronounced changes in the relatedness pattern as in the
151    situation where both mutation and recombination rates are low (Fig. 2,A). Interconnected clusters
152    do emerge under a high rate of recombination, often spanning across large parts of the entire
153    population (Fig. 2,C). The variability of the mismatch distribution at the equilibrium becomes
154    elevated under all regimes of baseline parameter values when microepidemics occur at a frequent
155    rate, as illustrated by the broader confidence intervals (Fig. 2,A-D). Both the frequency and size
156    distribution of the individual microepidemics influence how much probability mass is shifted

157    towards identical genotypes, but the change is also influenced by mutation and recombination rate
158    parameters (Supplementary Fig. 1).

159    The effect of migration rate on the allelic mismatch distribution within a subpopulation is a
160    complicated function of mutation, recombination and microepidemic rates in a structured
161    population, even if there are only two subpopulations (Fig. 3). We studied the combinations in
162    which a subpopulation undergoes microepidemic expansions at a moderate rate and is coupled with
163    another subpopulation where the rate varies from zero to twice that of the first subpopulation. An
164    increase of the migration rate between the two subpopulations by an order of magnitude leads either
165    to a substantial decrease of the genotypic diversity (Supplementary Fig. 2, i), an increase in the
166    genotypic diversity (Supplementary Fig. 2, a), or to no change at all (Supplementary Fig. 2, e),
167    depending on whether the subpopulation considered as a source experiences more, less, or an equal
168    amount of the microepidemics, compared with the target subpopulation. The effect of migration
169    remains equally complex for the between-subpopulations allelic mismatch distribution, which is
170    insensitive to a change in the migration rate by an order of magnitude for many combinations of
171    subpopulation dynamics (Supplementary Fig. 3). Population stratification combined with
172    asymmetric migration rates can produce patterns of relatedness which are otherwise unlikely under
173    the neutral model (Supplementary Fig. 4). For example, in all our simulations a characteristic U-
174    shaped allelic mismatch distribution only arose when the migration rate was highly asymmetric and
175    one subpopulation experienced considerable microepidemics while the other one had none
176    (Supplementary Figs. 5,6).

177    To obtain an analytical insight to the joint effect of microepidemic and migration rates on genotypic
178    diversity, we considered how the equilibrium probability of identical genotypes is affected by
179    introducing a change to the subpopulation based on either mechanism. Fraser et al. derived the
180    equilibrium probability of identical genotypes at $L$ unlinked loci, under the assumption of no
181    microepidemics(Fraser et al., 2005), which equals $p_0^L = \frac{1 + L\rho p_0^{L-1} p_0^1}{1 + L\theta + L\rho}$. Here $\theta = 2\mu N$, where $\mu$ is the
182    per locus mutation rate and $N$ is the population size. Furthermore, the recombination rate is defined
183    as $\rho = 2rN$, where $r$ is the rate per locus in relation to the mutation rate. Since this extension of the
184    classical equilibrium result by Kimura to allow for recombination is based on the assumption that in
185    any generation only a single event occurs, Fraser et al. handled the effect of microepidemics on a
186    population at equilibrium implicitly by introducing a probabilistic mixture where a single parameter
187    represents the increase in the probability $p_0^L$ caused by microepidemics. Consistent with this, we
188    quantify the change in the probability of identical strains by evaluating the expectation of the effect
189    of microepidemic and migration events when allowed at the equilibrium of a simpler population
190    experiencing only mutation and recombination events.

191    Consider first the effect of stochastic microepidemics occurring in a single generation. The expected
192    number of identical genotype pairs arising from them equals $(\gamma + 1)^2 N\omega$ , where $\omega$ is the scaled
193    rate at which microepidemics occur per generation and $\gamma$ is the expected size of each microepidemic
194    (Methods). The expected contribution to the probability of homozygous strains is then $\frac{(\gamma+1)^2 N\omega}{\binom{N}{2}}$,
195    which is an increasing function of both the expected size and rate of microepidemics. Next,
196    consider two subpopulations of sizes $N_1, N_2$, which at equilibrium become connected with migration

197  rates $N_1 m_{12}$, $N_2 m_{21}$, respectively, in addition to the effect of introducing microepidemics (Methods).
198  Each subpopulation is assumed to have its own set of parameters $\gamma_1^2 N_1 \omega_1$, $\gamma_2^2 N_2 \omega_2$ governing the
199  extent of microepidemics. Assume now that the subpopulations are of equal size $N_1 = N_2$. Then, the
200  expected contribution to the probability of identical strains in subpopulation 1 by an increase in the
201  migration rate $m_{21}$ depends on whether $\gamma_1^2 N_1 \omega_1 > \gamma_2^2 N_2 \omega_2$ or $\gamma_1^2 N_1 \omega_1 < \gamma_2^2 N_2 \omega_2$, since larger and
202  more frequent microepidemics in subpopulation 2 will increase the probability that the genotypes
203  migrating to subpopulation 1 are identical to each other. Conversely, increased migration from
204  subpopulation 2 will have expected effect of decreasing the probability when the extent of
205  microepidemics in subpopulation 2 is smaller than in subpopulation 1. A difference in the sizes of
206  the subpopulations can further amplify these effects since the rates of events are relative to them.

207  Global surveillance data based on MLST typing for several common nosocomial bacterial
208  pathogens (*S. aureus, S. epidermidis, E. faecalis, E. faecium*) generally match well with the
209  expected shape of the allelic mismatch distribution for the considered archetypical population types
210  (Fig. 4). eBURST diagrams provide additional insight into the structure of these populations (Fig.
211  5). *S. aureus* is known to have a very low recombination rate(Everitt et al., 2014) and its population
212  structure is mainly shaped by a combination of mutation rate and intensive clonal expansion of
213  distinct genotypes (Fig. 5, C). Conversely, its sister species *S. epidermidis* displays the bell-shaped
214  mismatch distribution typical for organisms with high recombination rate(Meric et al., 2015) (Fig.
215  4, D) and a large connected network of related genotypes (Fig. 5, D). The numerous distinct clusters
216  with short distances to the ancestral genotype observed in *S. aureus* population (clonal complexes
217  with single-locus variants) were not accurately predicted by the model, despite of an extensive
218  search over the parameter space. The main deviance arose from the inability to recapitulate a large
219  number of descendant genotypes connected with each single ancestral genotype. The most closely
220  matching neutral model predicts instead invariably that several further branches emerge from these
221  descendants during the timescale at which genotype clusters themselves emerge.

222  Contrasting the population structures of *E. faecium* and *E. faecalis* reveals marked differences,
223  where *E. faecium* forms large networks of related genotypes characteristic of highly
224  recombinogenic bacteria (Fig. 5, B) (Turner et al., 2007), despite a relatively low estimated
225  recombination rate(de Been et al., 2013). *E. faecalis* shows only limited clustering of genotypes
226  (Fig. 5, A) and a mismatch distribution typical for a population dominated by mutation, with a
227  slight increase of identical genotype pairs due to localized hospital transmission (Fig. 4, A).

228  The model parameter configurations leading to matching characteristics between the observed and
229  simulated population structure are given in Table 1 for the two species where the neutral model re-
230  capitulates the surveillance data well (*S. epidermidis, E. faecalis*). We compared genotype networks
231  using the standard measures of degree distribution and geodesic distances between nodes and found
232  a considerable agreement between the data and the simulations (Table 1, Supplementary Fig. 7,8).

233  Table 1. Population characteristics of genotype relatedness for real and simulated data.

|  | *S. aureus* | *S. epidermidis* | *E. faecalis* | *E. faecium* |
|---|---|---|---|---|
| *N* commensal | 555 | 120 | 225 | 126 |
| *N* hospital | 543 | 264 | 1003 | 1534 |
| Mean degree | 3.14 | 1.95 | 1.12 | 4.03 |

| Max degree | 34 | 10 | 9 | 47 |
|---|---|---|---|---|
| Mean geodesic distance | 2.15 | 3.54 | 1.82 | 4.12 |
| Max geodesic distance | 5 | 7 | 5 | 12 |
| Simulation settings | Not matching | $\omega_1 = 45$, $\gamma_1 = 30$, $\omega_2 = 10$ , $\gamma_2 = 10$ , $m_{12} = 0.01$ , $m_{21} = 0.001$ , $\theta = 0.0704$ , $r/m = 2$ | $\omega_1 = 10$ , $\gamma_1 = 20$, $\omega_2 = 15$ , $\gamma_2 = 20$, $m_{12} = 0.001$, $m_{21} = 0.001$, $\theta = 0.198$, $r/m = 1$ | Not matching |
| Mean degree | | 1.67 | 1.08 | |
| Max degree | | 20 | 13 | |
| Mean geodesic distance | | 3.1 | 1.94 | |
| Max geodesic distance | | 8 | 6 | |

234

235     Discussion

236     Previously described neutral models specified by mutation and recombination rate in combination
237     with microepidemics show a close fit to observed genotype survey data for several commensal and
238     pathogenic bacteria. This holds true for both short-term population evolution dominated by the local
239     dynamics of microepidemics (Fraser et al., 2005; Hanage et al., 2006) and for longer time scales
240     where recombination acts as a cohesive force keeping populations together(Fraser et al., 2007).
241     However, there is limited knowledge about how varying levels of isolation in host organisms, such
242     as human and different animal species (Fraser et al., 2009), might influence the evolutionary
243     dynamics and lead to structured populations. Here we introduce a neutral model incorporating
244     microepidemics and migration, which mimics a situation where ecological factors limit
245     transmission between subpopulations. By comparing the model predictions with MLST data large
246     scale genotyping surveys of four major human pathogens we find that for two species the
247     population structure is well delineated by the neutral assumptions, while different types of
248     deviations from the model predictions are observed for the remaining two.

249

250     The observed differences between *E. faecium* and *E. faecalis*, which colonize the gastrointestinal
251     tract, are particularly interesting since mutation and recombination rates have been estimated to be
252     similar for the two species based on both MLST and whole-genome data(de Been et al., 2013; Vos
253     and Didelot, 2009). Moreover, they are responsible for roughly equal frequencies of nosocomial
254     infections worldwide (Tedim et al., 2015; Willems et al., 2012). *E. faecalis* population structure
255     bears the hallmarks of either a high rate of mutation or drift (or both). *E. faecalis* is known to
256     colonize the vast majority of normal hosts within a population (Tedim et al., 2015), and therefore
257     can be considered as part of the physiological commensal microbiota of humans and many other
258     animals. Certainly, its population structure could be reflective of the evolutionary dynamics of a
259     generalist organism which regularly experiences a high level of drift and gene flow between
260     different host species.

261     On the basis of the predictions made by our model, *E. faecium* would need to have substantially
262     higher recombination rate than *E. faecalis* to lead to the observed pattern of genotype relatedness
263     under neutrality. Since there is evidence of the recombination rate not being substantially higher in
264     *E. faecium*, the only possibility for the large genotype networks to arise under our neutral model
265     would be unobserved population stratification. If unobserved sources experiencing very large clonal
266     expansions contributed continuously to the hospital subpopulation of *E. faecium*, the expected
267     allelic mismatch distribution would bear the characteristics of a subpopulation with high
268     recombination rate (Supplementary Fig. 3, i). It is known that intensive farming and animal
269     production practices provide opportunities for rapid clonal expansion of bacterial strains colonizing
270     the animal hosts. Given the known connection between strains from domesticated animals and the
271     hospital associated *E. faecium* (Lebreton et al., 2013; Willems et al., 2012), it is plausible that these
272     clonal expansions could manifest themselves as connected networks in the human hospital
273     subpopulation. However, the extensively connected network of *E. faecium* genotypes would still
274     remain unlikely unless the rate of recombination was substantial. An alternative explanation for the
275     extensive genotype relatedness is a marked deviation from neutrality, such that the connected
276     strains represent either a subpopulation adapted to the hospital environment, consistent with
277     previous studies(Lebreton et al., 2013; Willems et al., 2012), or an adaptation to different host
278     subpopulations (Faith et al., 2015). Further dense sampling will be required to characterize
279     mechanistically the role of hospital adaption for creating the observed relatedness patterns of *E.*
280     *faecium* strains.

281     *S. aureus* and *S. epidermidis* frequently colonize the skin, soft tissue and the nares of human hosts,
282     while also being ubiquitous in a range of animals. However, the overall population density and the
283     proportion of human or animal hosts colonized by *S. epidermidis* largely exceed that of *S. aureus*,
284     so that *S. epidermidis*, but not *S. aureus*, can be considered of a physiological commensal, part of
285     the normal microbiota. The human *S. aureus* population is characterized by several genetically
286     distinct clonal complexes, each sharing a single ancestral genotype. Such a population can arise
287     under the neutral mutation/drift driven evolutionary trajectory combined with a high rate of
288     localized transmission. In this scenario clonal complexes appear and proliferate for a time, to be
289     replaced by others arising through genetic drift at the operational timescale of decades or longer.
290     This has been previously described as an 'epidemic clonal' structure(Smith et al., 2000).

291     We may consider that *E. faecalis* and *S. epidermidis*, members of the normal microbiota, have an
292     "endemic polyclonal structure", where endemicity is assured by a highly frequent inter-host
293     migration (both vertical and horizontal), resulting in a minimal adaptive stress in colonization of
294     most hosts. On the contrary, *E. faecium* and *S. aureus* are less-adapted organisms to the generality
295     of potential hosts, thus requiring local adaptation, and migration being dependent of this local
296     success, an "epidemic clonal structure". Obviously, in hospitals due to the homogenization of
297     colonizable hosts (age, antibiotic exposure), and facilitation of host-to-host migration (hospital
298     cross-colonization, microepidemics) *E. faecium* and *S. aureus* might appear as "locally endemic",
299     and therefore are expected to locally evolve towards a more complex population structure.

300     Both the commensal and hospital subpopulations of *S. epidermidis* display a pattern of genetic
301     relatedness typical of a population where recombination is the dominant force generating population
302     structure. An exception to this can be seen in the higher fraction of maximally distinct commensal

303  genotypes, which could plausibly arise when novel strains infrequently migrate to the human
304  commensal population from several non-overlapping zoonotic sources(Meric et al., 2015).
305  However, our model was not able to accurately predict the persistence of the clonal complex
306  structure observed for *S. aureus*, which may be reflecting a deviance from neutrality.

307  The complexities of within- and between-subpopulation strain dependence, and the extent of
308  localized transmission and migration across ecological patch boundaries makes formal statistical
309  inference about microepidemics and migration rates difficult. A particular challenge is that, when a
310  population evolves within a drift dominated model, it is unlikely that reliable estimates of the
311  parameters driving the population dynamics can be obtained, since observed outcomes of the
312  population structure vary substantially. Similarly, as the consequences of migration events are
313  dependent on other stochastically varying factors across the subpopulations, high migration rates
314  may lead to a pattern of relatedness indistinguishable from those generated by low rates. It is
315  possible that these issues could be resolved using coalescent-based models developed mainly for
316  eukaryotic populations(Beerli and Felsenstein, 1999; Beerli and Felsenstein, 2001; Choi and Hey,
317  2011; Hey and Machado, 2003; Hey and Nielsen, 2004). However, robust generalization of such
318  models is challenging due to the specific features of bacterial metapopulations which, in general,
319  evolve by a complex combination of the stochastic forces of mutation, recombination, clonal
320  expansion and host switches. Another obstacle for using coalescent-based methods is the large
321  number of hosts that need to be explicitly considered in studies on large-scale bacterial pathogen
322  populations.

323  It is evident that a limited number of neutrally evolving core genes, such as those typically used in
324  the MLST typing schemes, limits the scope of models that can be fitted to genetic surveillance data.
325  However, our results imply that some evolutionary scenarios would remain unidentifiable even if
326  housekeeping loci were considered at the whole-genome scale, in particular if the data are mainly
327  cross-sectional even if densely covering the host population. Hence, one of our main conclusions is
328  that the optimal data for studying dynamics in this fashion are densely sampled longitudinal
329  surveillance data covering evolutionary events at whole-genome level(Croucher et al., 2013). This
330  highlights the importance of easy access online repositories of genomic variation as an extension of
331  the currently existing MLST databases and that sample metadata should be an equally important
332  focus of the data sharing principles. Using such a strategy in the near future may enable important
333  model-based predictions about the dynamics of existing and emerging pathogens that pose a
334  considerable global challenge for human and animal health.

335

338  Author contributions

339  J.C., E.N., M.G. developed and implemented the model, P.M. and M.S. provided additional
340  expertise for the model development and analyses, J.C, G.M., S.K.S, T.C., F.B., W.V.S., R.W.,

341   E.F., W.P.H. provided data, biological expertise and interpretation, J.C., E.N., E.F. and W.P.H.
342   wrote the manuscript. All authors approved the final manuscript.

343   References

344

345   Bamshad, M. J., Mummidi, S., Gonzalez, E., Ahuja, S. S., Dunn, D. M., Watkins, W. S., Wooding, S., Stone, A.
346          C., Jorde, L. B., Weiss, R. B., Ahuja, S. K., 2002. A strong signature of balancing selection in the 5' cis-
347          regulatory region of CCR5. Proc Natl Acad Sci U S A 99, 10539-44, doi:10.1073/pnas.162046399.
348   Beerli, P., Felsenstein, J., 1999. Maximum-likelihood estimation of migration rates and effective population
349          numbers in two populations using a coalescent approach. Genetics 152, 763-73.
350   Beerli, P., Felsenstein, J., 2001. Maximum likelihood estimation of a migration matrix and effective
351          population sizes in n subpopulations by using a coalescent approach. Proc Natl Acad Sci U S A 98,
352          4563-8, doi:10.1073/pnas.081068098.
353   Choi, S. C., Hey, J., 2011. Joint inference of population assignment and demographic history. Genetics 189,
354          561-77, doi:10.1534/genetics.111.129205.
355   Croucher, N. J., Finkelstein, J. A., Pelton, S. I., Mitchell, P. K., Lee, G. M., Parkhill, J., Bentley, S. D., Hanage,
356          W. P., Lipsitch, M., 2013. Population genomics of post-vaccine changes in pneumococcal
357          epidemiology. Nature Genetics 45, 656-+, doi:Doi 10.1038/Ng.2625.
358   de Been, M., van Schaik, W., Cheng, L., Corander, J., Willems, R. J., 2013. Recent recombination events in
359          the core genome are associated with adaptive evolution in Enterococcus faecium. Genome Biol
360          Evol, doi:10.1093/gbe/evt111.
361   Everitt, R. G., Didelot, X., Batty, E. M., Miller, R. R., Knox, K., Young, B. C., Bowden, R., Auton, A., Votintseva,
362          A., Larner-Svensson, H., Charlesworth, J., Golubchik, T., Ip, C. L., Godwin, H., Fung, R., Peto, T. E.,
363          Walker, A. S., Crook, D. W., Wilson, D. J., 2014. Mobile elements drive recombination hotspots in
364          the core genome of Staphylococcus aureus. Nat Commun 5, 3956, doi:10.1038/ncomms4956.
365   Ewens, W. J., 2004. Mathematical population genetics. Springer, New York.
366   Faith, J. J., Colombel, J. F., Gordon, J. I., 2015. Identifying strains that contribute to complex diseases
367          through the study of microbial inheritance. Proceedings of the National Academy of Sciences of the
368          United States of America 112, 633-640, doi:DOI 10.1073/pnas.1418781112.
369   Feil, E. J., Li, B. C., Aanensen, D. M., Hanage, W. P., Spratt, B. G., 2004. eBURST: inferring patterns of
370          evolutionary descent among clusters of related bacterial genotypes from multilocus sequence
371          typing data. J Bacteriol 186, 1518-30.
372   Francisco, A. P., Bugalho, M., Ramirez, M., Carrico, J. A., 2009. Global optimal eBURST analysis of multilocus
373          typing data using a graphic matroid approach. BMC Bioinformatics 10, 152, doi:10.1186/1471-
374          2105-10-152.
375   Fraser, C., Hanage, W. P., Spratt, B. G., 2005. Neutral microepidemic evolution of bacterial pathogens. Proc
376          Natl Acad Sci U S A 102, 1968-73, doi:10.1073/pnas.0406993102.
377   Fraser, C., Hanage, W. P., Spratt, B. G., 2007. Recombination and the nature of bacterial speciation. Science
378          315, 476-80, doi:10.1126/science.1127573.
379   Fraser, C., Alm, E. J., Polz, M. F., Spratt, B. G., Hanage, W. P., 2009. The bacterial species challenge: making
380          sense of genetic and ecological diversity. Science 323, 741-6, doi:10.1126/science.1159388.
381   Goh, K. I., Oh, E., Jeong, H., Kahng, B., Kim, D., 2002. Classification of scale-free networks. Proc Natl Acad Sci
382          U S A 99, 12583-8, doi:10.1073/pnas.202301299.
383   Hanage, W. P., Fraser, C., Spratt, B. G., 2006. The impact of homologous recombination on the generation
384          of diversity in bacteria. J Theor Biol 239, 210-9, doi:10.1016/j.jtbi.2005.08.035.
385   Harpending, H. C., 1994. Signature of ancient population growth in a low-resolution mitochondrial DNA
386          mismatch distribution. Hum Biol 66, 591-600.
387   Hartl, D. L., Clark, A. G., 2007. Principles of population genetics. Sinauer Associates, Sunderland, Mass.
388   Hey, J., Machado, C. A., 2003. The study of structured populations--new hope for a difficult and divided
389          science. Nat Rev Genet 4, 535-43, doi:10.1038/nrg1112.

390  Hey, J., Nielsen, R., 2004. Multilocus methods for estimating population sizes, migration rates and
391       divergence time, with applications to the divergence of Drosophila pseudoobscura and D.
392       persimilis. Genetics 167, 747-60, doi:10.1534/genetics.103.024182.
393  Hudson, R. R., 1987. Estimating the recombination parameter of a finite population model without
394       selection. Genet Res 50, 245-50.
395  Lebreton, F., van Schaik, W., McGuire, A. M., Godfrey, P., Griggs, A., Mazumdar, V., Corander, J., Cheng, L.,
396       Saif, S., Young, S., Zeng, Q. D., Wortman, J., Birren, B., Willems, R. J. L., Earl, A. M., Gilmore, M. S.,
397       2013. Emergence of Epidemic Multidrug-Resistant Enterococcus faecium from Animal and
398       Commensal Strains. MBio 4, doi:ARTN e00534-13DOI 10.1128/mBio.00534-13.

399  Meric, G., Miragaia, M., de Been, M., Yahara, K., Pascoe, B., Mageiros, L., Mikhail, J., Harris, L. G., Wilkinson,
400       T. S., Rolo, J., Lamble, S., Bray, J. E., Jolley, K. A., Hanage, W. P., Bowden, R., Maiden, M. C., Mack,
401       D., de Lencastre, H., Feil, E. J., Corander, J., Sheppard, S. K., 2015. Ecological Overlap and Horizontal
402       Gene Transfer in Staphylococcus aureus and Staphylococcus epidermidis. Genome Biol Evol 7,
403       1313-28, doi:10.1093/gbe/evv066.
404  Mousset, S., Derome, N., Veuille, M., 2004. A test of neutrality and constant population size based on the
405       mismatch distribution. Mol Biol Evol 21, 724-31, doi:10.1093/molbev/msh066.
406  Plucinski, M. M., Starfield, R., Almeida, R. P., 2011. Inferring social network structure from bacterial
407       sequence data. PLoS One 6, e22685, doi:10.1371/journal.pone.0022685.
408  Rogers, A. R., Harpending, H., 1992. Population growth makes waves in the distribution of pairwise genetic
409       differences. Mol Biol Evol 9, 552-69.
410  Smith, J. M., Feil, E. J., Smith, N. H., 2000. Population structure and evolutionary dynamics of pathogenic
411       bacteria. Bioessays 22, 1115-22, doi:10.1002/1521-1878(200012)22:12<1115::AID-BIES9>3.0.CO;2-
412       R.
413  Tedim, A. P., Ruiz-Garbajosa, P., Corander, J., Rodriguez, C. M., Canton, R., Willems, R. J., Baquero, F.,
414       Coque, T. M., 2015. Population biology of intestinal enterococcus isolates from hospitalized and
415       nonhospitalized individuals in different age groups. Appl Environ Microbiol 81, 1820-31,
416       doi:10.1128/AEM.03661-14.
417  Turner, K. M., Hanage, W. P., Fraser, C., Connor, T. R., Spratt, B. G., 2007. Assessing the reliability of eBURST
418       using simulated populations with known ancestry. BMC Microbiol 7, 30, doi:10.1186/1471-2180-7-
419       30.
420  Vos, M., Didelot, X., 2009. A comparison of homologous recombination rates in bacteria and archaea. ISME
421       J 3, 199-208, doi:10.1038/ismej.2008.93.
422  Willems, R. J., Top, J., van Schaik, W., Leavis, H., Bonten, M., Siren, J., Hanage, W. P., Corander, J., 2012.
423       Restricted gene flow among hospital subpopulations of Enterococcus faecium. MBio 3, e00151-12,
424       doi:10.1128/mBio.00151-12.
425  Wood, S. N., 2010. Statistical inference for noisy nonlinear ecological dynamic systems. Nature
426       466(7310):1102–1104.

427  Figure legends

428  Fig. 1. Allelic mismatch distributions for combinations of mutation and recombination rates in a population
429  with $N = 3000$. Bold line in green shows the mean mismatch probability over 20000 generations, sampled at
430  intervals of 100 generations. The green shaded area shows the 95% confidence interval and the colored lines
431  are examples of mismatch distributions at random time points. Vertical axis in each panel shows the
432  probability mass associated with the points of the curves across the values on the horizontal axis.
433  Distributions are shown as continuous curves for visual clarity.
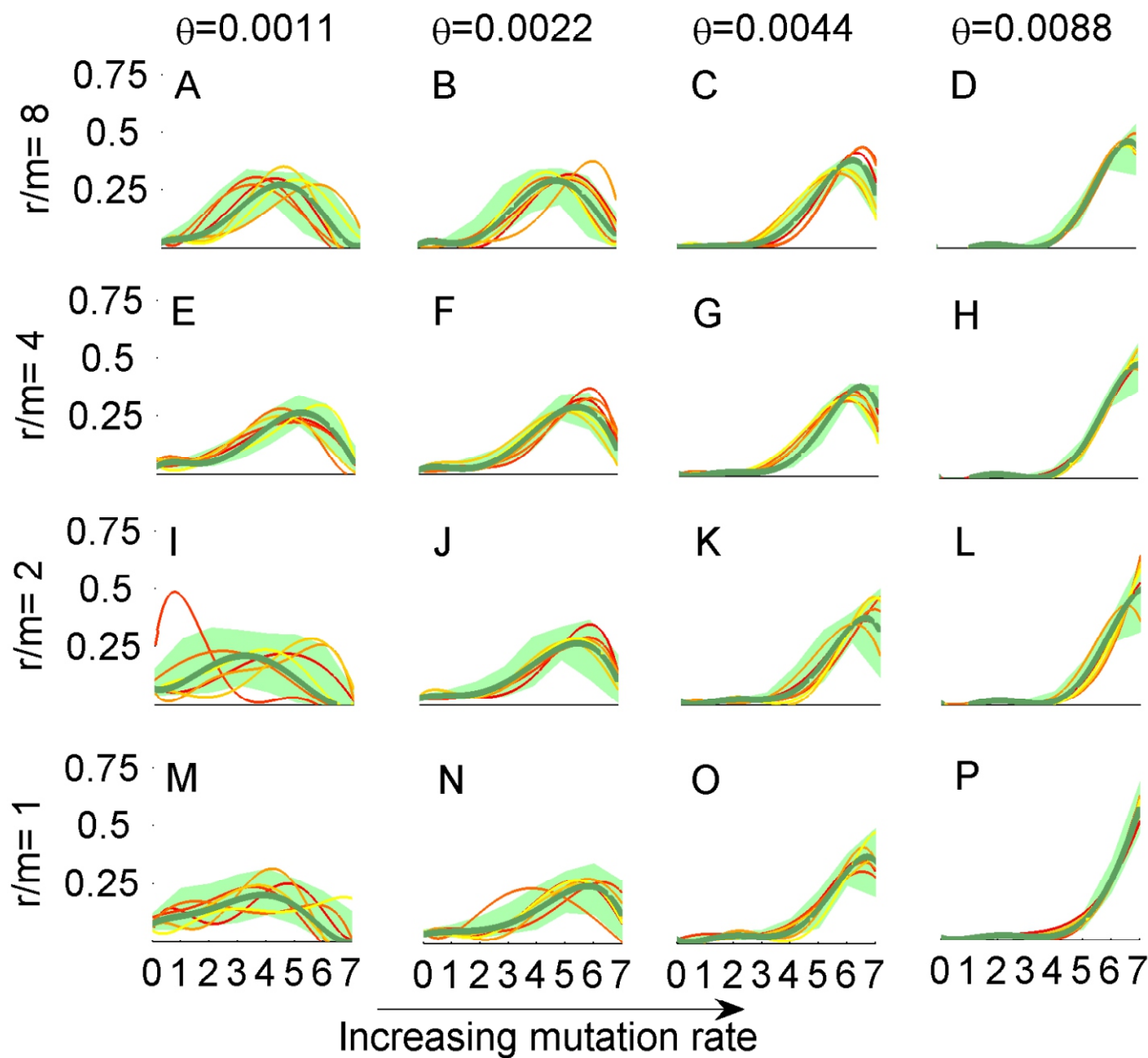
434  Fig. 2. eBURST networks and mismatch distributions for a population without (grey) and with (yellow)
435  microepidemics, where $\omega = 27$, $\gamma = 16$. The 95% confidence intervals are shown by shaded areas and are
436  defined as in Fig. 1. The mutation and recombination parameters used are: 0.0011, 1 (A), 0.0088, 1 (B),
437  0.0011, 8 (C), 0.0088, 8 (D).

438    Fig. 3. Schematic illustration of the combined effect of microepidemics and migration. The population on the
439    left is unstratified, in which case increasing rate ($\omega$) and size ($\gamma$) of microepidemics lead to decreased genetic
440    variation. In a stratified population with two subpopulations ($P_1$, $P_2$) the effect of increasing microepidemics
441    ($\omega_1$, $\gamma_1$) on genetic diversity in subpopulation $P_1$ depends both on the microepidemics in subpopulation $P_2$
442    ($\omega_2$, $\gamma_2$) and on the migration rate ($m_{21}$). The case with $m_{21} = 0$ leads to identical decrease of genetic variation
443    as in an unstratified population. The notation "<<" is used to indicate that the parameters on the left side of
444    the double inequality are much smaller than those on the right side.

445    Fig. 4. Mismatch distributions of commensal and hospital subpopulations of four common nosocomial
446    bacterial pathogens. The right-most column shows the between-subpopulation mismatch distributions.

447    Fig. 5. eBURST networks of the isolates used to calculate the mismatch distribution in Fig. 4; *E. faecalis*
448    (A), *E. faecium* (B), *S. aureus* (C), *S. epidermidis* (D).

Figure showing a 4×4 grid of plots. Columns labeled θ=0.0011, θ=0.0022, θ=0.0044, θ=0.0088. Rows labeled r/m=8, r/m=4, r/m=2, r/m=1 (Increasing recombination rate). X-axis: Increasing mutation rate (0 1 2 3 4 5 6 7). Panels labeled A–P. Y-axis values: 0.25, 0.5, 0.75.
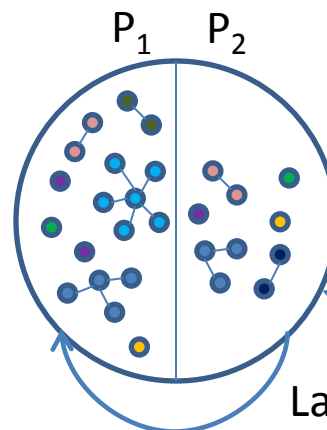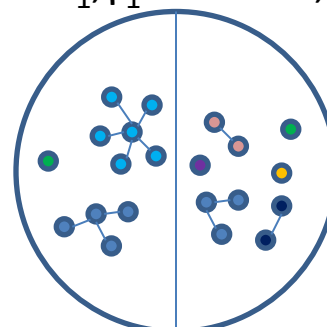
High level of genetic diversity

P$_1$  P$_2$

$\omega_1, \gamma_1$ increase
$\omega_2, \gamma_2 \ll \omega_1, \gamma_1$

P$_1$  P$_2$

Large m$_{21}$

Baseline level of genetic diversity

$\omega_1, \gamma_1$ increase, m$_{21}$ = 0
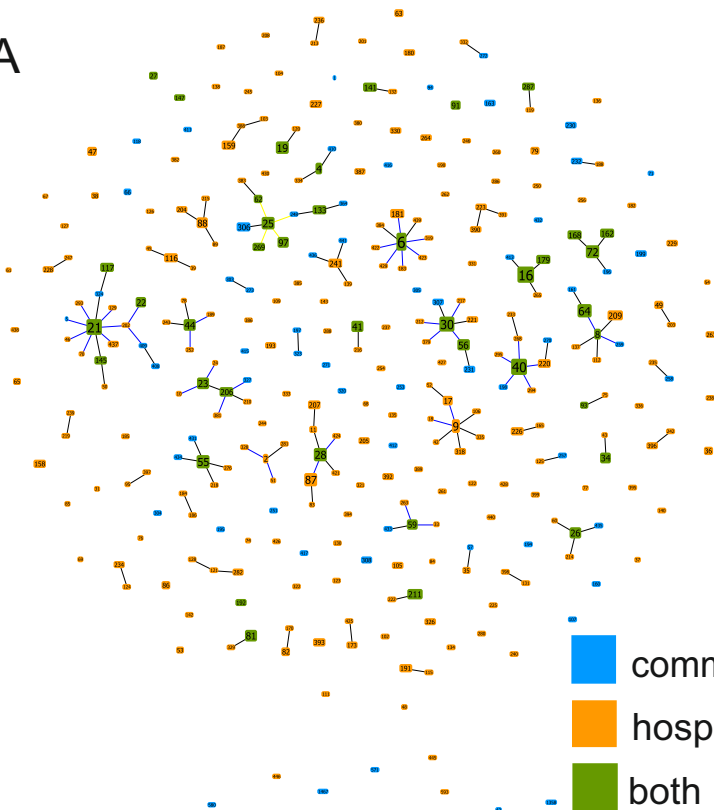
$\omega, \gamma$ increase

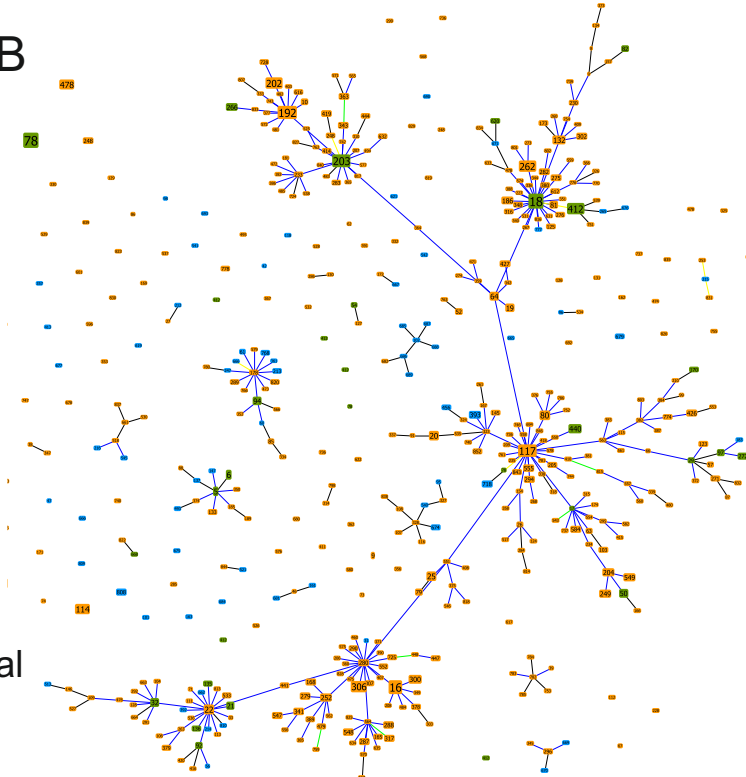Unstratified population

Stratified population

$\omega_1, \gamma_1$ increase
$\omega_2, \gamma_2 \gg \omega_1, \gamma_1$

P$_1$  P$_2$

Low level of genetic diversity

Large m$_{21}$

A

B

C

D

commensal

hospital

both