



Recent mixing of *Vibrio parahaemolyticus* populations

Chao Yang¹ · Xiaoyan Pei² · Yarong Wu¹ · Lin Yan² · Yanfeng Yan¹ · Yuqin Song¹ · Nicola M. Coyle³ · Jaime Martinez-Urtaza⁴ · Christopher Quince⁵ · Qinghua Hu⁶ · Min Jiang⁶ · Edward Feil³ · Dajin Yang² · Yajun Song¹ · Dongsheng Zhou¹ · Ruifu Yang¹ · Daniel Falush³ · Yujun Cui¹

Received: 24 January 2019 / Revised: 4 June 2019 / Accepted: 7 June 2019 / Published online: 24 June 2019
© The Author(s), under exclusive licence to International Society for Microbial Ecology 2019

Abstract

Humans have profoundly affected the ocean environment but little is known about anthropogenic effects on the distribution of microbes. *Vibrio parahaemolyticus* is found in warm coastal waters and causes gastroenteritis in humans and economically significant disease in shrimps. Based on data from 1103 genomes of environmental and clinical isolates, we show that *V. parahaemolyticus* is divided into four diverse populations, VppUS1, VppUS2, VppX and VppAsia. The first two are largely restricted to the US and Northern Europe, while the others are found worldwide, with VppAsia making up the great majority of isolates in the seas around Asia. Patterns of diversity within and between the populations are consistent with them having arisen by progressive divergence via genetic drift during geographical isolation. However, we find that there is substantial overlap in their current distribution. These observations can be reconciled without requiring genetic barriers to exchange between populations if long-range dispersal has increased dramatically in the recent past. We found that VppAsia isolates from the US have an average of 1.01% more shared ancestry with VppUS1 and VppUS2 isolates than VppAsia isolates from Asia itself. Based on time calibrated trees of divergence within epidemic lineages, we estimate that recombination affects about 0.017% of the genome per year, implying that the genetic mixture has taken place within the last few decades. These results suggest that human activity, such as shipping, aquatic products trade and increased human migration between continents, are responsible for the change of distribution pattern of this species.

Introduction

Hospitable environments for particular marine microbes can be separated by large distances but whether dispersal barriers substantially influence their distribution and evolution is unknown. There are many studies of distribution of

marine microbes e.g. [1–4], but these typically survey patterns of macro-scale diversity. Differences in species level or genus level composition between locations are as likely to reflect environmental heterogeneity as dispersal, making the patterns difficult to interpret. Recent spread of microbes between continents has been documented for lineages that cause pathogenic infection of humans, including notorious clonal groups within *Vibrio parahaemolyticus* and *Vibrio cholerae* [5–8]. However, these lineages are unusual in using humans as vectors, which might facilitate long-range dispersal as in the case of the Haitian cholera outbreak [9]. We currently have little information on rates of spread of

These authors contributed equally: Chao Yang, Xiaoyan Pei

Supplementary information The online version of this article (<https://doi.org/10.1038/s41396-019-0461-5>) contains supplementary material, which is available to authorized users.

✉ Ruifu Yang
ruifuyang@gmail.com

✉ Daniel Falush
danielfalush@gmail.com

✉ Yujun Cui
cuiyujun.new@gmail.com

¹ State Key Laboratory of Pathogen and Biosecurity, Beijing Institute of Microbiology and Epidemiology, Beijing 100071, China

² National Center for Food Safety Risk Assessment, Beijing 100022, China

³ University of Bath, Bath, Somerset, UK

⁴ The Centre for Environment, Fisheries and Aquaculture Science, Dorset DT48UB, UK

⁵ Warwick Medical School, University of Warwick, Warwick, UK

⁶ Shenzhen Centre for Disease Control and Prevention, Shenzhen 518055, China

the great majority of environmental organisms that do not colonize large-animal hosts.

V. parahaemolyticus prefers warm coastal waters and causes gastroenteritis in humans [10, 11]. Disease outbreaks became common from 1990s and became global, due to spread of particular clones which are responsible for the great majority of recognized human infections [5], which has been attributed to factors such as El Niño and climate change [12–14]. It is not clear to what extent this pattern is historically typical, or whether it instead reflects better surveillance and different patterns of usage of marine resources. These clones also make up a small fraction of the *V. parahaemolyticus* diversity and only a very small fraction of strains isolated during environmental sampling.

The *V. parahaemolyticus* genome undergoes high rates of homologous recombination with other members of the species [15, 16]. We have previously found evidence that the species is split into several populations [16]. Members of a population are not necessarily particularly related at the clonal level, for example they may have recombined their entire genomes since sharing a common cellular ancestor, but they are nevertheless on average more similar to each other than to members of other populations because they have acquired DNA from a common gene pool. Previously we found evidence of a single population with a well-mixed gene pool in Asian waters and for one or more differentiated populations in the US [16].

Here we use a larger and more broadly sampled collection of 1103 genomes of environmental and clinical isolates to examine the global population structure of the species. We find four populations with different but overlapping modern geographic distributions as well as a small number of hybrid strains. Under the assumption that genetic exchange between strains is constrained by geography, the current extent of overlap is too high to maintain the populations as distinct entities and we conclude that most of this mixing is likely to have taken place within the last few decades, possibly coinciding with the recent emergence of pandemic clones.

Results and Discussion

Distribution of *V. parahaemolyticus* populations

We analyzed genomes of 1103 strains including 392 new strains sequenced as part of this study. These strains were isolated from a mixture of sources during 1951–2016, including 634 (57%) environmental isolates and 457 (41%) clinical isolates, and covered 24 countries (Supplementary Fig. 1 and Supplementary Table 1). Clonal relationships between strains can be inferred from identifying long stretches of near-identity, corresponding to regions of the genome that have been inherited by direct descent since the

strains shared a common ancestor, or, more simply, by the strains having a small number of SNP differences between them genome wide. Genetic distances between isolates are visualized using a Neighbor-Joining (NJ) tree [17] (Fig. 1a). Based on criterion of high nucleotide identity (pairwise SNP distance less than 2000 SNPs), the dataset contains 13 clonal groups (more than 10 isolates), with 10 associated with human disease and 3 associated with the environment (Supplementary Table 1).

Notwithstanding the larger number of clonal groups, the overall genetic diversity of clinical isolates is high and similar to environmental isolates (Supplementary Fig. 1d), thus we included them to define the overall population structure more accurately. The presence of clonally related strains in the data complicates analysis of deeper population structure, so we first removed closely related isolates to make a “non-redundant” dataset of 469 strains, in which no sequence differed by less than 2000 SNPs in the core genome (Methods). We used fineSTRUCTURE to identify distinct populations [18]. In total, 115 populations were identified in this initial analysis, however most comprised only two or three strains (Supplementary Fig. 2a). These are likely to be sets of strains that are clonally related, so we removed all but one from each group and reran fineSTRUCTURE. After several iterations of the same procedure, we identified four populations with between 10 and 217 members and two singletons (Supplementary Fig. 2d). These singletons might be hybrids or representatives of otherwise unsampled populations.

To test the reproducibility of the results, we reselected the panel of 469 strains twice and reran the same procedure and got very similar results (Supplementary Figs. 2b, c). We also tested the assignment of 1103 strains using a leave-one-out approach by calculating which of the four populations they had lowest average genetic distance with. In all cases, the assignment is concordant with that obtained by fineSTRUCTURE.

The fineSTRUCTURE populations formed clades in both NJ and Maximum Likelihood (ML) trees with a handful of exceptions (Fig. 1a, Supplementary Fig. 3). Both NJ and ML trees imply that CG2 strains clustered with VppAsia strains. However, inspection of the pairwise SNP distance and the chunk count distance show that although the individually closest strain is a VppAsia strain (S058), this appears to be simply a statistical fluke since S058 does not contribute many chunks to CG2 and on average CG2 shares many more chunks with and is much closer in average pairwise SNP distance to VppX strains than to VppAsia strains (Supplementary Fig. 4a). Similar results are obtained for the other exceptions (Supplementary Fig. 4b, c). Most noticeable, while the VppUS1 population is paraphyletic in the trees, each of the VppUS1 strains share more chunks with each other than with strains from the other populations (Supplementary Fig. 4d).

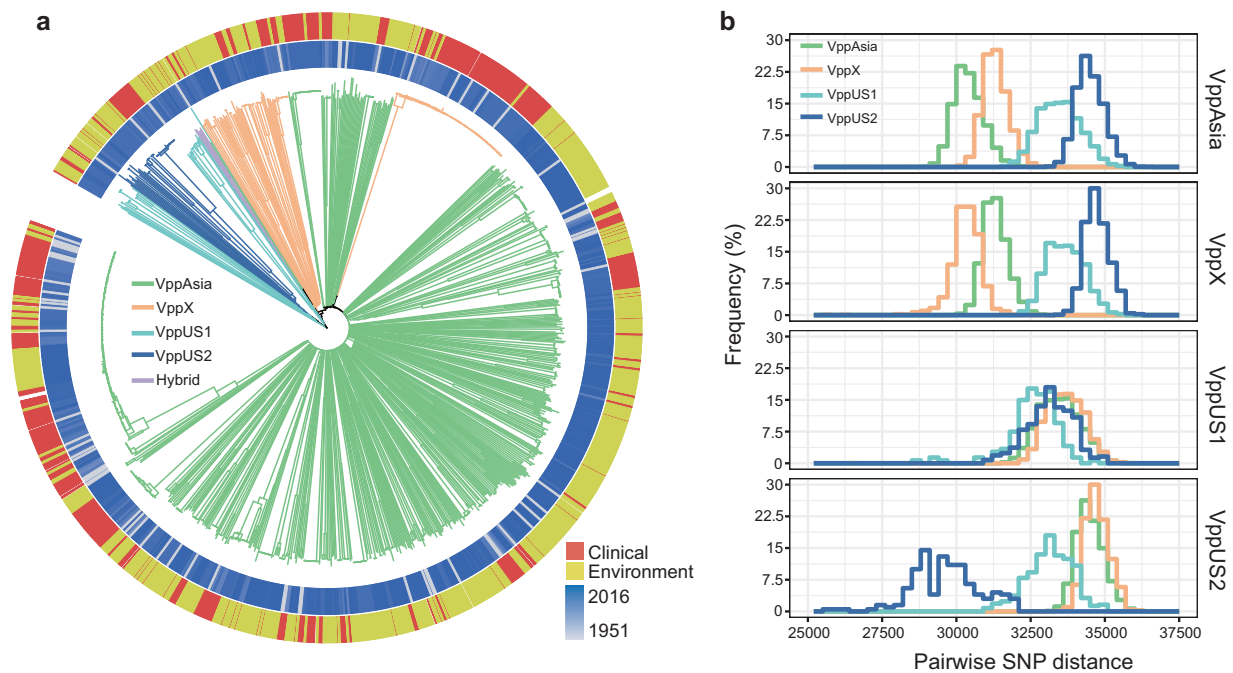


Fig. 1 Population structure of *V. parahaemolyticus* and relationships within and between populations. **a** NJ tree of 1103 *V. parahaemolyticus* strains based on 462,214 SNPs. Branch colors indicate populations defined by fineSTRUCTURE, green for VppAsia, orange for VppX, light blue for VppUS1, dark blue for VppUS2, purple for

hybrid strains. The ring colors from inner to outer indicate isolation time and sample type, respectively. The blank indicates information not available. **b** SNP distance within and between populations based on 469 non-redundant strains. Colors indicate populations and are consistent with branch colors of (a)

Overall, these results demonstrate the robustness of our iterative method to defining populations based on fineSTRUCTURE and confirm that as long as clonally related strains can be identified and removed, it represents a more appropriate approach for highly recombining species than phylogenetic methods.

The current distribution of the populations is shown in Fig. 2a. The great majority of isolates from Asia (574/600) are assigned to VppAsia, with all but one of the remainder (VppUS1, isolated from a shrimp farm in Thailand) being assigned to VppX. VppUS1 is found almost entirely in the US and is most common in the Mexican Gulf, with 13 out of the 29 VppUS1 strains are isolated from there. VppUS2 is most common on the US Atlantic coast (20 of 42) and has also been isolated several times in Northern Europe. VppX is most common on the Pacific coast and the Northern part of the US coast.

The geographic patterns found in our dataset are not predominantly determined by the spread of human disease clones or other clinical isolates, since similar patterns are observed if the dataset is restricted to the 634 environmental strains alone (Supplementary Fig. 5a), or indeed to other subsets of the data (Supplementary Fig. 5b, c). Furthermore, the distribution of CG1, the pandemic clonal group that mostly belongs to sequence type (ST) 3 [5], is similar to that of other VppAsia isolates, while CG2 (ST36), an epidemic group that is abundant in US and Canada [8], has a similar

distribution to that of VppX isolates, except that it has not been isolated from Asia (Supplementary Fig. 5d). Thus, human strains seem to have similar geographic patterns to other isolates (Supplementary Fig. 5c).

The 1103 genomes in this study have been collected for a variety of different purposes and do not represent a defined environmental or epidemiological cohort. Furthermore, sampling numbers in most locations are small and the coasts of Africa and Australia, for example are almost entirely unsampled. Nevertheless, our results demonstrate that at a global scale, geographic distributions of populations overlap considerably and that there is a substantial difference in the frequencies of the populations in the waters of Asia and those of the US Coast (Fig. 2a, Supplementary Fig. 5), for both environmental and clinical isolates.

Relationships amongst populations

The populations have a modest level of differentiation at the nucleotide level (Supplementary Table 2), with F_{st} values of around 0.1 approximately equivalent to that between humans living on different continents [19], implying that most common polymorphisms are shared between populations. VppUS1 is the most diverse and isolates are no more similar to each other in terms of mean SNP distance than they are to members of the other populations (Fig. 1b). However, according chromosome painting, which is based

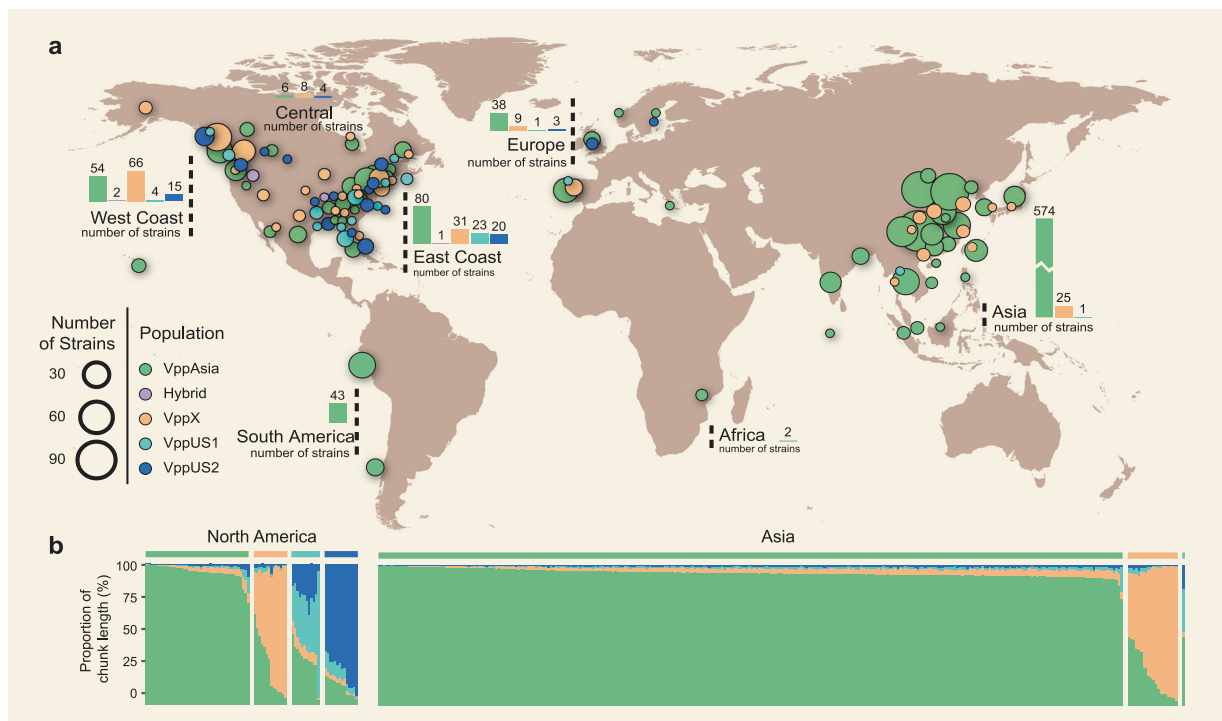


Fig. 2 Geographical distribution and admixture of *V. parahaemolyticus* populations. Colors in circle and bar plot indicate populations and are as in Fig. 1. Each circle indicates the population composition of a city/country, with radius in proportion to the sample size. Bar plot indicates the ancestry composition inferred by chromosome painting of two geographical regions: Asia and North

America. Each vertical bar represents one non-redundant strain and the proportion of color indicates the contribution of each population. Different populations are separated by blank vertical bar. Only strains with information of isolation location are included in (a) ($n = 1008$) and (b) ($n = 422$)

on haplotype similarity and therefore more sensitive in detecting sharing of DNA due to common descent, all of the members show substantially higher coancestry with other members of the population than any of the other isolates in the dataset (Fig. 2b), implying that the population consists of isolates that share ancestry, rather than being a collection of unassignable genomes. The other populations have consistently lower distances with members of their own populations and VppX and VppAsia are more closely related to each other than they are to VppUS1 and VppUS2.

One explanation for the high diversity of VppUS1 is that it has frequently absorbed genetic material from other populations. In order to test this hypothesis, while avoiding the effect of clonal relationships within the population itself on estimates of relationships with other populations, we painted the chromosomes of each of its members, using the members of the other three populations as donors. A high diversity of painting palettes was observed from VppUS1, with between 43 and 74% assigned to VppAsia and between 15 and 49% to VppUS2 (Supplementary Fig. 6a). By contrast, the other three populations showed lower levels of variation in assignment fractions in analogous paintings (Supplementary Fig. 6b-d). Thus, VppUS1 owes its high diversity to being a hub for admixture, with input from both VppUS2 and VppAsia. The members of VppUS1 in our

sample are all clearly distinct in ancestry profile from members of other populations (Fig. 2b), justifying the distinct population label, but if gene flow levels were higher, it seems likely that the population would lose its distinct identity and ancestry patterns would be better described by a continuum than discrete population labels.

We investigated differentiation between populations in core and accessory genes. Most of the variations show low level of differentiation, and the pattern is similar for core and accessory genes (Supplementary Fig. 7). There are relatively few highly differentiated loci, based on the F_{st} threshold of 0.8, we identified 192 core genes (825 SNPs) and 135 accessory genes (Supplementary Fig. 7, Supplementary Table 3). These genes are widely dispersed across many functions. Functional enrichment analysis highlighted elevated frequencies of the COG term of “amino acid transport and metabolism”, “transcription”, “inorganic ion transport and metabolism” and “secondary metabolites biosynthesis, transport, and catabolism” of differentiated accessory genes (Supplementary Fig. 8). It is also noteworthy that a type VI secretion system (VP1386-VP1414, T6SS1), which was proposed to be involved in the environment fitness of *V. parahaemolyticus* [20], is differentiated between groups, being present in 86% of VppX population and absent in VppUS2 and intermediate frequencies in the

other populations (Supplementary Fig. 9, Supplementary Table 3). Further, there is a four core gene cluster (VPA0679-VPA0682) which is highly differentiated between all pairs of populations, and thus can be considered as a “population indicator locus”. These genes include a *LysR* family transcriptional regulator, two arylsulfatase proteins and a hypothetical protein (Supplementary Table 3). These results together suggests that bacteria from the populations have evolved functional differences but further investigation will be required to elucidate their phenotypic or ecological consequences.

Recent mixing of *V. parahaemolyticus* populations

The observation of distinct populations is informative about patterns of migration in the past. Population genetic theory implies that differentiation between demes can only arise and persist if levels of migration between them are low, specifically on the order of magnitude of one migrant per generation or less [21]. The intuition behind the theory is that once a migrant arrives in a deme, it progressively imports DNA from other strains and becomes more and more similar to the other strains in its new deme. If too many strains are migrants, the demes will progressively lose their distinct genetic profiles and merge into a single gene pool. This theory has been developed for outbreeding eukaryotes [22] and bacterial populations deviate from several of the assumptions of the theory, in ways that are currently not well understood, making quantitative predictions impossible. Nevertheless, the qualitative expectation is that at equilibrium most isolates should have the ancestry profile of the region, with only a small fraction of the isolates having part ancestry from other locations.

The data differs from the qualitative predictions of migration-drift equilibrium because while there are few strains of clearly intermediate ancestry in the dataset, many locations have multiple strains from two or more of the four distinct populations that we have identified, making it not obvious what deme they belong to. Asia is clearly the most likely ancestral home range of VppAsia based on its high prevalence there but it is difficult to define boundaries of likely ancestral ranges for the other three populations with any confidence because the isolates assigned to those populations are too dispersed and they do not make up a clear majority anywhere. Thus the current distribution is qualitatively inconsistent with migration-drift equilibrium.

There are a number of factors which can in principle maintain subdivision when members of more than one population are found in the same location over long time periods. It is possible that the mechanism by which recombination occurs results in import occurring preferentially from members of the same population. For example, barriers to recombination due to homology

dependent mismatch repair has been proposed to account for the differentiation between phylogroups of *Escherichia coli*, despite high overall level of recombination [23] because the mechanism preferentially aborts recombination events between members of different phylogroup. Other mechanisms that can generate barriers to gene flow are lineage specific phage [24], restriction modification systems [25] or differences in an ecological niche [26].

The pattern of sharing of diversity is very different in *V. parahaemolyticus* to that found in *E. coli*, with high nucleotide diversity and low differentiation between them, so that barriers to recombination cannot be a simple function of recombination between strains from different populations being aborted due to low homology. None of the loci that are highly differentiated between populations are known to affect recombination, either through restriction modification or via modulating phage infection (Supplementary Fig. 7, Supplementary Table 3). Furthermore, to provide a general explanation for barriers to recombination, the mechanism would need to differentiate between all four populations and there is only one locus in the genome with that property, namely VPA0679-VPA0682. Another possibility is that functional differentiation might constrain recombination between populations. However, we have shown that more dramatic differentiation between ecogroups at selected loci has not resulted in any barrier to gene flow in the rest of the genome [26], making it highly unlikely that such a barrier between populations could maintain differentiation genome-wide.

We propose instead that the species is far from equilibrium because barriers to movement of strains have reduced recently. Under this hypothesis, it should be possible to approximately estimate the timescale on which mixing has taken place, based on the amount of introgression found in locations where the different populations now co-occur. Specifically, within our dataset, it is natural to compare the VppAsia isolates within Asia and in North America. Since Asia has been least affected by between continent migration (Fig. 2a), we predict that the VppAsia isolates in North America should have more ancestry from other sources, that they have acquired recently in their new locations. This prediction is borne out, a number of North America VppAsia isolates have high levels of VppUS1 and VppUS2 ancestry and on average the North America VppAsia isolates in the non-redundant set of 469 strains have 1.01% more (on average 2.97% in North America vs 1.96% in Asia) of their painting palette from VppUS sources than those from Asia (Fig. 2b and Fig. 3a).

In order to provide a timescale for the acquisition of non-Asian ancestry, we examined the evolution within the largest two clonal populations, CG1 and CG2. We removed recombination regions, then ran BEAST [27] to estimate a clock rate of 5.5×10^{-7} per site per year, with very similar

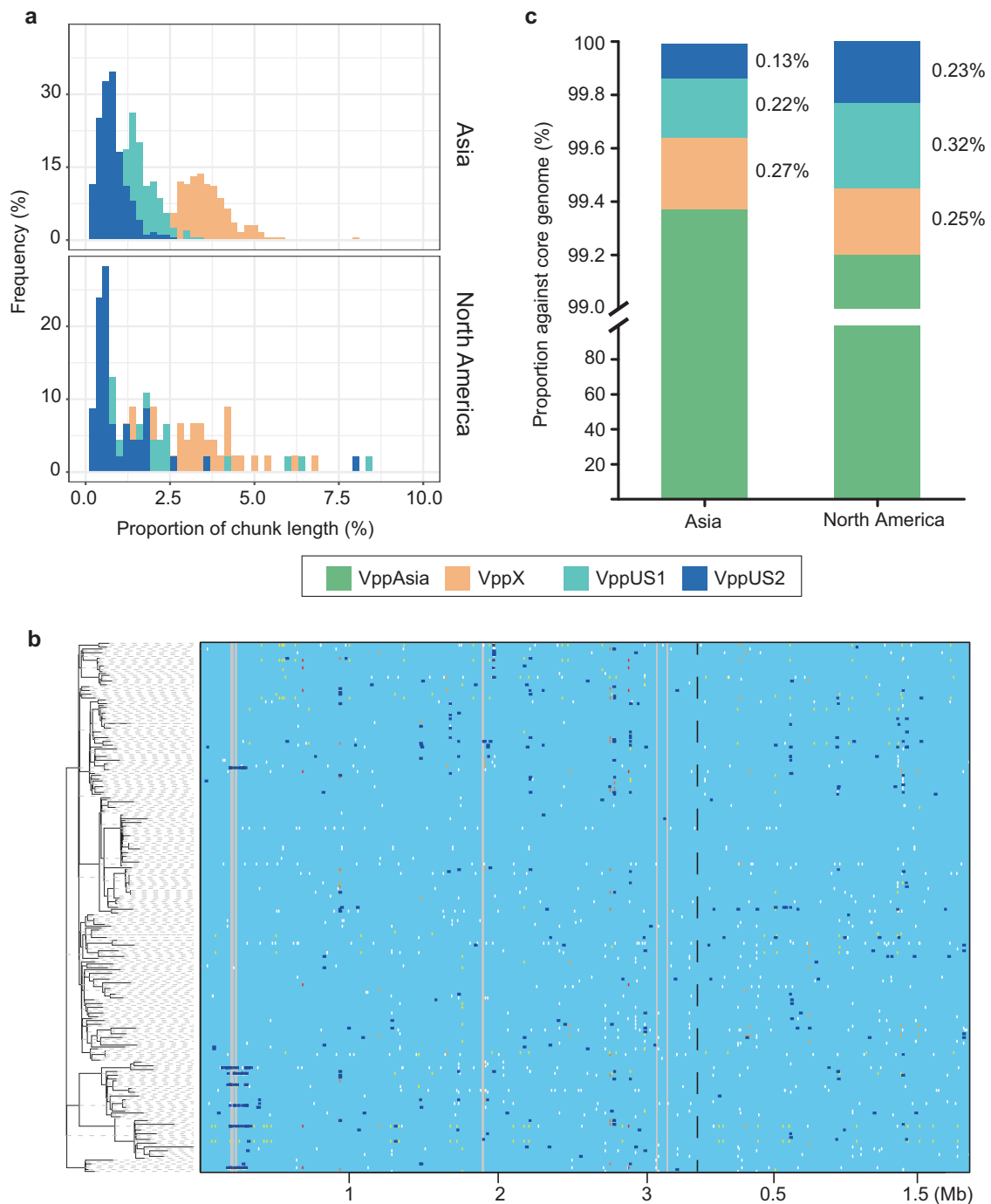


Fig. 3 Recent mixing of *V. parahaemolyticus* populations. **a** Ancestry composition of three other *V. parahaemolyticus* populations in VppAsia strains in different geographical regions. The contribution from other populations to the VppAsia is inferred by chromosome painting. X axis indicates the proportion of contributed chunk length of a population in one strain and Y axis indicates the corresponding frequency. **b** ClonalFrameML recombination analysis of 141 CG1 strains.

Left: ClonalFrameML reconstructed phylogeny. Right: dark blue horizontal bars indicate recombination events, grey areas indicate non-core regions. Two chromosomes are separated by dot line. **(c)** Source of recombination fragments of CG1 strains in different geographical regions. Y axis indicates the proportion of recombination fragments input from different population against core genome. Colors in **(a)** and **(c)** indicate four populations and are as in Fig. 1

values for the two clonal complexes (Supplementary Fig. 10). There are about 313 bases exchanged per mutation (Supplementary Fig. 11), so this implies a rate of recombination of 1.7×10^{-4} per site per year. Thus if all of the

import into the VppAsia bacteria was from US populations, then it would imply it would take about 59 years (with extreme lower and upper boundaries of 35–129 years) to acquire an extra 1.01% ancestry at this rate of import.

We also examined the origin of imports within CG1, the global pandemic clonal group. As for the VppAsia isolates, a higher fraction of the imports was from the two US populations amongst the isolates found in the North America than for the isolates found in Asia itself. This small difference in ancestry, corresponding to about 0.19% of the genome in total (Fig. 3c), has arisen during around 20 years since the beginning of global spread of CG1 in 1996.

These observations are consistent with a hypothesis that barriers to migration have become substantially weaker within the last few decades, but do not constitute direct evidence that patterns of gene flow between populations have changed. This hypothesis is empirically testable although we do not have a suitable strain collection to facilitate it. For example, if the Asian bacteria have arrived in large numbers in the US recently, then DNA from VppAsia bacteria should make up a higher proportion of recent genetic imports than older ones.

In order to explain why the pattern of dispersal has changed recently, it is necessary to first postulate reasons why dispersal was previously limited. We hypothesize that spread of bacteria between continents is limited by large distances between environments that are hospitable, making it rare that bacteria survive transportation between them. Large mammals, seabirds and other aquatic organisms travel large distances but do not necessarily provide habitats that *V. parahaemolyticus* can colonize for the days or weeks required to get from one continent to another. Thus, we propose that the necessary long-range dispersal did occur but was rare.

Humans have changed several aspects of the ocean environment, creating new habitats through effluent discharge, warming and acidifying the oceans through climate change, providing new mobile habitats on the hulls of ship and in ballast water and transporting copepods and other marine organisms deliberately to facilitate aquaculture or more accidentally through trade in marine products [28, 29]. Several of these could have facilitated transmission of bacteria between oceans. Furthermore, *V. parahaemolyticus* can adapt to colonize copepods [30] so that for example human-associated dispersal of species such as the manila clam from Asia to the America and Europe [31] could be responsible for the high frequency of Asian *V. parahaemolyticus* there. A single introduction via ballast water or introduction of shellfish for aquaculture would typically have low values of propagule pressure (a single event with few individuals), while recurring introductions through recently increased human activity may contribute in a regular basis introducing trans-ocean migration of *V. parahaemolyticus*, including for example by carriage by humans traveling for long distances, as for cholera [9].

Further work is required to narrow down the most important factors. To identify the frequency of *V. parahaemolyticus* reads in extensive metagenomic sampling of

the open ocean would provide knowledge on natural transmission of this bacterium. One objection to a human-mediated dispersal, rather than for example a role for climate change is that the absolute number of bacteria transported by ships or trade or humans themselves is likely to be small. However, this objection does not seem especially compelling. The absolute number of bacteria transported from one ocean to another does not need to be very large; if bacteria are fit in their new environment, they can multiply rapidly to constitute a substantial proportion of the bacteria in their new habitat.

Conclusions

Our results support our earlier conclusion that *V. parahaemolyticus* is subdivided into distinct geographical populations. We have identified 4 clearly differentiated populations, two of which appear to have foci in the US (VppUS1 and VppUS2). A third is predominant in Asia, while the ancestral home range of the fourth VppX is difficult to guess based on current sampling. However, these ranges pose a puzzle, in that they overlap substantially, both for environmental and human disease causing isolates, which show approximately similar patterns of distribution. Hybrids are rare, for example, amongst VppAsia isolates found in the US, most have ancestry profiles indistinguishable from strains found in Asia, while a handful have less than 10% introgression from either of the two US populations. The simplest and most parsimonious explanation is that previous barriers to migration have been reduced recently, allowing bacteria to disperse rapidly between continents but that because bacterial recombine relatively slowly (about 0.017% of their genome a year on average), there has not had sufficient time to generate hybrids.

These results have two major implications. Firstly, they suggest that recent human activity has disrupted long-standing barriers to genetic exchange in the oceans and that this has affected microbial population structure. Secondly, changing global patterns of *V. parahaemolyticus* disease incidence may be directly connected to changes in dispersal of the species, rather than being specific to the small number of clonal lineages that are responsible for most of the major outbreaks.

Materials and methods

Bacterial strains

In total 1103 strains were used in this research, including 392 newly sequenced and 711 publicly available strains (Supplementary Table 1). The newly sequenced strains

were isolated in China during daily food surveillance in 2014. The remaining 711 publicly available strains were downloaded from the NCBI database. The genomes of newly sequenced strains are available in GenBank with the accession numbers listed in Supplementary Table 1.

Sequencing and assembly

Strains were cultured in LB-2% NaCl agar at 37 °C, with genomic DNA extracted using a phenol/chloroform method. DNA was sequenced by using Illumina HiSeq 4000. A paired-end sequencing library with average insert size of 350 bp was built according to manufacturer's instructions (Illumina Inc., USA). The read length was set to 150 bp. On average, 500 Mb raw data were generated for each strain, corresponding to a sequencing depth of approximately 100 fold. The adaptor sequence and low quality reads were filtered and the clean reads were assembled by using SOAPdenovo v2.04 [32] as described previously [16]. The number of contigs and average size of assemblies are 263 and 5.1 Mb, respectively.

Variation detection

The SNPs were identified by aligning the *V. parahaemolyticus* genomes against reference genome (RIMD 2210633) by using MUMmer [33] as previously described [16], and only bi-allelic SNPs in the core genome (present in all isolates) were used in further analysis. As the size of the core genome and number of SNPs decreases as strains are added, we generated different core genome alignments for analyses of subsets of the data. In total 462,214 SNPs (2.4 Mb core genome) were identified from all 1103 genomes, 650,683 SNPs (3.3 Mb core genome) were from 469 non-redundant genomes, 355–8921 SNPs were identified within 13 clonal groups. To estimate the influence of different reference genomes to genome alignments, SNP identification and phylogeny construction, we tried two other reference genomes (FDA-R31 and 10329) for the analysis of 1103 genomes and got similar core genome size (2.38 and 2.39 Mb), SNP number (461,910 and 456,269) and phylogenetic trees. The distribution patterns of different populations in the trees were also similar (Supplementary Fig. 3).

Population structure

To illustrate the genetic distances between isolates, we constructed the Neighbor-Joining and Maximum-Likelihood trees using the TreeBest software (<http://treesoft.sourceforge.net/treebest.shtml>) and FastTree 2 [34] separately with default parameters based on sequences of concatenated

SNPs, and the trees were visualized by using online tool iTOL [35].

In order to define the population structure amongst the isolates while allowing for the effect of clonality, we used an iterative algorithm to successively remove strains with evidence of clonal relatedness. In the first step, we randomly removed isolates that differed by less than 2000 SNPs (the maximum SNP distance between isolates of pandemic clone CG1) from another, until a non-redundant genome set of 469 strains, each >2000 SNPs from the others in the dataset was obtained. We then ran fineSTRUCTURE [18] on this subset. The fineSTRUCTURE results included 115 populations with between 2 and 172 members as well as larger populations. The smaller populations consisted of strains with detectable clonal signal. We therefore randomly removed all but one of the strains and rerun fineSTRUCTURE. This procedure was continued iteratively a further 5 times until only populations with a single individual or more than 15 individuals remained, leading to a set of 260 genomes in total (Supplementary Fig. 2d).

We checked that these results are reproducible by repeating the same analysis but with two other independent randomizations. Equivalent population structure results were obtained in all three cases (Supplementary Fig. 2a-c). To investigate possible effect of sample size imbalance on clustering, we selected 60 strains, including 14–16 strains from each population and 2 hybrid strains, to repeat the fineSTRUCTURE analysis (Supplementary Fig. 2e). The results further verify the population structure of *V. parahaemolyticus* species but implies that VppUS1 population can be subdivided into subpopulations with different mixture profiles.

Population assignment based on fineSTRUCTURE was consistent with NJ tree (Fig. 1a, Supplementary Fig. 3) except for two strains, PCV08-7 and TUMSAT_H01_S4, and one epidemic group, CG2. Strain PCV08-7 and TUMSAT_H01_S4 were assigned to VppAsia and VppUS1 respectively by fineSTRUCTURE analysis, but in the NJ tree they are more closely related with VppX strains. The CG2 strains were all assigned to VppX populations by fineSTRUCTURE, but in NJ tree it was grouped with VppAsia strains. The length of chunks were extracted from the output file of Chromosome painting based on 469 non-redundant strains (Fig. 2b).

New designation of *V. parahaemolyticus* populations

In previous study, we designated four *V. parahaemolyticus* populations, named Asia-pop, US-pop 1, Hyb-pop 1, and Hyb-pop 2, separately, based on dataset of 157 genomes [16]. Here based on a larger sample, we found a new

population that was mostly isolated from US, and the previously defined Hyb-pop 1 were known as just several hybrid strains, or representatives of otherwise unsampled populations. As *V. parahaemolyticus* populations are geographical clustered, we propose a new nomenclature for them, VppAsia, VppX, VppUS1 and VppUS2. The ‘Vpp’ is abbreviation of ‘*V. parahaemolyticus* population’. The first three populations correspond to the previously defined Asia-pop, Hyb-pop 2 and US-pop 1, and VppUS2 is newly identified in this study.

Identification of pan-genome and restriction modification system genes

We re-annotated all the genomes of 1103 isolates using Prokka [36], and the annotated results were used in Roary [37] to identify the pan-genome and gene presence/absence. In total 41,052 pan-genes were identified, including 966 core genes (present in 1103 isolates) and 40,086 accessory genes.

To identify the restriction modification system (R-M system) genes of *V. parahaemolyticus*, BLASTP was used to scan the pan-genes of 1103 isolates for homologs (e -value $< 10^{-3}$, coverage > 0.5) against the standard protein sequences of R-M system genes retrieved from REBASE [38] (last accessed in 2013). 1132 genes was identified as possible R-M system genes, including 15 core genes and 1117 accessory genes. The pan-gene protein sequences were also used to BLAST (BLASTP) against the Clusters of Orthologous Groups of proteins (COG) database to get the COG annotations.

Inference of substitution rate using BEAST

Two clonal groups with large sample size, CG1 ($n = 153$, global pandemic group, also known as O3:K6 and its sero-variants group) and CG2 ($n = 92$, an epidemic group that popular in US, also known as serotype O4:K12), were selected to calculate molecular clock respectively by using BEAST v1.83 [27]. The variants that were attributed to recombination by our pipeline were excluded in substitution rates analysis. 10 out of 153 CG1 strains revealed an atypically large number of strain-specific SNPs and unusual long branches in the NJ tree, even after removing variants that could be attributed to recombination, and a similar pattern was observed in 1 of the 92 CG2 strains (Supplementary Fig. 12). These 11 strains with unusual high number of SNPs, and 22 strains with unknown isolation time, were excluded from BEAST analysis. We implemented analysis under GTR + Γ substitution model and relaxed clock model with constant size coalescent. The MCMC chain was run for 10^8 and sampling for every 5000

generations. The effective sample sizes of all inferred parameters were above than 200 in our results. The estimated molecular clock based on CG1 genomes is 5.6×10^{-7} with 95% confidence interval (CI) of $4.3\text{--}6.7 \times 10^{-7}$ per site per year, and 5.4×10^{-7} with 95% CI of $3.6\text{--}7.2 \times 10^{-7}$ for CG2 genomes. Here the average value, 5.5×10^{-7} , was used as the most likely estimate of *V. parahaemolyticus* molecular clock. The extremes of 95% CI based on two clonal groups, i.e., $3.6\text{--}7.2 \times 10^{-7}$, were used as lower and upper 95% bounds to be as conservative as possible.

Recombination detection and inference of recombination rate

In total, 13 clonal groups with more than 10 strains (Supplementary Table 1), defined by intra-group paired-distance less than 2000 SNPs, were selected to be used in detection of recombination events. We firstly used previously pipeline to detect recombination [16]. Briefly, we recalled the SNPs for each clonal group because different datasets had different core-genomes, and these SNPs were used to construct a NJ tree. Then PAML software package [39] was used to determine the SNPs of each branch. Assuming neutrality and no recombination, the observed SNP density of a given region should follow a binomial distribution. We used the sliding window method to identify regions that rejected the null hypothesis ($P < 0.05$) and all SNPs in such windows were treated as recombined SNPs. We also used ClonalFrameML [40], a software based on maximum likelihood method, to detect bacterial recombination within the same dataset. Sequence alignments of genomes for each clonal group and the corresponding maximum-likelihood tree constructed using PHYML with HKY model [41], were used as input files and non-core regions were ignored during calculation. The inferred recombination regions using ClonalFrameML are mostly consistent with our in-house method (Supplementary Fig. 11).

Two sets of r/μ estimates (ratio of size of recombination regions to the number of point mutations) were obtained through different methods, namely 331 with 95% CI of 257–405 for our in-house method and 295 with 95% CI of 217–373 for ClonalFrameML. The average value, 313, was selected as r/μ of the *V. parahaemolyticus*. Based on an estimate of the average mutation rate for the two clonal groups of 5.5×10^{-7} per site per year obtained using BEAST above, the most likely recombination rate of *V. parahaemolyticus* is 1.7×10^{-4} per site per year. We selected the extremes of 95% CI of r/μ and mutation rate (217–405, $3.6\text{--}7.2 \times 10^{-7}$, separately) to estimate confidence intervals of $7.8 \times 10^{-5}\text{--}2.9 \times 10^{-4}$ per site per year. Accordingly, the time to obtained 1.01% of genome

fragments would be 59 years with confidence limits of 35–129 years.

Ancestry inference for pandemic genome recombination events

We assigned CG1 strains into two groups according to their isolated location, with one isolated from Asia and another isolated from North America. By using Clonal-FrameML [40], we inferred the recombination fragments that occurred on each strain. Totally 81 fragments were found in Asia CG1 strains and 65 fragments were found in North America CG1 strains, with total size of 221 kb and the median length of 1035 bp. By aligning these recombination fragments against with 468 non-redundant genomes (excluding the CG1 genome from the dataset) using BLASTN, with a threshold of coverage $\geq 80\%$ and identity $\geq 99.5\%$, we identified possible donors genomes. For each recombination fragment, we calculated its observed donor frequency in each of the four populations, separately. And then for all recombination fragments carried by Asia CG1 genomes, we calculated their average value of observed donor frequency in a population, to employ as the contribution proportion from the corresponding population to the CG1 strains in Asia. Similarly, we obtained the contribution proportion of each population to the North America CG1 genomes (Fig. 3c). We also try the relaxed identity thresholds (99.0%) in BLASTN and acquired similar results.

Acknowledgements We gratefully acknowledge Dr. Narjol Gonzalez-Escalona for contributing *V. parahaemolyticus* genomes. This work is supported by the National Key Research & Development Program of China (No. 2017YFC1601503, 2016YFC1200100 and 2017YFC1200800), Sanming Project of Medicine in Shenzhen (No. SZSM201811071) and the National Natural Science Foundation of China (No. 31770001). J. Martinez-Urtaza is funded by Natural Environment Research Council (NERC) project (No. NE/P004121/1). D.F. is funded by a Medical Research Council Fellowship as part of the MRC CLIMB consortium for microbial bioinformatics (grant number MR/M501608/1).

Author's contributions YC, DF, and RY designed the study and coordinated the project; XP, LY, JM, QH, and DY contributed strains for analysis; CY, XP, YW, NC, YQS, YJS, YY, MJ, CQ, DF, and YC analyzed the data; EF, JM, and DZ provided insightful comments, DF and YC wrote the manuscript. All authors approved the final version of the manuscript.

Compliance with ethical standards

Conflict of interest The authors declare that they have no conflict of interest.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

References

- Brown MV, Ostrowski M, Grzymalski JJ, Lauro FM. A trait based perspective on the biogeography of common and abundant marine bacterioplankton clades. *Mar Genomics*. 2014;15:17–28.
- Yilmaz P, Yarza P, Rapp JZ, Glockner FO. Expanding the world of marine bacterial and archaeal clades. *Front Microbiol*. 2015;6:1524.
- Kent AG, Dupont CL, Yooseph S, Martiny AC. Global biogeography of *Prochlorococcus* genome diversity in the surface ocean. *ISME J*. 2016;10:1856–65.
- Hellweger FL, van Sebille E, Calfee BC, Chandler JW, Zinser ER, Swan BK, et al. The role of ocean currents in the temperature selection of plankton: insights from an individual-based model. *PLoS ONE*. 2016;11:e0167010.
- Nair GB, Ramamurthy T, Bhattacharya SK, Dutta B, Takeda Y, Sack DA. Global dissemination of *Vibrio parahaemolyticus* serotype O3:K6 and its serovariants. *Clin Microbiol Rev*. 2007;20:39–48.
- Mutreja A, Kim DW, Thomson NR, Connor TR, Lee JH, Kariuki S, et al. Evidence for several waves of global transmission in the seventh cholera pandemic. *Nature*. 2011;477:462–5.
- Weill FX, Domman D, Njamkepo E, Tarr C, Rauzier J, Fawal N, et al. Genomic history of the seventh pandemic of cholera in Africa. *Science*. 2017;358:785–9.
- Martinez-Urtaza J, van Aarle R, Abanto M, Haendiges J, Myers RA, Trinanes J, et al. Genomic variation and evolution of *Vibrio parahaemolyticus* ST36 over the course of a transcontinental epidemic expansion. *mBio*. 2017;8:pii: e01425–17.
- Chin CS, Sorenson J, Harris JB, Robins WP, Charles RC, Jean-Charles RR, et al. The origin of the Haitian cholera outbreak strain. *N Eng J Med*. 2011;364:33–42.
- Yeung PS, Boor KJ. Epidemiology, pathogenesis, and prevention of foodborne *Vibrio parahaemolyticus* infections. *Foodborne Pathog Dis*. 2004;1:74–88.
- Su YC, Liu C. *Vibrio parahaemolyticus*: a concern of seafood safety. *Food Microbiol*. 2007;24:549–58.
- Ansedo-Bermejo J, Gavilan RG, Trinanes J, Espejo RT, Martinez-Urtaza J. Origins and colonization history of pandemic *Vibrio parahaemolyticus* in South America. *Mol Ecol*. 2010;19:3924–37.
- Martinez-Urtaza J, Trinanes J, Gonzalez-Escalona N, Baker-Austin C. Is El Niño a long-distance corridor for waterborne disease? *Nat Microbiol*. 2016;1:16018.
- Baker-Austin C, Trinanes J, Gonzalez-Escalona N, Martinez-Urtaza J. Non-cholera vibrios: the microbial barometer of climate change. *Trend Microbiol*. 2017;25:76–84.
- Yan Y, Cui Y, Han H, Xiao X, Wong HC, Tan Y, et al. Extended MLST-based population genetics and phylogeny of *Vibrio parahaemolyticus* with high levels of recombination. *Int J Food Microbiol*. 2011;145:106–12.
- Cui Y, Yang X, Didelot X, Guo C, Li D, Yan Y, et al. Epidemic clones, oceanic gene pools, and eco-LD in the free living marine pathogen *Vibrio parahaemolyticus*. *Mol Biol Evol*. 2015;32:1396–410.
- Mihaescu R, Levy D, Pachter L. Why neighbor-joining works. *Algorithmica*. 2009;54:1–24.
- Lawson DJ, Hellenthal G, Myers S, Falush D. Inference of population structure using dense haplotype data. *PLoS Genet*. 2012;8:e1002453.
- Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, Zhivotovsky LA, et al. Genetic structure of human populations. *Science*. 2002;298:2381–5.
- Salomon D, Gonzalez H, Updegraff BL, Orth K. *Vibrio parahaemolyticus* type VI secretion system 1 is activated in marine

- conditions to target bacteria, and is differentially regulated from system 2. *PLoS one*. 2013;8:e61086.
21. Wright S. Evolution in Mendelian populations. *Genetics*. 1931;16:97–159.
 22. Whitlock MC, McCauley DE. Indirect measures of gene flow and migration: F_{ST} not equal to $1/(4Nm+1)$. *Heredity*. 1999;82:117–25.
 23. Didelot X, Meric G, Falush D, Darling AE. Impact of homologous and non-homologous recombination in the genomic evolution of *Escherichia coli*. *BMC Genom*. 2012;13:256.
 24. Rodriguez-Valera F, Martin-Cuadrado A-B, Rodriguez-Brito B, Pašić L, Thingstad TF, Rohwer F, et al. Explaining microbial population genomics through phage predation. *Nat Rev Microbiol*. 2009;7:828.
 25. Jeltsch A. Maintenance of species identity and controlling speciation of bacteria: a new function for restriction/modification systems? *Gene*. 2003;317:13–16.
 26. Cui Y, Yang C, Qiu H, Wang H, Yang R, Falush D. The landscape of coadaptation in *Vibrio parahaemolyticus*. *bioRxiv*. 2018; <https://doi.org/10.1101/373936>.
 27. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol*. 2012;29:1969–73.
 28. Ruiz GM, Rawlings TK, Dobbs FC, Drake LA, Mullady T, Huq A, et al. Global spread of microorganisms by ships. *Nature*. 2000;408:49–50.
 29. Martinez-Urtaza J, Baker-Austin C, Jones JL, Newton AE, Gonzalez-Aviles GD, DePaola A. Spread of Pacific Northwest *Vibrio parahaemolyticus* strain. *N Engl J Med*. 2013;369:1573–4.
 30. Martinez-Urtaza J, Blanco-Abad V, Rodriguez-Castro A, Ansedo-Bermejo J, Miranda A, Rodriguez-Alvarez MX. Ecological determinants of the occurrence and dynamics of *Vibrio parahaemolyticus* in offshore areas. *ISME J*. 2012;6:994–1006.
 31. Chiesa S, Lucentini L, Freitas R, Nonnis Marzano F, Breda S, Figueira E, et al. A history of invasion: COI phylogeny of Manila clam *Ruditapes philippinarum* in Europe. *Fish Res*. 2017;186:25–35.
 32. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. *GigaScience*. 2012;1:18.
 33. Delcher AL, Salzberg SL, Phillippy AM. Using MUMmer to identify similar regions in large sequence sets. *Curr Protocol Bioinform*. 2003;Chapter 10:Unit 10 3.
 34. Price MN, Dehal PS, Arkin AP. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS ONE*. 2010;5:e9490.
 35. Letunic I, Bork P. Interactive tree of life (iTOL)v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res*. 2016;44:W242–5.
 36. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*. 2014;30:2068–9.
 37. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, et al. Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics*. 2015;31:3691–3.
 38. Roberts RJ, Vincze T, Posfai J, Macelisi D. REBASE—a database for DNA restriction and modification: enzymes, genes and genomes. *Nucleic Acids Res*. 2014;43:D298–9.
 39. Yang Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol*. 2007;24:1586–91.
 40. Didelot X, Wilson DJ. ClonalFrameML: efficient inference of recombination in whole bacterial genomes. *PLoS Comput Biol*. 2015;11:e1004041.
 41. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst Biol*. 2003;52:696–704.