



Minerva Access is the Institutional Repository of The University of Melbourne

Author/s:

Wan, Y;Wick, RR;Zobel, J;Ingle, DJ;Inouye, M;Holt, KE

Title:

GeneMates: an R package for detecting horizontal gene co-transfer between bacteria using gene-gene associations controlled for population structure

Date:

2020-09-24

Citation:

Wan, Y., Wick, R. R., Zobel, J., Ingle, D. J., Inouye, M. & Holt, K. E. (2020). GeneMates: an R package for detecting horizontal gene co-transfer between bacteria using gene-gene associations controlled for population structure. *BMC Genomics*, 21 (1), pp.658-. <https://doi.org/10.1186/s12864-020-07019-6>.

Persistent Link:

<https://hdl.handle.net/11343/251620>

License:

[CC BY](#)

SOFTWARE

Open Access



GeneMates: an R package for detecting horizontal gene co-transfer between bacteria using gene-gene associations controlled for population structure

Yu Wan^{1*} , Ryan R. Wick^{1,2}, Justin Zobel³, Danielle J. Ingle^{4,5}, Michael Inouye^{6,7} and Kathryn E. Holt^{2,8}

Abstract

Background: Horizontal gene transfer contributes to bacterial evolution through mobilising genes across various taxonomical boundaries. It is frequently mediated by mobile genetic elements (MGEs), which may capture, maintain, and rearrange mobile genes and co-mobilise them between bacteria, causing horizontal gene co-transfer (HGcoT). This physical linkage between mobile genes poses a great threat to public health as it facilitates dissemination and co-selection of clinically important genes amongst bacteria. Although rapid accumulation of bacterial whole-genome sequencing data since the 2000s enables study of HGcoT at the population level, results based on genetic co-occurrence counts and simple association tests are usually confounded by bacterial population structure when sampled bacteria belong to the same species, leading to spurious conclusions.

Results: We have developed a network approach to explore WGS data for evidence of intraspecies HGcoT and have implemented it in R package GeneMates (github.com/wanyuac/GeneMates). The package takes as input an allelic presence-absence matrix of interested genes and a matrix of core-genome single-nucleotide polymorphisms, performs association tests with linear mixed models controlled for population structure, produces a network of significantly associated alleles, and identifies clusters within the network as plausible co-transferred alleles. GeneMates users may choose to score consistency of allelic physical distances measured in genome assemblies using a novel approach we have developed and overlay scores to the network for further evidence of HGcoT. Validation studies of GeneMates on known acquired antimicrobial resistance genes in *Escherichia coli* and *Salmonella Typhimurium* show advantages of our network approach over simple association analysis: (1) distinguishing between allelic co-occurrence driven by HGcoT and that driven by clonal reproduction, (2) evaluating effects of population structure on allelic co-occurrence, and (3) direct links between allele clusters in the network and MGEs when physical distances are incorporated.

Conclusion: GeneMates offers an effective approach to detection of intraspecies HGcoT using WGS data.

Keywords: Horizontal gene transfer, Acquired genes, Mobile genetic elements, Physical linkage, Population structure, Association analysis, Linear mixed models, Principal components, Network approach, R package

*Correspondence: wanyuac@126.com

¹Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, University of Melbourne, Parkville, Victoria, 3010 Australia

Full list of author information is available at the end of the article



© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Background

Horizontal gene transfer (HGT) or lateral gene transfer accelerates bacterial genome innovation and evolution [1]. It contributes to gene flows across taxonomic boundaries and thus variation in accessory gene content both within and between species [2, 3]. Mobile genetic elements (MGEs), such as plasmids, bacteriophages, and transposons, are common vectors of acquired genes in HGT [4]. Particularly, when acquired genes are physically linked (namely, co-localised in an MGE or otherwise physically close in DNA molecules), they can be horizontally co-transferred between bacteria, causing positive gene-gene associations known as genetic linkage [5, 6].

The rapid accumulation of bacterial whole-genome sequencing (WGS) data in the most recent two decades [7] enables us to study horizontal gene co-transfer (HGcoT) at the population level. Since bacteria reproduce asexually and HGT can occur across different levels of taxonomic boundaries [2], gene-gene associations that cannot be completely explained by bacterial population structure (which determines the distribution of co-inherited genes) suggests HGcoT [8]. Consequently, it is usually trivial to identify candidates of interspecies HGcoT using simple association tests (such as chi-squared tests and simple logistic regression), whereas for detecting intraspecies HGcoT, we must overcome two related challenges arising from the presence of population structure within a species: (1) how to control for population structure in association tests; and (2) how to accurately estimate or represent the population structure to be controlled for.

Univariate linear mixed models (LMMs), which have been widely used in human genome-wide association studies (GWAS) [9] and recently applied to bacterial GWAS [10], provide a solution to address both challenges. Each model explains a response variable using a fixed effect of an independent variable and mixed random effects of population structure and environmental factors. For each LMM, the population structure is represented by a relatedness matrix, whose principal components (PCs) can be used for an orthonormal transformation of genetic variation underlying the population structure [11]. McVean demonstrates that not only do these PCs simplify computations, but also correlate with bacterial genealogies [12]. Compared to phylogeny-based association tests, such as phylogenetically independent contrasts and phylogenetic generalised least squares [13], LMMs do not rely on specific phylogenetic models nor assume that mutation rather than recombination dominates genetic variation.

Here, we introduce R package GeneMates, which implements a novel network approach to the identification of HGcoT within a bacterial species. This approach takes as input specific information extracted from bacterial WGS data and produces a network showing allele-level evidence

of HGcoT. We validated GeneMates using published WGS data from two bacterial species, *Escherichia coli* and *Salmonella enterica*, with which we identify horizontally co-transferred antimicrobial resistance (AMR) genes. We also provide helper scripts to assist users in preparing necessary input files from standard formats. Since GeneMates is theoretically applicable to any kind of acquired genes in bacteria, it has the potential to also be used for investigating structures and dissemination of other mobile gene clusters of interest, such as physically linked virulence genes.

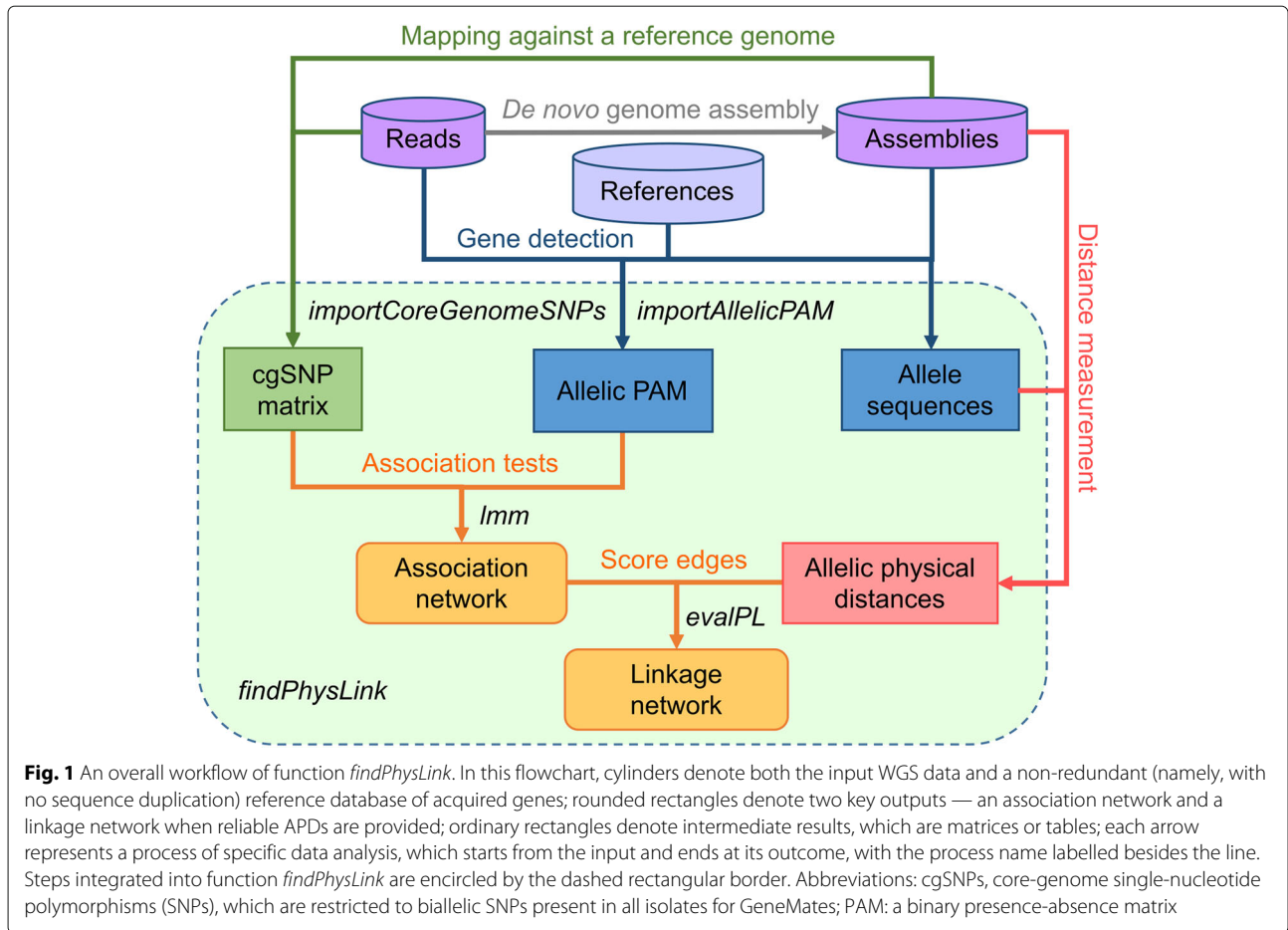
Implementation

GeneMates consists of R functions performing network construction, topological analysis, and data visualisation. It works at the allele level to track recent HGcoT. Particularly, we assume that bacterial isolates under study will typically be collected in a period that is too short to accumulate any mutation in recently acquired genes of interest, such as AMR genes. In addition, for convenience we assume that every one of these acquired genes has at most one allele in each genome because it is not possible to reliably resolve co-occurring alleles of the same gene using short reads, which by far account for the majority of WGS data.

We developed a network approach to integrate and visualise evidence of physical linkage between alleles of acquired genes in bacteria. In our network, nodes represent alleles and weighted directed edges reflect the strength of evidence. GeneMates produces such a network in each run. Let vector e denote an edge, then it represents a linear model $Y \sim X$, where scalar variables Y and X denote presence-absence of alleles Y and X in an isolate, respectively. This model explains the distribution of allele Y (response allele) using that of allele X (explanatory allele) and covariates. For GeneMates, in particular, the covariates are isolate projections from an orthonormal transformation of the population structure based on a core-genome relatedness matrix. Correspondingly, the edge is directed from node X to node Y . Users may filter edges of a resulting network based on an association score $s_a(e)$ and a distance score $s_d(e)$ in order to identify edges showing strong evidence of physical linkage, and carry out topological analysis on the filtered network afterwards. In following subsections, we describe key elements of our approach. Additional details of implementation are provided in Section 3 of Additional file 1.

Network construction

Figure 1 illustrates our work flow for network construction — the core functionality of GeneMates. In order to integrate functions implementing this work flow, we created a wrapper function *findPhysLink* (find physical



linkage), which takes as input a binary allelic presence-absence matrix (PAM) of target genes across genomes, a matrix of biallelic core-genome single-nucleotide polymorphisms (cgSNPs), and optionally, a table of allelic physical distances (APDs) for target genes. These inputs can be extracted from read alignments and genome assemblies using our helper scripts (released with GeneMates). Function *findPhysLink* determines nodes and edges of a resulting network and produces tables that can be exported to Cytoscape [14] as node and edge attributes for network visualisation.

Node generation

Assuming m alleles of target genes are detected in n bacterial isolates, GeneMates function *importAllelicPAM* imports an $n \times m$ binary PAM $A = (a_{ij})$, where entry $a_{ij} = 1$ if the j -th allele is found in the i -th isolate, and $a_{ij} = 0$ otherwise. This function offers two optional filters to discard alleles of insufficient frequencies and/or co-occurrence frequencies. In order to reduce the number of tests for allele-allele associations, we followed the implementation of R package BugWAS [10] and coded *importAllelicPAM* to de-duplicate each group of identical columns of PAM into a binary vector called an allelic

distribution pattern (Section 3.1.2 of Additional file 1). The function uses column means to apply a column-wise zero-centring to the pattern matrix, which is a common technique used for simplifying matrix algebra without affecting the distribution of data points [15].

Edge weights

GeneMates evaluates two kinds of evidence for inference of horizontal co-transfer of alleles X and Y . The first evidence is a significant positive fixed effect of X on presence-absence of Y when controlling for bacterial population structure. Testing for this effect is an analogue of GWAS, which test for genotype-phenotype associations. Specifically, for edge e in an association network, GeneMates function *Imm* estimates parameters of a univariate LMM using a residual maximum likelihood (REML) approach to test for the fixed effect of X (Section 3.1 in Additional file 1), and another function *evalPL* transforms the estimated effect size $\hat{\beta}$ and its Bonferroni-corrected p -value into an association score $s_a(e)$ with possible values 1 (significant positive association), -1 (significant negative association), and 0 (insignificant association). See Figure s1 and Section 3.3.1 in Additional file 1 for details.

The second evidence comes from consistent physical distances between alleles X and Y (namely, APDs) in different bacterial genomes as structural variation is likely to only occur at a limited level within a mobile gene cluster in a short period. For instance, the same AMR gene cluster *sul2-strA-strB* has been circulating amongst Gram-negative bacteria for decades due to its association with plasmids and transposons [16]. Since APDs are measured in genome assemblies, whose completeness determines the amount of measurable APDs when X and Y are co-localised in the same genomic region, for edge e , we also consider its distance measurability $m_{in}(e)$ — the percentage of genomes in which the APDs between X and Y are actually measurable. This measurability value is calculated by GeneMates function *summariseDist*, which also evaluates the consistency of APDs included for the distance assessment and assigns a consistency score $c(e)$ with values -1 (evidence against physical linkage), 0 (insufficient evidence) and 1 (evidence supporting physical linkage). Notably, this function estimates the probability of distance identity-by-descent (IBD) and compares it to a user predefined threshold (default: 0.9) for the assignment of each consistency score. (Section 3.3.2 and Figure s2 in Additional file 1) A summary distance score is thereby defined for edge e as the consistency score weighted by measurability: $s_d(e) = m_{in}(e)c(e)$. Finally, the association and distance scores are summed to get a linkage score $s(e)$, reflecting the evidence of physical linkage, where $-2 \leq s(e) \leq 2$ because $0 \leq m_{in}(e) \leq 1$, and $s(e)$ has five levels (Table s1). Particularly, we define a linkage network as an association network in which weights of each edge consist of a fixed-effect size and a linkage score (Section 3.3 in Additional file 1).

Network visualisation and topological analysis

GeneMates comprises several functions used for displaying resulting networks and exploring network topology for evidence of HGcoT under given conditions. For instance, function *mkNetwork* (make network) extracts user-specified node attributes (such as the frequency and associated AMR phenotype of each allele) and edge attributes (such as the estimated fixed effect size $\hat{\beta}$ and linkage score s of each edge) from result tables of *findPhysLink* and prints them to Cytoscape-compatible text files for network visualisation. Function *extractSubgraphs* follows a predefined node list to pull out subnetworks from a parental network produced by *findPhysLink*. Particularly, the subnetworks can be maximal cliques identified using function *max_cliques* of R package *igraph* [17]. In addition, function *countNeighbours* lists the number of neighbours per node in a network and function *getClusterMemberCooccurrence* identifies isolates in which member alleles of a selected subnetwork are co-occurring.

Results

We assessed the performance of GeneMates and validated our methodology using published WGS data sets of two bacterial pathogens of great clinical concern: multidrug-resistant *E. coli* and *S. enterica* serovar Typhimurium. Genomes in these well studied data sets have distinct population structures, harbour diverse AMR genes and MGEs, and show known gene-gene associations that we expected GeneMates to identify. See Section 4 of Additional file 1 for details of materials and methods.

Characteristics of example data sets

Both the *E. coli* and *Salmonella* data sets consisted of paired-end Illumina WGS reads, generated from 169 *E. coli* isolates collected during the Global Enteric Multicentre Study [18, 19] and 359 isolates of typical *S. Typhimurium* Definitive Type 104 [20], respectively. See Additional file 2 for detailed isolate information.

AMR gene content In genomes of the 169 *E. coli* isolates, we identified 178 alleles of 33 AMR genes conferring resistance to eight antimicrobial classes (Figure s3). The four known intrinsic AMR genes of *E. coli* (*ampH*, *ampC1*, *ampC2*, and *mrda*) displayed higher frequencies (> 87%) than the 29 acquired AMR genes (< 63%). Altogether, we detected 67 alleles of acquired AMR genes, including 45 alleles showing frequencies less than 3%. In genomes of the 359 *Salmonella* isolates, we identified 57 alleles of 24 AMR genes (conferring resistance to six antimicrobial classes), including a single allele of the known intrinsic AMR gene (*aac6-1aa*) of *S. enterica* and 56 alleles of 23 acquired AMR genes (Figure s5). Notably, five acquired AMR genes (*sul1*, *aadA*, *bla_{CARB}*, *tet(G)*, and *floR*) that are known to be frequently present in *Salmonella* genomic island 1 (SGI1) [21] were only detected in the dominant lineage of collected *Salmonella* genomes (Figure s6), whereas alleles of other acquired AMR genes were sparsely distributed across a core-genome phylogenetic tree of these *Salmonella* genomes (Figure s7). Further, we identified four and one clusters of identically distributed alleles in *E. coli* and *Salmonella* genomes, respectively (Figure s8).

Core-genome SNP sites We analysed chromosomal single-nucleotide polymorphisms (SNPs) of each species using the method described in Section 4.2 of Additional file 1. Particularly, we define cgSNP sites as SNP sites detected in all chromosomes and outside of repetitive or prophage regions. Altogether, numbers of cgSNP sites used for correcting for population structure (namely, biallelic cgSNP sites) of the 169 *E. coli* genomes and 359 *Salmonella* genomes were 209,097 and 2,316, respectively. The percentage of total genetic variation captured by PCs

of each cgSNP matrix is illustrated in Figure s9a. Furthermore, the REML estimate of parameter λ (which reflects the effect of population structure on the distribution of the response allele) in the null LMM for presence-absence of each allele of AMR genes perfectly predicted whether ≥ 5 PCs had significant effects (Bonferroni-corrected p -values ≤ 0.05) on the distribution of this allele (Figure s9b).

Effects of controlling for population structure

In order to evaluate the effect of controlling for population structure using LMMs when measuring associations between alleles of acquired AMR genes, we compared unadjusted p -values of fixed effects estimated using the LMMs to those estimated using simple penalised logistic models (PLMs) [22] (Fig. 2). Particularly, we considered a fixed effect estimated with either kind of models significant if its Bonferroni-corrected p -value was ≤ 0.05 . Figure 3 illustrates a directed comparative network for

detected alleles of AMR genes in each example data set. In this network, each edge starts from an explanatory allele and terminates at a response allele, representing a significant fixed effect of the explanatory allele in an LMM or PLM or both. For the 67 alleles in *E. coli* genomes, 3,364 LMMs were applied to 60 allele distribution patterns, and 118 significant pairwise associations were detected. Simple PLMs for the same allele patterns of AMR genes in *E. coli* genomes identified 70 significant associations, 50 of which overlapped with those from LMMs. The resulting comparative network consisted of 45 nodes, 138 edges, and two connected graph components (Fig. 3, *E. coli*). Of note, one of these components was comprised of rare alleles and was only identified by LMMs (Figure s10). Regarding the 56 alleles in *Salmonella* genomes, 2,040 LMMs were applied to 48 allele distribution patterns and 112 significant pairwise associations were detected. Simple PLMs for the same allele patterns of AMR genes in *Salmonella* genomes

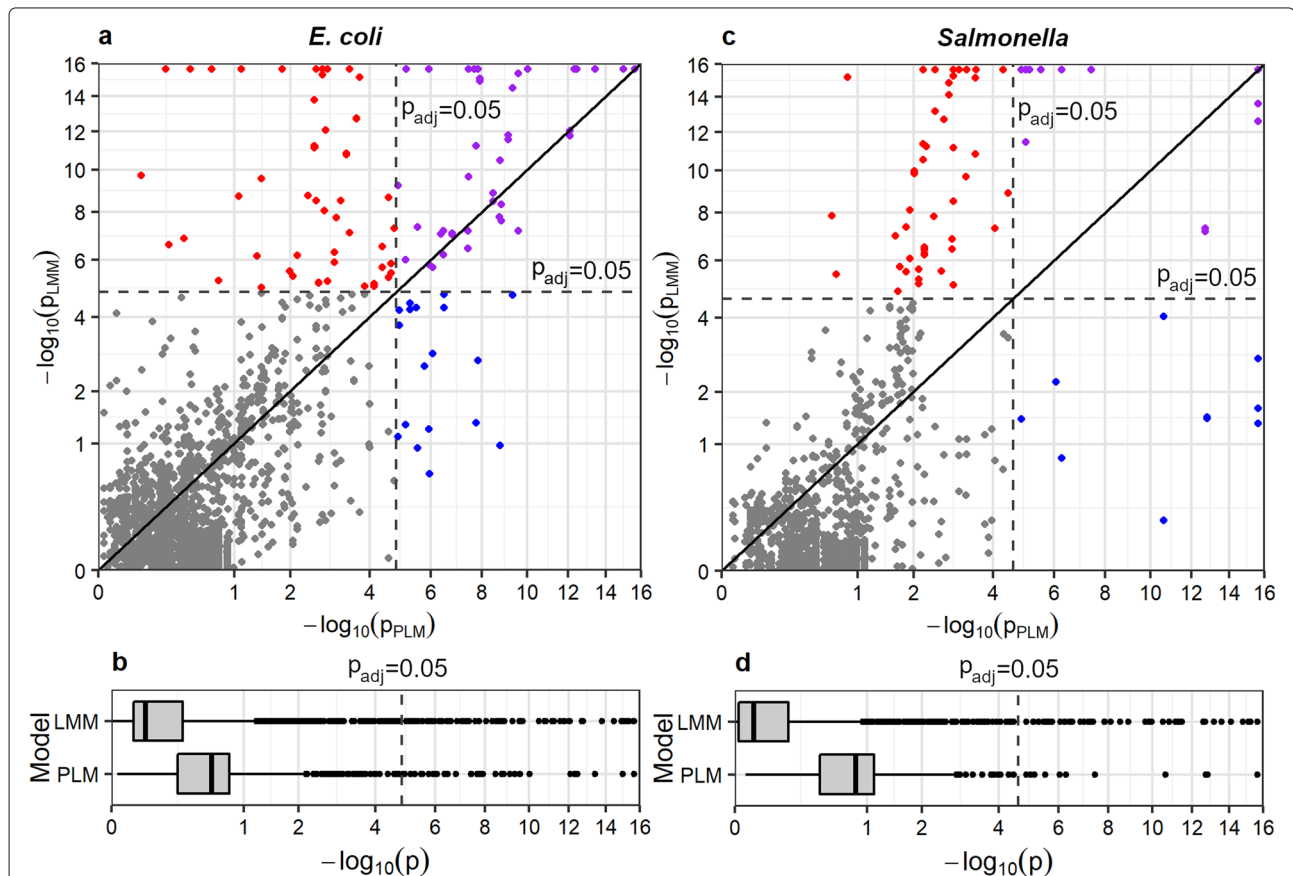


Fig. 2 A comparison of unadjusted p -values from PLMs and LMMs for the same pairs of allelic distribution patterns of AMR genes. For the *E. coli* data set (panels a and b) and *Salmonella* data set (panels c and d), respectively, a scatter plot and box plot of paired p -values are drawn on square-root transformed axes to compare the p -values. Any p -value less than 2.2×10^{-16} is rounded to 2.2×10^{-16} owing to the smallest precise floating number in our computer. In panels a and c, black diagonal lines indicate equality between p -values from these two kinds of linear models, and grey dashed lines indicate the p -value corresponding to the Bonferroni corrected p -value of 0.05, which is used in this study as a cut-off for significant associations. Associations that were only significant in PLMs, only significant in LMMs, significant in both PLMs and LMMs, and significant in neither kind of models, are represented by blue, red, purple, and grey circles, respectively

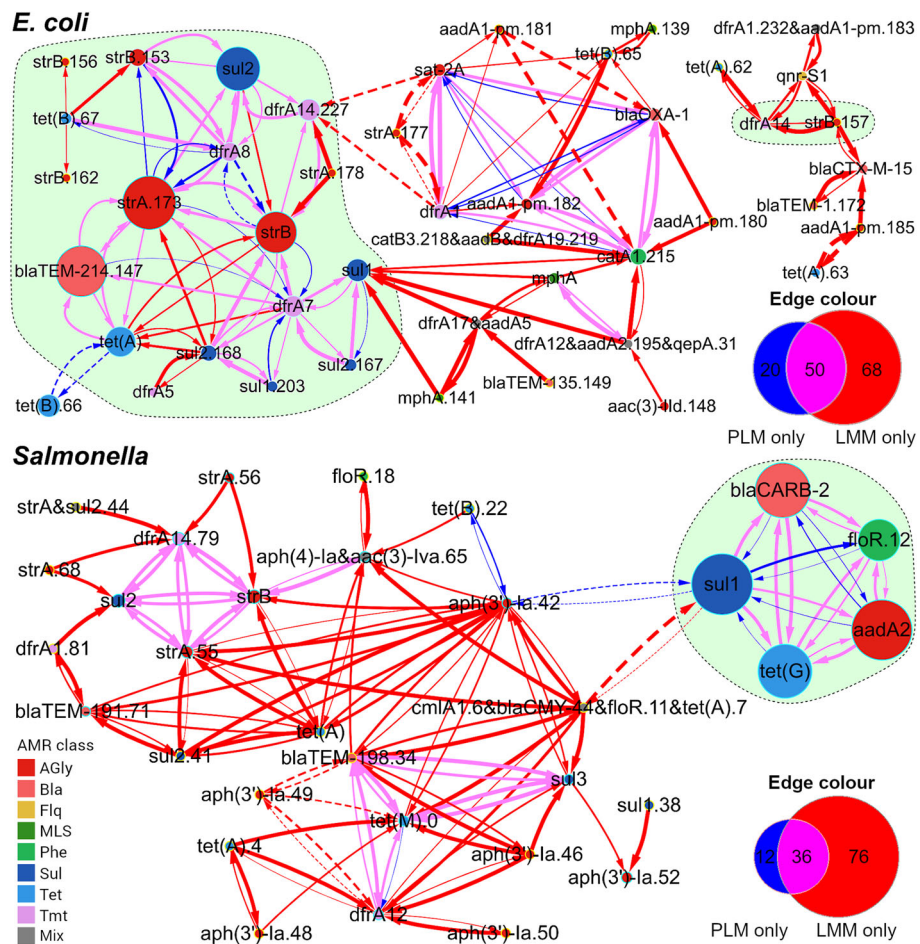


Fig. 3 Comparative networks for detected alleles of acquired AMR genes in 169 *E. coli* and 359 *Salmonella* genomes. Each node represents an allele or a cluster of identically distributed alleles, with a diameter proportional to the allele frequency and a fill colour indicating the AMR phenotype encoded. The yellow node border denotes an allele or allele cluster detected in one genome. The edge width is proportional to the strength of a significant association ($|\hat{\beta}|$) determined using an LMM. Edge colours indicate significant associations from LMMs or PLMs or both. Solid and dashed edges represent significant positive ($\hat{\beta} > 0$) and significant negative associations ($\hat{\beta} < 0$), respectively, identified using models indicated by edge colours. The Venn diagram of each panel counts edges under each colour. The shaded area (light green) encircles alleles of known co-transferred AMR genes. AMR classes defined by antimicrobials that bacteria were resistant to: AGly (aminoglycosides), Bla (beta-lactams), Flq (fluoroquinolones), MLS (macrolides, lincosamides, and streptogramins), Phe (phenicols), Sul (sulfonamides), Tet (tetracyclines), Tmt (trimethoprim), and Mix (multiple classes of antimicrobials)

identified 48 associations, 36 of which overlapped with those from LMMs. The resulting comparative network consisted of 32 nodes, 124 edges, and one connected graph component (Fig. 3, *Salmonella*). In addition, for both *E. coli* and *Salmonella* data sets, estimates of fixed effects in LMMs displayed a complete sign identity to those in PLMs given a maximum type-I error rate of 0.05.

For most allele pairs, correcting for population structure using LMMs yielded greater *p*-values than did PLMs: the majority of fixed effects (74% for *E. coli* and 85% for *Salmonella*) in LMMs became less significant than those in PLMs (Fig. 2b, d). By contrast, associations in some allele pairs became more significant (that

is, of smaller *p*-values) after adjusting for population structure, and in most of these pairs (93% and 99% for *E. coli* and *Salmonella*, respectively), LMMs identified a moderate to strong structural random effect (namely, parameter estimate $1 < \hat{\lambda}_0 \leq 10^5$ in an LMM under the null hypothesis of no fixed effect from an explanatory allele) underlying presence-absence status of one or both paired alleles across isolates. Further, using one-sided two-proportions Z-tests, we found that LMMs showing moderate to strong structural random effects were less likely to show increased *p*-values than were other LMMs when compared to PLMs, with contrasts of proportions 70% versus 90% (*p*-value ≈ 0) in LMMs for the *E. coli* data set and 85% versus

92% (p -value = 0.06) in LMMs for the *Salmonella* data set.

Effects of adding APDs to association networks

Every APD was measured by a shortest-path distance (SPD), defined as the smallest distance between two given loci in an assembly graph (constructed using Unicycler [23]). This approach allows us to recover query sequences that are split into adjacent contigs and to measure APDs between loci that are located in different contigs, potentially increasing distance measurability of the graph, defined as the percentage of measured APDs in all possible APDs of a complete circular genome. Nonetheless, since SPDs are affected by the topology of each assembly graph, which is usually tangled and partially resolved when only short reads are used for genome assembly, it is necessary to determine appropriate criteria for filtering out inaccurate SPDs. Therefore, we downloaded from GenBank [24] 10 complete reference genomes of multidrug-resistant *E. coli* and *S. Typhimurium*, respectively (Tables s2, s3 and Additional file 2), and compared SPDs measured in *de novo* assembly graphs (constructed from simulated Illumina reads) to true physical distances extracted from original circularised complete genomes (Section 4.6 of Additional file 1). Specifically, we used Bandage [25] to identify BLAST hits of random coding sequences (CDSs) in each assembly graph and to extract SPDs for each pair of hits. As expected, Bandage always recovered more query CDSs from the assembly graph than it did from contigs of the same bacterial genome (Tables s4 and s5).

Filters determined for removing inaccurate SPDs We considered an SPD accurate if its error fell within a given tolerance range (for instance, ± 2 kbp), and hence defined the accuracy rate as the percentage of accurate SPDs in all SPDs, either filtered or not. In practice, any parameter for BLAST hits and SPD measurement can be taken into account for excluding SPDs. In this study, we assessed the accuracy rate of SPDs measured within various maximum distances and node numbers in each assembly graph when confining BLAST hits to a minimum query coverage and nucleotide identity of 95%. Across *E. coli* and *S. Typhimurium* genomes, we constantly saw that the accuracy rate reached $\geq 90\%$ under an error tolerance ± 1 kbp when SPDs were measured within 250 kbp and no more than two nodes (Figures s11 and s12).

Moreover, since the accuracy rate of SPDs measured within contigs stayed above 90% when tolerating errors within ± 1 kbp (Figures s13 and s14), we implemented prioritisation of SPDs based on their sources (namely, contigs or assembly graphs) in order to exclude inaccurate SPDs from assembly graphs where repeats had not been resolved by the genome assembler. Specifically, when the

SPD between two CDSs is measurable in both a contig and an assembly graph of the same genome, this method overrides the graph-based SPD with its corresponding contig-based SPD, thereby taking advantage of both the high accuracy rate of contig-based SPDs and high measurability of graph-based SPDs. As shown in Tables s6 and s7, the prioritisation method led to an accuracy rate of 100% for the majority of SPDs measured between acquired AMR genes, which are often embedded within tangled sub-graphs owing to surrounding repeats.

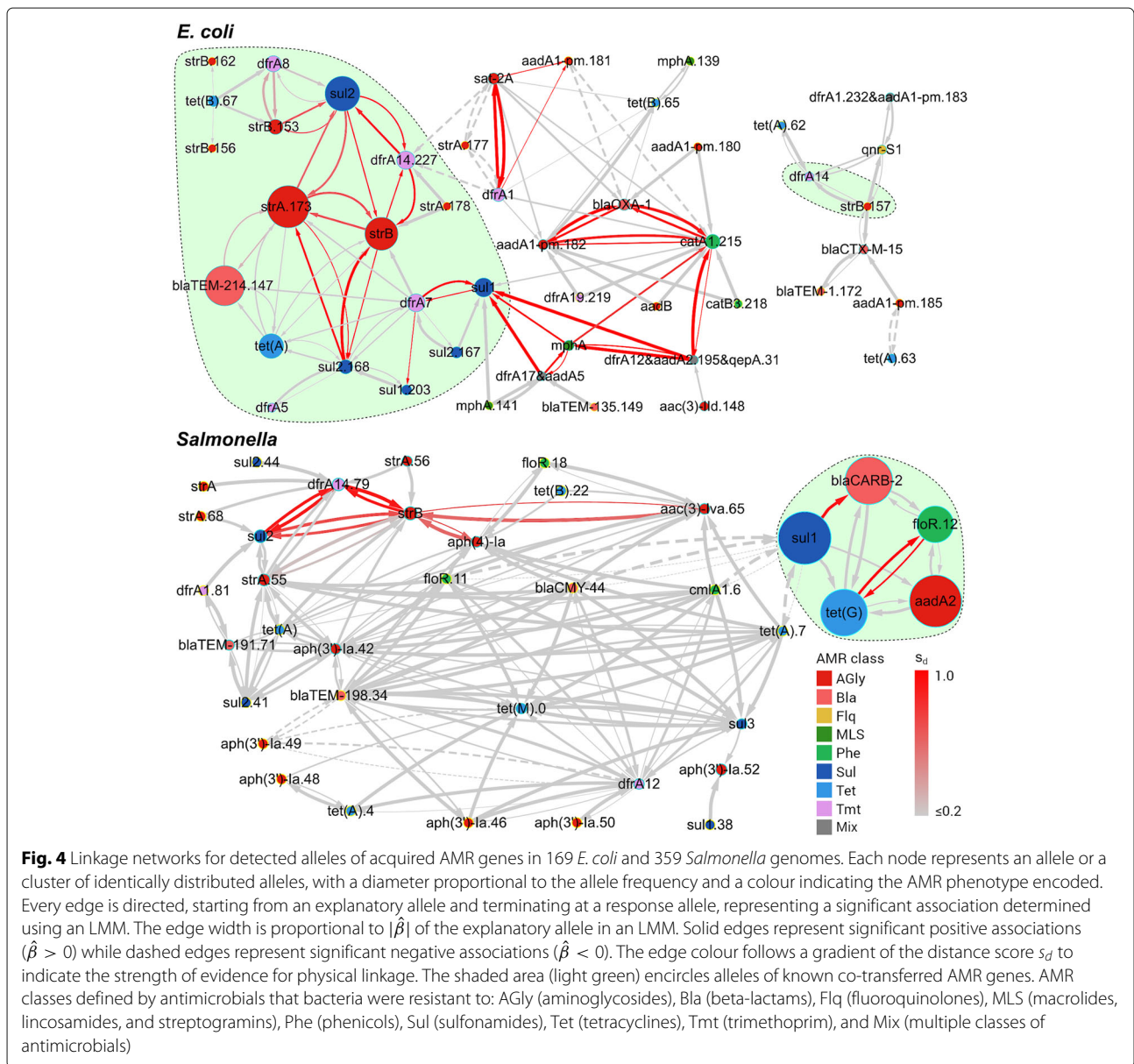
We found that inaccurate SPDs measured across two nodes of *Salmonella* assembly graphs (Table s7) were caused by chimeric alleles (that is, highly similar alleles of the same gene were mistakenly assembled into one allele) in the assembly, owing to the limited capacity of short reads in resolving repeats. We addressed this issue through increasing the threshold for both the nucleotide identity and query coverage of BLAST hits from 95% to 99% (BLAST hits to chimeric alleles were hence discarded) and thereby improved the accuracy rate from two genomes at the cost of reducing distance measurability (Table s8). Accordingly, we applied this adjusted threshold to the measurement of SPDs between acquired AMR genes in the 358 *Salmonella* draft genomes. Besides, we directly calculated SPDs for strain DT104, whose complete genome is publicly accessible on GenBank.

SPDs between alleles of acquired AMR genes From *de novo* assemblies of the 169 *E. coli* genomes, we obtained 1,550 SPDs from 301 allele pairs that were tested for associations (Hence alleles of each pair did not have identical presence-absence status across genomes) and 20 SPDs from nine pairs of identically distributed alleles. The largest SPD was 59,433 bp (measured across 17 nodes) and the greatest node traversal to measure an SPD was across 39 nodes (yielding an SPD of 4,628 bp). The exclusion of SPDs measured across more than two nodes resulted in 673 (43%) SPDs reliably measured for 163 tested allele pairs and 18 (90%) SPDs for eight pairs of identically distributed alleles. From *de novo* assemblies of the 359 *Salmonella* genomes, we obtained 2,880 SPDs from 224 allele pairs tested for associations, including 2,322 SPDs between alleles of five SG11-borne AMR genes (Table s9), and obtained another 10 SPDs from five pairs of identically distributed alleles. The largest SPD was 710,475 bp (measured across 44 nodes) and the greatest node traversal to measure an SPD was across 57 nodes (yielding an SPD of 687,051 bp). We saw large SPDs (> 56 kbp, which were extraordinarily larger than the others) when the node number exceeded 14. The exclusion of SPDs that were greater than 250 kbp and measured across more than two nodes resulted in 994 SPDs from 59 tested allele pairs and eight SPDs from three pairs of identically distributed alleles.

Overall, positively associated alleles of acquired AMR genes in *E. coli* and *Salmonella* genomes showed higher measurability of SPDs than those measured between negatively associated alleles. Specifically, as for pairs of positively associated alleles, 64 (72%) out of 89 pairs in *E. coli* and 25 (26%) out of 97 pairs in *Salmonella* had at least two SPDs measured, respectively (Tables s10 and s11). Moreover, after removing SPDs that were greater than 250 kbp and measured across SPDs more than two nodes, 31 (35%) out of the 89 allele pairs in *E. coli* and 10 (10%) out of the 97 allele pairs in *Salmonella* had distance measurability above 75%. On the contrary, no SPD was measurable between negatively associated alleles as these alleles did not co-occur in any genome.

Linkage networks showing support of SPDs to HGcoT

For each species, we created a linkage network, in which nodes represent alleles or clusters of identically distributed alleles of acquired AMR genes and directed edges indicate significant associations obtained from LMMs (Fig. 4). The edge width is proportional to an estimated effect size $\hat{\beta}$, solid lines and dashed lines indicate positive associations and negative ones, respectively, and the edge colour indicates the distance score s_d . No filter was applied to distance scores. The linkage network for *E. coli* consisted of 122 edges linking 46 nodes corresponding to 52 alleles of 26 acquired AMR genes. The distance score followed a bimodal distribution, with 57 out of 83 allele pairs (69%) having $s_d = 0$ (84 edges) and 23 (28%) allele



pairs having $s_d > 0.5$ (38 edges). Considering only the edges that yielded $s_d > 0.5$ and connected alleles encoding distinct kinds of resistance phenotypes, aminoglycoside resistance alleles linked to 14 alleles — the largest number of connections, followed by sulfonamide resistance alleles (linked to 11 alleles). By contrast, bla_{OXA-1} , the only beta-lactam resistance allele having edges with distance scores above 0.5, was linked to two alleles ($aadA1$ -*pm*.182 and $catA1$.215). Notably, as shown in Fig. 4, five alleles ($dfxA14$.227, $strB$, $sul2$.168, $strA$.173, and $sul2$) formed a cluster that was interconnected by bidirectional edges with high distance scores ($s_d > 0.6$).

The linkage network for *Salmonella* consisted of 37 alleles (21 AMR genes) connected by 162 edges (Fig. 4). No identically distributed alleles could be collectively represented by a single node in this network due to absence of perfect measurability or consistency of SPDs. The distance score again followed a bimodal distribution, with 95 out of 104 (91%) pairs having $s_d < 0.05$, and seven (7%) pairs (altogether, nine alleles) having $s_d > 0.5$ (13 edges), all of which corresponded to significant positive associations. Considering only the 13 edges having distance scores above 0.5, $strB$ (aminoglycoside resistance) linked to the largest number of alleles (four, altogether), followed by $sul2$ and $dfxA14$.79, each connected to two alleles.

Reasons for inconsistency in measured physical distances In both the *E. coli* and *Salmonella* data sets, a lack of consistency in SPDs (namely, $c = 0$) measured between several positively associated alleles of acquired AMR genes was observed (Tables s10 and s11). We investigated this issue based on GeneMates outputs and identified two common explanations.

First, in many cases diverse genetic structures were found carrying the same combination of alleles of AMR genes. For example, based on assembly graphs of *E. coli* genomes, we recovered six distinct genetic structures linking alleles $bla_{TEM-214}$.147 and $tet(A)$, which showed significant positive associations in both LMMs and PLMs but had six distinct SPDs (Figure s15). We found that these SPDs followed a lineage-specific distribution. As illustrated in Figure s16, allele $bla_{TEM-214}$.147 was carried by transposon Tn2, which was common to all the six structures, either in its complete or truncated form. The variety of insertion sites and orientations of this transposon relative to allele $tet(A)$ in *E. coli* genomes, as well as plausible gene gain/loss events (inferred from structural comparisons), resulted in differences in SPDs measured between these two alleles.

Second, the GeneMates algorithm depreciates physical distances showing IBD for scoring the distance consistency (Section 3.3.2 of Additional file 1). For instance, two alleles of SGI1-borne AMR genes $sul1$ and $aadA2$

in *Salmonella* were frequent amongst the 359 *Salmonella* isolates, with an occurrence count of 328 (91%) and 323 (90%), respectively (Table s9). Co-occurrence of these two alleles were also frequent, with a count of 318 (89%) in total. SPDs between $sul1$ and $aadA2$ were obtained from genome assemblies of 295 (93%) out of the 318 isolates where the alleles were co-occurring. After removing the only SPD measured across more than two nodes (504 bp, across three nodes), we obtained 294 SPDs, consisting of 293 SPDs measured in either contigs or assembly graphs and one SPD measured in the complete chromosome genome of reference strain DT104. All the filtered SPDs were 504 bp, except the one from the complete genome (9,964 bp). Despite this consistency in SPDs, the consistency score $c = 0$ as all the 294 SPDs were obtained from the same lineage highlighted in Figures s17a (IBD probability of reliable SPDs: 95%) and were possibly resulted from the same genetic and assembly structures related to these two alleles (Figures s17b–d).

Validation of GeneMates

We validated our approach through identification of known and novel physical clusters of mobile AMR genes in the example data sets. Networks were constructed at the allele level for these genes using GeneMates function *findPhysLink*.

Identifying known clusters of mobile AMR genes The first set of positive controls in our validation study consisted of 28 pairs of AMR genes that are known to be co-mobilised by MGEs between *E. coli* genomes [19]. As shown in Table 1, LMMs and PLMs identified significant positive associations at the allele level in 19 (68%) and 16 (57%) pairs, respectively. The second set of positive controls for validation consisted of five AMR genes ($aadA2$, $floR$, $tet(G)$, bla_{CARB-2} , and $sul1$) that are co-localised in the acquired multidrug-resistant element SGI1 in *Salmonella* genomes [20]. For these genes, LMMs and PLMs identified significant positive associations between eight and ten allele pairs, respectively (Table 2). The exclusion of allele pairs having $s_d < 0.6$ led to a substantial reduction of co-mobilisation candidates, with 12 out of 33 (36%) allele pairs in *E. coli* genomes and two out of 20 (10%) allele pairs in *Salmonella* genomes passed this filter.

Identifying novel clusters of AMR genes Some edges in linkage networks suggested novel physical linkage between several alleles of AMR genes in *E. coli* and *Salmonella* genomes. For instance, 20 novel edges in the linkage network for *E. coli* (Fig. 4) had the maximum association score $s_a = 1$ ($0.12 < \hat{\beta} \leq 1$ in LMMs) and high distance scores ($0.75 \leq s_d \leq 1$), and these edges formed five maximal cliques each consisting of three alleles. Similarly, a three-allele clique showing high distance

Table 1 A comparison of significant positive associations in the linkage network for *E. coli* genomes to known co-localisation of mobile AMR genes

Gene_1	Gene_2	Allele_1	Allele_2	LMM	PLM	s_d	p_{IBD}
tet(B)	strB	tet(B).67	strB.153	●	○	0	-
		tet(B).67	strB.156	●	○	0	-
		tet(B).67	strB.162	●	○	0	-
tet(B)	sul2	-	-	○	○	-	-
tet(B)	strA	-	-	○	○	-	-
dfrA14	strB	dfrA14.227	strB	●	●	0.93	0.5
		dfrA14	strB.157	●	○	0	-
dfrA14	sul2	dfrA14.227	sul2	●	●	0.87	0.5
		dfrA5	sul2.168	●	○	0	-
dfrA7	strA	dfrA7	strA.173	●	●	0	0.35
dfrA7	strB	dfrA7	strB	●	●	0	0.35
dfrA7	sul2	dfrA7	sul2.167	●	○	0	-
		dfrA7	sul2.168	●	●	0	-
dfrA7	sul1	dfrA7	sul1	●	●	1	0.46
		dfrA7	sul1.203	●	●	1	0.04
		dfrA17	sul1	●	○	1	0.26
dfrA7	tet(A)	dfrA7	tet(A)	●	○	0	-
tet(A)	strA	tet(A)	strA.173	●	○	0	1
tet(A)	strB	tet(A)	strB	●	○	0	1
tet(A)	sul2	tet(A)	sul2.168	●	○	0	-
tet(A)	strA	-	-	○	○	-	-
bla _{TEM-198}	tet(A)	bla _{TEM-214} .147	tet(A)	●	●	0	0.5
bla _{TEM-198}	sul1	-	-	○	○	-	-
bla _{TEM-198/191}	strA	bla _{TEM-214} .147	strA.173	●	○	0.29	0.5
bla _{TEM-198/191}	strB	-	-	○	○	-	-
bla _{TEM-198/191}	sul2	-	-	○	○	-	-
dfrA8	strA	dfrA8	strA.173	○	○	0.26	0.12
dfrA8	strB	dfrA8	strB.153	●	●	0.33	0.12
dfrA8	sul2	dfrA8	sul2	●	●	0.23	0.07
strA	strB	strA.173	strB	●	○	0.68	0.5
		strA.173	strB.153	○	●	0.67	0.5
		strA.178	strB	●	○	0	-
strA	sul2	strA.173	sul2	●	●	0.62	0.5
		strA.173	sul2.168	●	○	1	0.38
strB	sul2	strB	sul2.168	●	●	1	0.38
		strB.153	sul2	●	○	0.65	0.5
		strB	sul2	●	○	0.84	0.5
sul1	strA	-	-	○	○	-	-
sul1	strB	-	-	○	○	-	-
sul1	sul2	sul1	sul2.167	○	○	0	-
		sul1.203	sul2.168	●	●	0	-

Co-localisation of AMR genes in MGEs were previously determined by Ingle, et al. [19]. Each pair of significantly associated alleles (denoted by alleles 1 and 2, regardless of their roles in a linear model) is also identifiable in the comparative network shown in Fig. 3. Directionality of associations is omitted in this table, hence each pair of alleles only appears once in the table, although the alleles may mutually associated in linear models. Note that an AMR gene may have multiple alleles (whose names are listed in column Allele_1 or Allele_2) or no allele (denoted by a dash sign in the table) present in a network. Abbreviations: LMM, linear mixed model; PLM, penalised logistic model; IBD: identity by descent. Symbols indicating whether a significant association is identified by either an LMM or a PLM: (●), yes; (○), no. s_d : the distance score, which takes into account the distance consistency, measurability, and the probability of IBD. p_{IBD} : an estimate of the probability that APDs used for calculating the s_d are in IBD. This probability does not exist when no APD is available

scores ($0.7 \leq s_d \leq 0.9$) was identified in the linkage network for *Salmonella* (Fig. 4). As a further validation of GeneMates, we investigated two maximal cliques through inferring their plausible genetic structures and vectors from genome assemblies.

The first clique consisted of alleles *dfrA1*, *aadA1-pm.181*, and *sat-2A* detected in *E. coli* genomes. Between these alleles, LMMs identified significant positive associations (Figure s18a), and SPDs (Table s12) showed complete identity as well as perfect measurability (namely, $s_d = 1$ for every edge of this clique). Co-occurrence of all three alleles was identified in two distantly related genomes (Figure s18b). Assembly graphs revealed co-localisation of these alleles in the same array of gene cassettes in a class-2 integron (Fig. 5a, b). We confirmed that this

Table 2 LMM-based significant associations between five alleles of AMR genes in SG11

Gene_Y	Gene_X	Allele_Y	Allele_X	$\hat{\lambda}$	LMM	PLM	s_d	p_{IBD}
<i>floR</i>	<i>sul1</i>	<i>floR.12</i>	<i>sul1</i>	64.98	○	●	0.03	0
<i>sul1</i>	<i>floR</i>	<i>sul1</i>	<i>floR.12</i>	$\geq 10^5$	○	●	0.03	0
<i>aadA2</i>	<i>floR</i>	<i>aadA2</i>	<i>floR.12</i>	137.49	●	●	0.02	0
<i>floR</i>	<i>aadA2</i>	<i>floR.12</i>	<i>aadA2</i>	82.73	●	●	0.02	0
<i>bla_{CARB}</i>	<i>floR</i>	<i>bla_{CARB-2}</i>	<i>floR.12</i>	110.28	●	●	0.04	0
<i>floR</i>	<i>bla_{CARB}</i>	<i>floR.12</i>	<i>bla_{CARB-2}</i>	57.27	●	●	0.04	0
<i>floR</i>	<i>tet(G)</i>	<i>floR.12</i>	<i>tet(G)</i>	64.45	●	●	0.99	0.5
<i>tet(G)</i>	<i>floR</i>	<i>tet(G)</i>	<i>floR.12</i>	91.3	●	●	0.99	0.5
<i>bla_{CARB}</i>	<i>sul1</i>	<i>bla_{CARB-2}</i>	<i>sul1</i>	70.79	●	●	0.92	0.85
<i>sul1</i>	<i>bla_{CARB}</i>	<i>sul1</i>	<i>bla_{CARB-2}</i>	$\geq 10^5$	○	●	0.92	0.85
<i>aadA2</i>	<i>bla_{CARB}</i>	<i>aadA2</i>	<i>bla_{CARB-2}</i>	130.1	○	●	0	0
<i>bla_{CARB}</i>	<i>aadA2</i>	<i>bla_{CARB-2}</i>	<i>aadA2</i>	123.6	○	●	0	0
<i>bla_{CARB}</i>	<i>tet(G)</i>	<i>bla_{CARB-2}</i>	<i>tet(G)</i>	244.33	○	●	0.04	0
<i>tet(G)</i>	<i>bla_{CARB}</i>	<i>tet(G)</i>	<i>bla_{CARB-2}</i>	160.14	●	●	0.04	0
<i>sul1</i>	<i>tet(G)</i>	<i>sul1</i>	<i>tet(G)</i>	$\geq 10^5$	○	●	0.03	0
<i>tet(G)</i>	<i>sul1</i>	<i>tet(G)</i>	<i>sul1</i>	55.9	○	●	0.03	0
<i>aadA2</i>	<i>tet(G)</i>	<i>aadA2</i>	<i>tet(G)</i>	118.93	●	●	0.04	0
<i>tet(G)</i>	<i>aadA2</i>	<i>tet(G)</i>	<i>aadA2</i>	77.76	●	●	0.04	0
<i>aadA2</i>	<i>sul1</i>	<i>aadA2</i>	<i>sul1</i>	73.42	●	●	0	0.95
<i>sul1</i>	<i>aadA2</i>	<i>sul1</i>	<i>aadA2</i>	$\geq 10^5$	○	●	0	0.95

Allele_Y and Allele_X denote the response allele and explanatory allele in an LMM $Y \sim X$, respectively. An association is denoted by a filled circle in the column LMM when it is significant, otherwise, an unfilled circle is drawn. Directionality is shown in this table for comparing the value of $\hat{\lambda}$, which denotes a REML estimate of parameter λ for evaluating structural random effects in an LMM. s_d : score of APDs. p_{IBD} : an estimate of the probability that APDs used for calculating the s_d are in IBD

integron was carried by variants (100% coverage and 99% nucleotide identity) of transposon Tn7 (GenBank accession: KX159451). Moreover, each Tn7 variant was interrupted by a distinct insertion sequence (IS), as illustrated in Fig. 5c. In summary, we saw strong evidence for transposon-mediated co-transfer of these alleles between *E. coli* lineages.

The second clique consisted of alleles *strA.55*, *strB*, *dfrA14.79*, and *sul2* detected in *Salmonella* genomes (Fig. 4). Both LMMs and PLMs identified that associations between all of these alleles were significantly positive ($0.69 \leq \hat{\beta} \leq 0.98$ in LMMs). Edges between these alleles had distance scores between 0.7 and 0.9 except edges linking *strA.55* ($0 \leq s_d \leq 0.25$). As shown in Fig. 6a, these four alleles co-occurred in 13 *Salmonella* genomes from distantly related clades. Using Bandage and nucleotide BLAST, we confirmed co-localisation of these alleles in a 3,084 bp region that was present in 10 out of the 13 genomes with a 100% nucleotide identity and coverage, and these 10 genomes were sparsely distributed across the phylogenetic tree in Fig. 6a. Furthermore, we saw an insertion of allele *dfrA14.79* into *strA.55*, splitting the latter allele into two segments that covered 65.8% and 34.5% of its original length, respectively. In the assembly graph of one of the 10 genomes, ERR026101, we found the 3,084 bp multidrug-resistance (MDR) region in a 6,790 bp node, which appeared as a self-circularised sequence independent to other graph components (Fig. 6b). Using megaBLAST under its default parameters, a sequence search of this MDR region against the NCBI nucleotide database of the *Enterobacteriaceae* group (taxid: 91347,

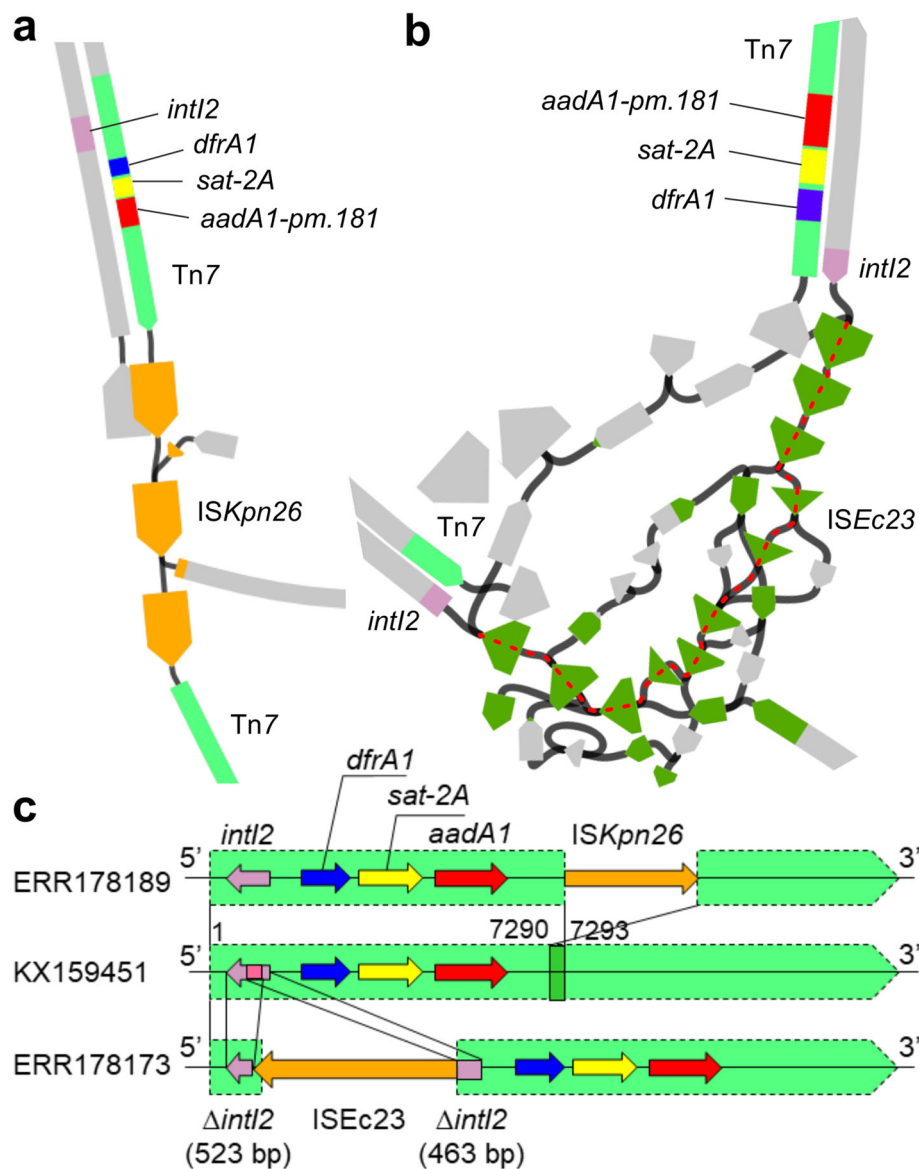
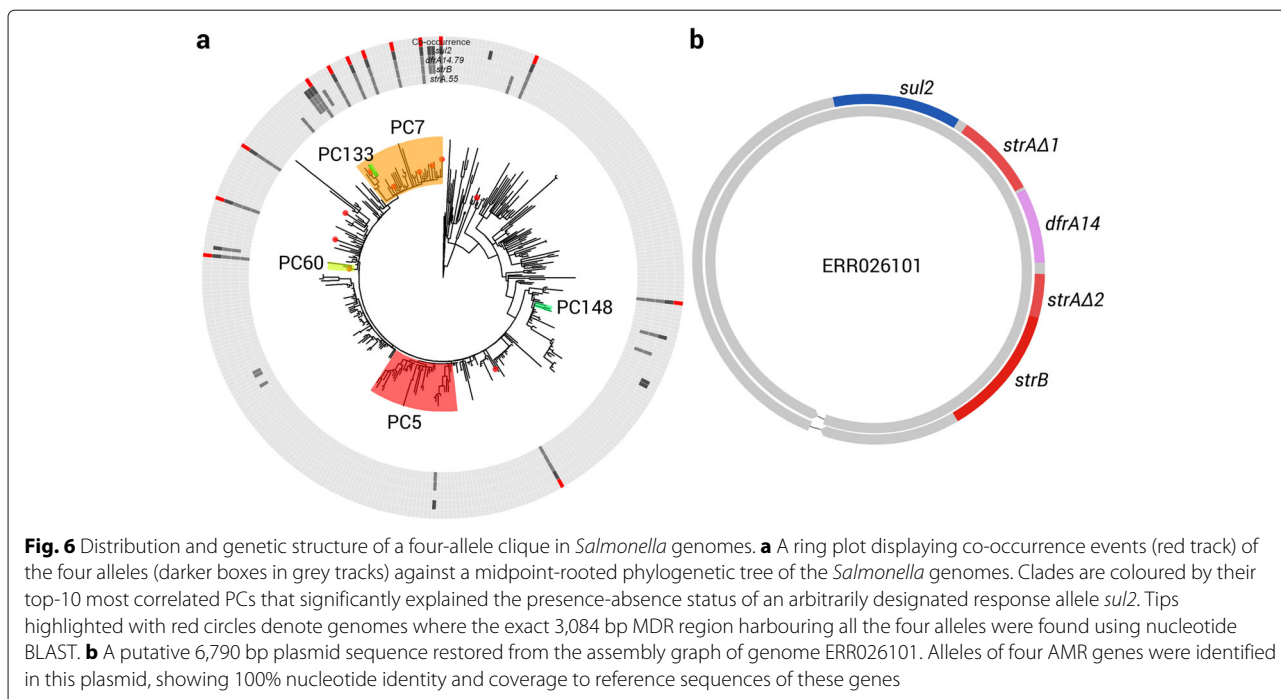


Fig. 5 Putative physical linkage between three alleles of AMR genes in two *E. coli* genomes. **a** A path constituting an inferred MGE in the assembly graph of genome ERR178189. This plot is drawn in Bandage with a double-strand style, where the orientation of each DNA strand is indicated by an arrow-like node end. The width of each node is proportional to its read depth determined by Unicycler. Some nodes not contributing to any MGE-related path were deleted from the original assembly graph for visual conciseness. **b** A path constituting the other inferred MGE (following the red dashed line) in the assembly graph of genome ERR178173. This plot was drawn in the same way as panel **a**. **c** Alignment of the two Tn7 variants we identified in *E. coli* genomes (ERR178189 and ERR178173) to a reference Tn7 sequence (GenBank accession: KX159451, denoted by green shaded areas). Two direct repeats flanking the ISs, including inverted repeats, are denoted by green and pink boxes, respectively. Reference DNA sequences of ISKpn26 (1,196 bp) and ISEc23 (2,532 bp) were retrieved from database ISFinder [29] in January 2018. Each IS in the resolved region showed a 100% coverage and 99% nucleotide identity to its reference

accessed in April, 2018) showed exact matches (100% nucleotide identity and coverage) to a known and widely distributed MDR plasmid pCERC1 (GenBank accession: JN012467) as well as a number of plasmids widely distributed in bacteria of *Enterobacteriales* (Table s13). Therefore, this MDR region was shared by a great variety of plasmids.

Discussion

GeneMates implements a novel network approach to detection of intraspecies HGcoT between bacteria. Compared to existing methods relying on co-occurrence counts, association coefficients, or simple regression models of bacterial genes, GeneMates enables us to test for gene-gene associations when controlling for



population structure, the main confounding factor in bacterial GWAS [26], through incorporating PCs into LMMs. Results in “Effects of controlling for population structure” section show that LMMs of GeneMates retain statistical power of association tests, which is in line with equivalent LMMs from literature (cf., Section 3.4.1 in Additional file 1) [10, 27]. In our examples, association networks constructed using GeneMates reveal extensive associations between alleles of horizontally acquired AMR genes (Fig. 3). Moreover, as expected, the majority of p -values from association tests became greater after correcting for population structure using LMMs, while the other p -values saw a reduction, indicating increased significance of associated alleles (Fig. 2). LMMs provide us with another advantage: we only need to estimate three parameters (β , γ , and τ) besides the intercept term α to fit a model, thereby circumventing the problem of overfitting as well as relaxing the requirement for sample sizes. Notably, long-tailed curves of cumulative percentages of total genetic variations captured by PCs (Figure s9) indicate the necessity of including all PCs obtained from the cgSNP matrix for accurate modelling in association analysis.

On the grounds of examples where stable structures of acquired AMR genes are shared between bacteria via MGEs in a short period of time [16, 19, 20, 28], we have implemented another innovation in GeneMates — the evaluation of APD consistency for evidence of HGcoT. Since the physical distances between two loci in bacterial genomes are correlated with both locus co-occurrence

and PCs representing population structure, APDs cannot be directly incorporated into LMMs or PLMs as additional covariates. Instead, GeneMates scores the variation of APDs while taking the population structure into account. As for APDs, it is self-evident that they can be precisely calculated from genetic coordinates in finished-grade genome assemblies, which remain the minority of available bacterial genome sequences. In order to overcome some of the limitations of measuring APDs in draft genomes, we measured APDs in the form of SPDs in assembly graphs, and developed a simulation-based approach to determining reliability filters of SPDs. Using this approach, the accuracy of SPDs from *de novo* assembly graphs of reference multidrug-resistant *E. coli* or *S. Typhimurium* genomes consistently exceeded 90% when the distances were only measured across one or two nodes (Figures s11 and s12), implying a universal filter for other bacterial genomes. Nevertheless, as summarised in “Effects of adding APDs to association networks” section and displayed in Fig. 4, short-read genome assemblies intrinsically confine both the measurability and accuracy of SPDs, causing a loss of evidence for real physical linkage in HGT.

In “Validation of GeneMates” section, the identification of known and novel physical clusters of mobile AMR genes suggests that maximal cliques in linkage networks are useful starting points for recovering structures of horizontally co-transferred genes. A plausible reason is that loci in strong physical linkage tend to predict the presence of each other (namely, mutual

positive associations) in HGT, and strong physical linkage is often related to close physical proximity, which in turn increases the measurability of their physical distances, leading to a greater distance score given the same consistency score. In practice, users may apply other filters to the network to identify edges addressing specific questions.

GeneMates also offers a framework (Fig. 1) for further development, such as adding new modules and introducing other statistical models. Our methodology and analysis demonstrated in example studies are applicable to other kinds of acquired genes in haploid genomes of the same species, as long as we can accurately determine alleles of interested genes in each genome. Since GeneMates is designed specifically to overcome challenges of identifying intraspecies HGcoT, we did not test its performance in detecting interspecies HGcoT. GeneMates could theoretically be applied to analyse data sets including multiple species. However, users should consider the input SNP matrix carefully, as it determines what scale of population structure is captured in the covariance matrix used in LMMs. In cross-species analyses, the SNP matrix would necessarily be restricted to core-genome sites shared across species, which may represent only a small fraction of individual genomes (depending on the taxonomic breadth of included species). In such cases the PCs would likely capture differences between species but may not effectively capture intraspecies population structure. Furthermore, the inability of short reads to resolve repeats, either through read mapping or *de novo* assembly, may cause false negatives and errors in allele calls when homologues of a target gene coexist in a genome. Therefore, we expect a better performance of GeneMates in analysing high-quality complete genomes or long-read sequencing data that are able to resolve at least most of the repeats. Finally, biological experiments are necessary as a gold standard for validating candidates of HGcoT.

Conclusions

We have developed R package GeneMates and helper scripts for analysing associations between alleles of acquired genes and for inferring physical linkage between these alleles in HGT. We have also demonstrated utilities of this package using publicly available WGS data of 169 *E. coli* isolates and 359 *Salmonella* isolates. In our study, functions of GeneMates identified clusters of co-transferred AMR alleles in known and novel MGEs, enabling further investigations of HGcoT amongst a wider range of bacterial species. GeneMates differs from contemporary methods for bacterial GWAS in three aspects. First, it focuses on gene-gene associations rather than genotype-to-phenotype associations. Second, it only performs association tests for acquired genes rather than genome-wide SNPs. Third, it evaluates

evidence of physical distances between associated loci for inference of physical linkage, although users may opt to turn this utility off to only analyse gene-gene associations. GeneMates offers a scalable and versatile approach that is readily applicable to various kinds of horizontally acquired genes. It is, however, confined by limitations of short-read genome assembly, and its power will increase in the future as we are accumulating complete genomes and enhancing our ability in resolving repeats using sequencing data.

Availability and requirements

- **Project name:** GeneMates
- **Project home page:** github.com/wanyuac/GeneMates
- **Programming language:** R
- **Operating system(s):** platform independent
- **Other requirements:** GEMMA v0.96
- **Licence:** Apache License, Version 2.0

Supplementary information

Supplementary information accompanies this paper at <https://doi.org/10.1186/s12864-020-07019-6>.

Additional file 1: A PDF document consisting of supplementary tables and figures, and a rigorous mathematical justification of statistical methods implemented in package GeneMates.

Additional file 2: An Excel spreadsheet of sample information, including database accessions, sampling locations and years, sequencing layout, and so forth. It consists of four tables with names *Ecoli_reads* (Illumina read sets of *E. coli* isolates), *STyphimurium_reads* (Illumina read sets of *Salmonella* isolates), *MDR_Ecoli_ref* (GenBank records for complete genomes of 10 multidrug-resistant *E. coli* isolates), and *MDR_STyphimurium_ref* (GenBank records for complete genomes of 10 multidrug-resistant *S. Typhimurium* isolates).

Additional file 3: An Excel spreadsheet about prophages identified in reference chromosome sequences and statistics of read mapping and *de novo* genome assemblies.

Acknowledgements

We thank Kelly Wyres and David Edwards (Monash University), Claire Gorrie (University of Melbourne), and Gittan Blezer (University of Melbourne and Utrecht University) for their suggestions on data collection, quality control and read mapping. We appreciate Guoqi Qian (University of Melbourne) for his comments on our statistical methods. We also express our gratitude to the maintenance team of GEMMA (github.com/genetics-statistics/GEMMA) and the GitHub community for constructive discussions about software usage and statistics.

Authors' contributions

YW derived algebra of our methodology, developed and validated GeneMates and its helper scripts, curated and analysed the test data, interpreted results, and prepared the manuscript for publication. RRW developed and implemented the method for measuring APDs in assembly graphs, and provided innovative suggestions on short-read simulation and sequence analysis. KEH and MI conceived the methodology and this project. KEH, MI and JZ supervised this project and substantially contributed to research strategy and result interpretation. DJJ participated in study design and result interpretation. All authors read and approved the final manuscript as well as contributing to revisions.

Funding

YW was supported by a Melbourne International Research Scholarship from the University of Melbourne. KEH was supported by the Bill & Melinda Gates

Foundation, Seattle and a Senior Medical Research Fellowship from the Viertel Foundation of Australia.

Availability of data and materials

All sequence data analysed during this study are publicly available in NCBI databases. See Additional file 2 for accession numbers.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, University of Melbourne, Parkville, Victoria, 3010 Australia. ²Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne, Victoria, 3004 Australia. ³School of Computing and Information Systems, University of Melbourne, Parkville, Victoria, 3010 Australia. ⁴Microbiological Diagnostic Unit Public Health Laboratory, Department of Microbiology and Immunology, University of Melbourne at The Peter Doherty Institute for Infection and Immunity, Parkville, Victoria, 3010 Australia. ⁵National Centre for Epidemiology and Population Health, Australian National University, Canberra, Australian Capital Territory, 2601 Australia. ⁶Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne, Victoria, 3004 Australia. ⁷Cambridge Baker Systems Genomics Initiative, Department of Public Health and Primary Care, University of Cambridge, Cambridge, England, CB1 8RN UK. ⁸Department of Infection Biology, London School of Hygiene & Tropical Medicine, London, WC1E 7HT UK.

Received: 9 March 2020 Accepted: 20 August 2020

Published online: 24 September 2020

References

- Jain R, Rivera MC, Moore JE, Lake JA. Horizontal Gene Transfer Accelerates Genome Innovation and Evolution. *Mol Biol Evol.* 2003;20(10):1598–602. <https://doi.org/10.1093/molbev/msg154>.
- Popa O, Dagan T. Trends and barriers to lateral gene transfer in prokaryotes. *Curr Opin Microbiol.* 2011;14(5):615–23. <https://doi.org/10.1016/j.mib.2011.07.027>.
- Medini D, Donati C, Tettelin H, Massignani V, Rappuoli R. The microbial pan-genome. *Curr Opin Genet Dev.* 2005;15(6):589–94. <https://doi.org/10.1016/j.gde.2005.09.006>.
- Frost LS, Leplae R, Summers AO, Toussaint A. Mobile genetic elements: the agents of open source evolution. *Nat Rev Micro.* 2005;3(9):722–32.
- Baker-Austin C, Wright MS, Stepanauskas R, McArthur JV. Co-selection of antibiotic and metal resistance. *Trends Microbiol.* 2006;14(4):176–82. <https://doi.org/10.1016/j.tim.2006.02.006>.
- Chang H-H, Cohen T, Grad YH, Hanage WP, O'Brien TF, Lipsitch M. Origin and Proliferation of Multiple-Drug Resistance in Bacterial Pathogens. *Microbiol Mol Biol Rev.* 2015;79(1):101–16. <https://doi.org/10.1128/mbr.00039-14>.
- Land M, Hauser L, Jun S-R, Nookaew I, Leuze MR, Ahn T-H, Karpinets T, Lund O, Kora G, Wassenaar T, Poudel S, Ussery DW. Insights from 20 years of bacterial genome sequencing. *Funct Integr Genomics.* 2015;15(2):141–61. <https://doi.org/10.1007/s10142-015-0433-4>.
- Suzuki M, Shibayama K, Yahara K. A genome-wide association study identifies a horizontally transferred bacterial surface adhesin gene associated with antimicrobial resistant strains. *Sci Rep.* 2016;6:37811. <https://doi.org/10.1038/srep37811>.
- Kang HM, Zaitlen NA, Wade CM, Kirby A, Heckerman D, Daly MJ, Eskin E. Efficient Control of Population Structure in Model Organism Association Mapping. *Genetics.* 2008;178(3):1709–23.
- Earle SG, Wu C-H, Charlesworth J, Stoesser N, Gordon NC, Walker TM, Spencer CCA, Iqbal Z, Clifton DA, Hopkins KL, Woodford N, Smith EG, Ismail N, Llewelyn MJ, Peto TE, Crook DW, McVean G, Walker AS, Wilson DJ. Identifying lineage effects when controlling for population structure improves power in bacterial association studies. *Nat Microbiol.* 2016;1:16041.
- Hoffman GE. Correcting for Population Structure and Kinship Using the Linear Mixed Model: Theory and Extensions. *PLoS ONE.* 2013;8(10):75707.
- McVean G. A Genealogical Interpretation of Principal Components Analysis. *PLoS Genetics.* 2009;5(10):1000686.
- Blomberg SP, Lefevre JG, Wells JA, Waterhouse M. Independent Contrasts and PGLS Regression Estimators Are Equivalent. *Syst Biol.* 2012;61(3):382–91.
- Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* 2003;13(11):2498–504.
- Zhou X, Carbonetto P, Stephens M. Polygenic Modeling with Bayesian Sparse Linear Mixed Models. *PLoS Genet.* 2013;9(2):1003264.
- Anantham S, Hall RM. pCERC1, a Small, Globally Disseminated Plasmid Carrying the *dfraA14* Cassette in the *strA* Gene of the *sul2-strA-strB* Gene Cluster. *Microb Drug Resist.* 2012;18(4):364–71. <https://doi.org/10.1089/mdr.2012.0008>.
- Csardi G, Nepusz T. The igraph software package for complex network research. *InterJournal.* 2006;Complex Systems:1695.
- Kotloff KL, Nataro JP, Blackwelder WC, Nasrin D, Farag TH, Panchalingam S, Wu Y, Sow SO, Sur D, Breiman RF, Faruque ASG, Zaidi AKM, Saha D, Alonso PL, Tamboura B, Sanogo D, Onwuchekwa U, Manna B, Ramamurthy T, Kanungo S, Ochieng JB, Omere R, Oundo JO, Hossain A, Das SK, Ahmed S, Qureshi S, Quadri F, Adegbola RA, Antonio M, Hossain MJ, Akinsola A, Mandomando I, Nhampossa T, Acácio S, Biswas K, O'Reilly CE, Mintz ED, Berkeley LY, Muhsen K, Sommerfelt H, Robins-Browne RM, Levine MM. Burden and aetiology of diarrhoeal disease in infants and young children in developing countries (the Global Enteric Multicenter Study, GEMS): a prospective, case-control study. *Lancet.* 2013;382(9888):209–22. [https://doi.org/10.1016/S0140-6736\(13\)60844-2](https://doi.org/10.1016/S0140-6736(13)60844-2).
- Ingle DJ, Levine MM, Kotloff KL, Holt KE, Robins-Browne RM. Dynamics of antimicrobial resistance in intestinal *Escherichia coli* from children in community settings in South Asia and sub-Saharan Africa. *Nat Microbiol.* 2018;3(9):1063–73. <https://doi.org/10.1038/s41564-018-0217-4>.
- Mather AE, Reid SWJ, Maskell DJ, Parkhill J, Fookes MC, Harris SR, Brown DJ, Coia JE, Mulvey MR, Gilmour MW, Petrovska L, de Pinna E, Kuroda M, Akiba M, Izumiya H, Connor TR, Suchard MA, Lemey P, Mellor DJ, Haydon DT, Thomson NR. Distinguishable Epidemics Within Different Hosts of the Multidrug Resistant Zoonotic Pathogen *Salmonella* Typhimurium DT104. *Science.* 2013;341(6153):1514–7. <https://doi.org/10.1126/science.1240578>.
- Boyd D, Peters GA, Cloeckaert A, Boumedine KS, Chaslus-Dancla E, Imberechts H, Mulvey MR. Complete Nucleotide Sequence of a 43-Kilobase Genomic Island Associated with the Multidrug Resistance Region of *Salmonella enterica* Serovar Typhimurium DT104 and Its Identification in Phage Type DT120 and Serovar Agona. *J Bacteriol.* 2001;183(19):5725–32. <https://doi.org/10.1128/jb.183.19.5725-5732.2001>.
- Firth D. Bias reduction of maximum likelihood estimates. *Biometrika.* 1993;80(1):27–38.
- Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies from short and long sequencing reads. *PLoS Comput Biol.* 2017;13(6):1005595.
- Benson DA, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. *Nucleic Acids Res.* 2015;43(Database issue):30–5. <https://doi.org/10.1093/nar/gku1216>.
- Wick RR, Schultz MB, Zobel J, Holt KE. Bandage: interactive visualization of *de novo* genome assemblies. *Bioinformatics.* 2015. <https://doi.org/10.1093/bioinformatics/btv383>.
- Power RA, Parkhill J, de Oliveira T. Microbial genome-wide association studies: lessons from human GWAS. *Nat Rev Genet.* 2017;18(1):41–50.
- Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies. *Nat Genet.* 2012;44(7):821–4.
- Partridge SR, Tsafnat G, Coiera E, Iredell JR. Gene cassettes and cassette arrays in mobile resistance integrons. *FEMS Microbiol Rev.* 2009;33(4):757–84.
- Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M. ISfinder: the reference centre for bacterial insertion sequences. *Nucleic Acids Res.* 2006;34: <https://doi.org/10.1093/nar/gkj014>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.