1    **Title:** Population genomics of Vibrionaceae isolated from an endangered oasis

2            reveals local adaptation after an environmental perturbation.

3

4    **Authors:** Mirna Vazquez-Rosas-Landa[1,2], Gabriel Yaxal Ponce-Soto[1], Jonás A.

5    Aguirre-Liguori[1], Shalabh Thakur[3], Enrique Scheinvar[1], Josué Barrera-Redondo[1],

6    Enrique Ibarra-Laclette[2], David S. Guttman[3,4], Luis E. Eguiarte[1] and Valeria

7    Souza[1*]

8

9    **Addresses:** [1]Departamento de Ecología Evolutiva, Instituto de Ecología,

10   Universidad Nacional Autónoma de México, Ciudad Universitaria, 04510, Ciudad

11   de México, México.

12   [2]Red de Estudios Moleculares Avanzados, Instituto de Ecología, A.C. – INECOL,

13   Clúster Científico y Tecnológico BioMimic®, Carretera antigua a Coatepec 351, El

14   Haya, 91070 Xalapa, Veracruz, México.

15   [3]Department of Cell and Systems Biology, University of Toronto, Toronto, Ontario,

16   Canada.

17   [4]Centre for the Analysis of Genome Evolution and Function, University of Toronto,

18   Toronto, Ontario, Canada.

19

20   * Correspondence: Valeria Souza, Instituto de Ecología, UNAM, souza@unam.mx,

21   souza.valeria2@gmail.com.

22

23

24

25

26

**Abstract (323 words).**

**Background:** In bacteria, pan-genomes are the result of the evolutionary "tug of war" between selection and horizontal gene transfer (HGT). High rates of HGT increase the genetic pool and the effective population size, resulting in open pan-genomes. In contrast, selective pressures can lead to local adaptation by purging the variation introduced by HGT, resulting in closed pan-genomes and clonal lineages. In this study, we explored both hypotheses elucidating the pan-genome of Vibrionaceae isolates after a perturbation event in the endangered oasis of Cuatro Ciénegas Basin (CCB), Mexico, and looking for signals of adaptation to the environments in their genomes.

**Results:** We obtained 42 genomes of Vibrionaceae distributed in six lineages, two of them did not showed any close reference strain in databases. Five of the lineages showed closed pan-genomes and were associated to either water or sediment environment; their high effective population size ($N_e$) estimates suggest that these lineages are not from a recent origin. The only clade with an open pan-genome was found in both environments and was formed by ten genetic groups with low $N_e$, suggesting a recent origin. The recombination and mutation estimators (r/m) ranged from 0.0052 to 2.7249, which are similar to oceanic Vibrionaceae estimations; however, we identified 367 gene families with signals of positive selection, most of them found in the core genome; suggesting that despite recombination, natural selection moves the Vibrionaceae CCB lineages to local adaptation purging the genomes and keeping closed pan-genome patterns. Moreover, we identify 598 SNPs associated with an unstructured environment; some of the genes under this SNPs were related to sodium transport.

**Conclusions:** Different lines of evidence suggest that the sampled Vibrionaceae, are part of the rare biosphere usually living under famine conditions. Two of these lineages were reported by the first time. Most Vibrionaceae lineages of CCB are adapted to their microhabitats rather than to the sampled environments. This pattern of adaptation agrees with the association of closed pan-genomes and local adaptation.

58    **Keywords:** (4 to 6)

59    Pan-genome, genomics, Vibrionaceae, recombination, selection, effective

60    population size.

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

## Background

78     Comparative genomics analyses have shown a wide range of genomic
79 variation within bacteria from different phylogenetic groups [1-3]. This variation
80 range has been explained in part by the wide ecological niche occupied by different
81 bacterial groups [4-8]. Bacterial genomes, in contrast to eukaryotic genomes,
82 usually maintain constant genome sizes [9, 10], suggesting that while horizontal
83 gene transfer (HGT) increases the genome size by adding new genes, selection
84 maintains the genome size by removing deleterious, non-functional or non-useful
85 genes [11-13]. Therefore, bacteria can present very different genomic
86 compositions even within a species, with HGT creating a flexible genome and
87 natural selection purging or maintaining it [10, 14].

88     Thus, the type of pan-genome is an indication of the evolutionary "tug of
89 war" between selection and HGT. As a prediction, if there are high rates of HGT,
90 the total genetic pool will increase, as well as the effective population size,
91 generating an open pan-genome maintained by natural selection [15]. However, if
92 there is a selective pressure towards local adaptation, the genetic diversity
93 introduced by HGT will be purged, resulting in a closed pan-genome and clonal
94 lineages [14].

95     To start understanding the reasons why some pan-genomes are open while
96 others are closed, we can analyze the rate and type of recombination. On the one
97 hand, homologous recombination homogenizes populations, keeping them
98 genetically cohesive in a closed pan-genome [16, 17]. On the other hand, non-
99 homologous recombination brings new genetic material, offering new evolutionary
100 opportunities for diversification and generating an open pan-genome [18-21].
101 Recombination decreases linkage disequilibrium among genes, allowing selection
102 and the related Hill-Robertson effect to operate in specific genes and avoiding the
103 purged of genetic diversity along with the genome [22, 23]. As a result of this
104 diversity, species with higher recombination levels maintain a large historical
105 effective population size [15, 24, 25]. In contrast, highly clonal populations with low
106 or no HGT evolve mostly by mutation and genetic drift, because the efficiency of

107 selection is hampered by the Hill-Robertson effect that also reduces the standing
108 levels of variation in the population and the historical effective population sizes [23,
109 26].

110     In this study, we explored the probable role of different evolutionary forces
111 shaping the genetic diversity of Vibrionaceae in the oasis of the Cuatro Ciénegas
112 Basin (CCB), Mexico. CCB is composed of several aquatic systems that have a
113 significant unbalance of the nutrient stoichiometry [27]. Population genetic studies
114 of *Pseudomonas* spp., *Exiguobacterium* spp. and *Bacillus* spp. isolated from CCB
115 ponds and rivers in general have shown low recombination levels [28-30]. These
116 patterns suggest that nutrient constraints in CCB may work as an ecological filter,
117 reducing recombination maybe due to the cost of replicating new DNA, and leading
118 to local adaptation [27, 31, 32].

119     We tested whether the environmental nutrient constraint would affect the
120 genetic structure of *Vibrio* spp. lineages at CCB. Members of *Vibrio* spp. has been
121 characterized in general as a highly recombinant [33, 34]. We analyzed the genetic
122 structure of Vibrionaceae in a particular site of CCB, Pozas Rojas (Figure 1). This
123 site was the most stoichiometrically unbalanced (N:P 156:1) in our first sampling on
124 2008 [35]. Later, Pozas Rojas was naturally perturbed with intense rains
125 associated with hurricane Alex in 2010. The runoff detritus and water, caused the
126 nutrients ratios to change from extremely unbalanced stoichiometry to a ratio
127 similar the standard values in the sea (N:P 20:1; compared to the Redfield
128 standard N:P 16:1 values of the sea [36]). Given the change in stoichiometry ratios,
129 we asked the following questions: 1) How did a naturally recombinant lineage like
130 some members of Vibrionaceae respond to this perturbation? 2) Did Vibrionaceae
131 lineages maintained their local adaptation to this unique site by restricting
132 recombination, and maintaining their pan-genomes closed? Alternatively, 3) Is it
133 possible that *Vibrio* spp. developed open pan-genomes with large effective
134 population sizes, similar to the lineages in the ocean to deal with this stoichiometric
135 change? [33, 34].

136  Herein we analyzed the role of the evolutionary forces that have shaped

137  Vibrionaceae at CCB by performing a comparative genomics analysis of five

138  reference and 42 strains isolated from two different local environments (i.e., water

139  and sediments) in perturbed Pozas Rojas. Contrary to what we expected, our

140  results show that most CCB Vibrionaceae lineages had similar levels of

141  recombination compared to their oceanic relatives, and much higher levels of

142  recombination than other genera in the CCB [28-30]. However, since most of the

143  analyzed lineages had closed pan-genomes, we suggest that most of such

144  recombination is homologous. This type of recombination should promote

145  reproductive isolation and generate local adaptation. We did not observe a clear

146  pattern of adaptation to either water or sediment environments, suggesting that

147  there may be other environmental variables that we were not able to measure that

148  could be driving local adaptation among these lineages.

149  **Results.**

150  **Nutrients concentrations.** Based on Kruskal-Wallis statistical test, the total

151  nutrient concentrations (Carbon (C), Nitrogen (N), and Phosphorus (P)) of the

152  Pozas Rojas were not significantly different between sample points,(C: $p= 0.8815$;

153  N: $p= 0.2256$ and P: $p= 0.9624$; Additional file 1: Table 1) however, they were

154  statistically significant between type of environment (water vs. sediment) (C: $p=$

155  $3.486e-4$; N: $p= 0.03798$ and P: $p= 3.461e-4$).

156  The proportion of C:N:P was on average 350:9:1 for water, and 258:21:1 for

157  sediment (Additional file 1: Table 2). This ratios indicate a stoichiometric "balance"

158  (i.e., similar to Redfield standard ratios) in Pozas Rojas during 2013, due to higher

159  P availability, compared with the extreme stoichiometric imbalance observed in

160  most of CCB sites, and in particular in Pozas Rojas microbial mat during summer

161  2008 (i.e., 15,820:157:1)[35], previous to the hurricane Alex perturbation.

162  **Phylogenetic Diversity and the environmental association.** The phylogenetic

163  relationships of 16S rRNA gene (700 bases) of the 174 cultivated isolates from

164  Pozas Rojas, showed that the strain collection was dominated by Vibrionaceae

165  (63%), followed by Aeromonadaceae (14%) and Halomonadaceae (9.7%;

6

166   Additional file 1: Table 3). Among Vibrionaceae, we identified two different genera;

167   most strains belong to *Vibrio* spp. (93.6%) and far less to the related

168   *Photobacterium* spp*.* genus (6.3%).

169        The aligned sequences were used to construct a maximum likelihood tree

170   with PhyML (Additional file 1: Figure 1). Based on the previous taxonomic

171   assignment, the sequences of *V. alginolyticus, V. parahaemolyticus, V.*

172   *anguillarum, V. metschnikovii*, and *Photobacterium* spp*.* were included as

173   references. This analysis reveals seven different cultivated Vibrionaceae lineages

174   in Pozas Rojas.

175        In order to characterize the relationship between water/sediment

176   environments and Vibrionaceae lineages, we performed an AdaptML analysis [37].

177   The analysis showed that strains are structured according to the environment

178   where they were isolated, i.e., water or sediment, and not by pond (Additional file

179   1: Figure 2). While most clades were specialist either to water (higher nutrient

180   condition) or to sediment (lower nutrient condition), the most abundant lineage had

181   no preference for any environment. Based on the AdaptML analysis, we selected

182   42 isolates for further sequencing; these isolates were chosen as representatives

183   from the different lineages and environments.

184   **Genome features.** Among the 39 CCB sequenced *Vibrio* spp. genomes, we found

185   variation in terms of genome size, ranging from 3.1 Mbp to 5.1 Mbp, while the three

186   CCB *Photobacterium* spp. genomes had an average genome size of 4.5 Mbp.

187   Despite this variation, when we compared the CCB strains genomes to their

188   closest reference strain, we found similar genome sizes (Additional file 1: Table 4).

189   Moreover, for each of the assembled genomes, we evaluated their completeness

190   with BUSCO [38]. We found that 92.8% of the genomes contained more than 95%

191   of the 452 near-universal single-copy orthologs evaluated by the program

192   (Additional file 1: Table 5), suggesting that the observed variation in genome sizes

193   could be due to intrinsic characteristics of each strain and not to a sequencing bias.

194   **Pan-genome analyses of CCB Vibrionaceae and lineages description.** The

195   pan-genome analysis of 39 CCB *Vibrio* spp., 3 CCB *Photobacterium* spp., and 5

7

196     *Vibrio* spp. references strains involved a total of 20,121 orthologous gene families.

197     The genes that were present in at least 95% of the genomes conformed the core

198     genome, including reference genomes, composed by 1,254 gene families. The

199     accessory genome is far more substantial, consisting of a total of 14,072 genes

200     families that were found in at least two of the obtained genomes. The rest 4,795

201     genes families were strain-specific.

202     Using the core genes families, we reconstructed a lineage phylogeny

203     (Figure 2). In the core phylogeny we found seven lineages, of which six of them

204     were previously identified in the 16s rRNA gene tree, and one was represented by

205     a unique strain of marine *V. furnissii* sp. Nov. 4 stran (NCTC 11218) [39].

206     Reference strain *V. anguillarum* 775, isolated from a Coho salmon [40] clusters

207     within the large generalist Clade II, while reference strain *V. metschnikovii* CP 69-

208     14, which was isolated in marine systems, is basal to Clade III. Basal to Clade VI

209     are reference *V. parahaemoliticus* BB22OP, a pre-pandemic strain [41] associated

210     with seafood-borne gastroenteritis in humans and *V. alginolyticus* NBRC 15630 =

211     ATCC 17749, an aquatic organism that can cause bacteremia. Clades IV and V

212     are likely to be exclusive to CCB, given that there is no closely related strain

213     sequenced on databases. Finally, Clade I is related to *Photobacterium* spp. (Figure

214     2).

215     From the six clades identified, only Clade II presented an open pan-genome

216     as suggested by the Heaps law analysis [42] (alpha= 0.7913). The rest of the

217     clades displayed closed pan-genome patterns (i.e., alpha values >1.0; Table 1).

218     We performed random sub-samplings of genomes per clade to verify the effect of

219     sample size and we re-calculate alpha values of each clade with the minimum

220     sample size, and in all cases, we found the same results of closed pan-genomes

221     for the specialist clades and an open pan-genome for generalist clade (Additional

222     file 1: Figure 3).

223     **Genetic diversity and recombination estimates.** General estimators of genetic

224     diversity were obtained for each clade and Sub-clade (Table 2). We found that

225     nucleotide diversity values for Clades III, IV, and V were the lowest within sample,

226    ranging from 2.86E-05 to 0.0051, while Clades I, II, and VI had higher levels of

227    genetic variation, in the range of 0.011 to 0.046. When estimating the nucleotide

228    diversity for Sub-clades belonging to Clade II (described below, see Additional file

229    1: Figure 4), we found lower values in the range of 1.61E-06 to 5.47E-06. This

230    same pattern was observed for the θw values (Table 2). Due to the number of

231    individuals we could not obtain Tajima's *D* estimator for Clades I and VI. For the

232    rest of the clades Tajima's *D* values were negative, except for Clade II that had

233    positive values.

234        Since most lineages present a closed pan-genome, we tracked the

235    footprints of recombination by using two different approaches. The first approach

236    consisted of assessing the recombination in each ortholog group. The second

237    involved the identification of recombination signals based on a whole genome

238    alignment. With the first approach, we found that from the 15,380 ortholog clusters

239    analyzed, only the 11% (1,759) showed significant signal of recombination

240    (Additional file 2: Table 6). These recombination events occurred more frequently

241    among isolates of the same environment and pond, suggesting reproductive

242    isolation associated to an environmental variable (Figure 3). However and despite

243    we considered in our calculations the pan-genome size, number of strains per

244    clade and branch length, it is also true that most clades are conformed by only

245    isolates of water or sediment. Therefore, we propose that the frequency of

246    recombination events is mostly restricted to occur within clades (Figure 3;

247    Additional file 1: Figure 5).

248        In the case of the generalist Clade II, we found sub-structure. Using Nei's

249    genetic distances, we identified ten genetic groups (that we will call Sub-clades

250    therefore) with distances greater than 0.001. The discriminant function shows the

251    same structure as the Nei distances, reflecting a broader relationship between

252    Sub-clades A, D, F and G and B with C and E. Meanwhile, H, I and J Sub-clades

253    had dissimilar sub-structures (Additional file 1: Figure 4). Since only three of the

254    Sub-clades contained more than two isolates, further analyses were just performed

255    with the larger Sub-clades (A, D and G).

256    Following the second approach, we evaluated the impact of homologous

257    recombination and mutation within lineages estimating *r/m* using the clonal frame

258    software [43]. This measure reflects the ratio of probabilities that a given

259    polymorphism is explained by either recombination (*r*) or by mutation (*m*).  Clade VI

260    displayed the lowest *r/m* values=0.0052, while Clade I (i.e., *Photobacterium* spp.)

261    had the highest value in our dataset, *r/m* = 2.72 (Table 3). We also performed the

262    same analysis on *V. parahaemolyticus, V. ordalii, V. anguillarum*, and *P. leiognathi*

263    reference genomes, all isolated from marine environments. For the marine

264    samples, *r/m* estimates were within the range of CCB strains, except *V.*

265    *anguillarum*, which had the highest values (Table 3). This analysis also shows that

266    some recombination events are shared with *Vibrio* spp. references strains

267    (Additional file 1: Figure 6) supporting the hypothesis of ancient origin of these

268    recombination events even though more recent recombination events were

269    detected only among CCB strains. This indicates that homologous recombination is

270    a constant source (albeit relatively infrequent) of polymorphism in the analyzed

271    strains.

272    **Estimates of effective population sizes.**

273    Using a simulation approach with the Fastsimcoal2 program [44] we estimated the

274    posterior distribution of the effective population size ($N_e$) of each of the six clades.

275    We found large population sizes (Table 4) ranging from millions in the specialist

276    Clades I ($N_e$ = 12,822,270), III ($N_e$ = 15,018,880), V ($N_e$ = 9,594,874) to

277    intermediate in the range of thousands in the Clades IV ($N_e$ = 383,067) and VI *($N_e$*

278    =141,870), and to far smaller in the Sub-clades of the locality common generalist

279    Clade II ( i.e., Sub-clade A  $N_e$ = 55,938; Sub-clade D $N_e$ = 20,849;  Sub-clade G

280    $N_e$ = 29,791) reinforcing the idea of recent diversification in these Sub-clades.

281    **Selection analyses**

282    FUBAR uses a codon-based model of evolution that allows the identification

283    of evolving sites under positive or purifying selection in protein-coding genes

284    through a Markov chain Monte Carlo (MCMC) routine. From a total of 15,380

285    ortholog clusters analyzed, only 367 (2.3 %) had a significant signal of positive

286  selection according to FUBAR results. Of these ortholog gene families, 297

287  belonged to the flexible genome, while 70 are part of the core genome. However,

288  when we considered the universe of ortholog genes that conform the flexible

289  genome (14,072), only 2.1% of the flexible genome had signals of positive

290  selection, while in the core genome, composed by 1,254 genes, 5.6% of the genes

291  are positive selected (Additional file 2:  Table 7). A GO enriched analysis was

292  performed in order to identify those biological functions overrepresented given

293  those ortholog clusters with positive selection. Seven Gene Ontology (GO) terms

294  were enriched within these families. (Table 5). Moreover, based on a whole

295  genome alignment, we obtained 38,533 SNPs variants, from which 26,663 were bi-

296  allelic characters that were used in an UPGMA analysis of genetic distances. This

297  analysis produced the same clustering as the core genome phylogeny (Figure 2).

298  As well, with this SNPs we performed a membership probability test, which show

299  that all the isolated had the same probability of been isolated from any pond and

300  environment (Additional file 1: Figure 7).

301  We found on average 2,473 private (unique) SNPs for each one of the nine

302  ponds, 33,655 private SNPs for water or sediment environments, and 29,141,

303  private SNPs for each of the six clades. This abundance of private SNPs suggests

304  an effect of the environment, either by local adaptation (selection) or by genetic

305  drift (low effective sizes or little or no gene flow).

306  We removed the SNPs with a minor allele frequency < 0.05 (771 SNPs

307  removed) and we kept the alleles that were found in at least three individuals, for a

308  total of 25,892 SNPs. Within those SNPs we detected a total of 598 SNPs with an

309  association to the sediment environment. A UPGMA analysis of these 598 SNPs

310  was performed in order to infer the similarity between samples (Figure 4) finding

311  most of the clusters previously observed with the core genome phylogeny (Figure

312  2), except Clade III, which appears inside Clade II. Moreover, the mixed isolates of

313  Clade III fall among the Sub-clade G of Clade II, most of them were isolated from

314  water environment, as well as members of Clade III (Figure 4), suggesting a

315  preference for diluted, unstructured environments.

11

316       To analyze the distribution of the SNPs, we mapped the above detected 598

317    SNPs to their positions in the genome alignment from where they were obtained,

318    moving in 1 Kb windows. A total of 144 genomic regions containing SNPs were

319    inspected, and we found 237 ortholog gene families in these regions. From these

320    ortholog gene families, only 24 showed recombination signals, while 18 had

321    selection signals (Additional file 2: Table 8). Within those SNPs we performed a

322    test for GO-term enrichment with TopGO [45]. From the 24 ortholog genes families

323    with recombination signals, we detected four enriched GO, while we found only

324    one enriched GO-term in the 18 ortholog gene families with selection signals

325    (Table 6).

326       Besides those analyses, based on pan-genome information, we looked for

327    specific coding sequences that could be private (unique) to a specific pond,

328    environment, or clade. There were no specific genes associated with a particular

329    environment or pond, but we did identify ortholog gene clusters exclusive per

330    clade. From Clades I to VI, we observed 1280, 10, 72, 23, 72, and no exclusive

331    ortholog gene families, respectively. For each clade with exclusive ortholog gene

332    families, we looked for enriched GO terms. On Clade I the term related with

333    bacteriocin immunity was enriched;  Clade II were enriched with terms associated

334    to siderophore transport; in Clade III the category related to the biosynthesis of

335    lipopolysaccharides was enriched; and on Clades IV and V there were enriched

336    terms related to tRNA biosynthesis (Additional file 1: Table 9).

## Discussion

338       In this study we performed comparative genomic analyses to understand

339    how evolutionary forces shaped the pan-genome of 42 Vibrionaceae strains

340    isolated from CCB, where environmental filtering is believed to increase local

341    adaptation due to extreme stoichiometric bias [27]. In our study we described how

342    a natural perturbation lead to a temporal balanced stoichiometry, allowing six

343    lineages of Vibrionaceae to prosper under a "feast-famine" cycle. Most of these

344    lineages present large population sizes as well as recombination rates comparable

345    to their oceanic counterparts. However, their pan-genomes remained closed

346    probably due to selection purging HGT events external to each clade where

347    genetic isolation has maintained clade specific selective events. Clade II is the

348    exception, this large clade shows an open pan-genome with evidence of

349    substructure with small effective sizes suggesting early stages of diversification.

350    **Ecology and microbial diversity in CCB.** During the past 20 years, one of the

351    main questions surrounding CCB bacterial hyper-diversity has been related to the

352    roles of ecology and evolution promoting and maintaining its remarkable microbial

353    diversity [27, 46]. According to Souza *et al.* "lost world" hypothesis, the extreme

354    unbalanced stoichiometry (i.e., very low P availability) of CCB not only keeps the

355    "ancestral niche" of many bacterial lineages, but also works as a semipermeable

356    barrier to migration, restricting migration and keeping these ancient bacterial

357    lineages alive and thriving in CCB [27]. As a result of these ecological and

358    evolutionary conditions, CCB lineages are generally clonal [28-30] displaying an

359    ancient marine ancestry [27, 32, 47]. Paradoxically, this extremely unbalanced

360    stoichiometry seems to be in part the reason behind CCB high microbial

361    endemicity and local differentiation: "No food, no sex, no travel"  [27, 31, 32],

362    allowing for local adaptation and broad differentiation between sites.

363        In this study we explored the evolutionary dynamics after a natural

364    perturbation (in this case a flood) changed the ecological condition in CCB in a

365    particular site (Pozas Rojas), generating a temporarily more "balanced"

366    stoichiometric proportions (i.e., N:P 20:1). We know by meteorological data that

367    similar floods occur at CCB sporadically, due to the low incidence of intense storms

368    (i.e., three since 1940 [48]). The flood moved to this low land a large amount of

369    debris that with time, generated an increase in nutrients, in particular phosphorus

370    that opened opportunities for the "rare biosphere", represented by standing

371    bacterial lineages usually found at very low proportions, like the rare members of

372    Vibrionaceae that normally are not common at standard low nutrient conditions [49-

373    51]. Given this change in resources, we proposed two hypotheses when we started

374    this study: Vibrionaceae from CCB would show as their ocean counterparts, an

375    open pan-genome, showing high levels of recombination and genetic variation, as

13

376    well as a high $N_e$. Alternatively, due to local adaptation in each lineage of CCB,

377    Vibrionaceae would display closed pan-genomes, and a strong genetic structure,

378    generated by high clonality and low genetic variation probably related to periodic

379    selection and small effective population sizes among lower levels of genetic

380    variation.

381    **Vibrionaceae in CCB.** In a previous study at Pozas Rojas using both cultivated

382    strains and metagenomic data, Bonilla-Rosso *et al*. found that *Vibrio* spp. was

383    either very rare or absent [49]. In their study, the authors found mostly

384    Pseudomonads among the cultivated strains [49]. This result was confirmed with

385    metagenomics, where Pseudomonadales, Burkholderiales, and Bacillales

386    represented 50% of the metagenome reads. As a result of this previous

387    knowledge, in the 2013 sampling, we first used PIA media to analyze the effect of

388    the 2010 flood in the previously abundant genera, however, we found that this

389    lineage was replaced in the cultures by *Vibrio* spp. In other words, the increased

390    levels of nutrients and the perturbation reduced the abundances of *Pseudomonas*

391    and related genera in CCB. This effect was corroborated later in another system in

392    CCB (Churince) with a nutrient enrichment experiment [50, 52]. Among the

393    analyzed genomes, we found two clades of Vibrionaceae, Clades III and IV, that

394    had not been isolated previously and could be endemic to the basin.

395    **Recombination, pan-genomes, and selection in Vibrionaceae.** Diversity

396    measures, π and θw showed lower diversity than cosmopolitan *E. coli* [53],

397    nevertheless, for Clades I, II and VI, those values are comparable to the ones

398    observed in pathogenic *Vibrio* spp. [54, 55] suggesting similar demographic

399    dynamics. Tajima's D was in most cases negative, except for Clade II, but none of

400    the values were statistically significant. This could suggests bottlenecks in the

401    process of diversification explaining the extremely low effective population size and

402    diversity in those Sub-clades. Negative values of Tajima's D suggest high content

403    of rare alleles, which is in agreement with the private allele test we performed [56].

404    In the same way, it could be the result of selective sweeps or recent demographic

405    expansion as a result of the new nutrient conditions (feast).

14

406        This study corroborates the importance of recombination in Vibrionaceae,

407   supporting the recombinant nature of the genomes in the family [33, 34]. Elevated

408   recombination rates are maintained in all the lineages from Pozas Rojas,

409   supporting our first hypothesis that *Vibrio* spp. from CCB would have similar

410   evolutionary process and genetic structure than marine lineages. Further scrutiny

411   revealed an unexpected result: even if recombination rate is similar to their oceanic

412   counterparts, homologous recombination and selection apparently maintain the

413   adaptation to the local environment. Even in Clade II, recombination is more

414   abundant among related strains, suggesting that this clade is in an actively

415   diversifying process, allowing their different Sub-clades to adapt to different

416   environments within CCB, as it is the case of aquatic Sub-clade G that shares

417   similar SNPs under selection than aquatic Clade III.

418        We believe that the natural disturbance at Pozas Rojas generated by an

419   increase in nutrient availability relaxed selection against HGT. Nevertheless

420   recombination is kept within close lineages resulting in large effective population

421   sizes and a closed pan-genome in most of the lineages, allowing selection to act in

422   response to environmental pressures [57-59]. The closed pan-genome of these

423   lineages contrast to what has been reported in oceanic *Vibrio* spp. where

424   populations sizes are large and pan-genomes are kept open due to HGT [60].

425   Even though Clade II is the only one with an open pan-genome, its internal

426   substructure suggest a recent process of diversification where each of its Sub-

427   clades show again a closed pan-genome, with smaller $N_e$ and low genetic diversity.

428   **Selection and adaptation in Pozas Rojas.** We found 367 gene families that have

429   signals of positive selection, most of them regarding the whole group of ortholog

430   genes found in the core genome (2.05% of the flexible genome and 5% of the core

431   genome; Additional file 2: Table 9). This result suggests that selection purges the

432   genes that are in the flexible genome, closing the pan-genomes. Among the

433   detected genes with selection signals, seven functional GO terms were enriched,

434   one of them was the term GO:0007156, which is associated with cell-cell adhesion;

435   within this category, most of the genes annotated were related to cadherin domains

436   that have been associated to biofilm formation [61]. In natural environments, biofilm

437   formation allows bacteria to cope with environmental changes, protects the cell,

438   provides mechanical stability, and provides cellular adhesion with other cells or

439   with surfaces.  It has been observed that biofilm formation is a persistent

440   characteristic among bacteria from CCB in both water and sediment, and also

441   under different nutrient conditions [52].

442       When we performed a genome-wide association study (GWAS) test to

443   analyze the association of the SNPs to either water or sediment environment, we

444   identified 598 SNPs related to sediment. The UPGMA analysis showed a similar

445   clustering pattern as the core genome (Figure 4), suggesting a clade effect.

446   However, Cluster III grouped among the Sub-clade G of Clade II, and most of the

447   isolates of this Sub-clade as well as Clade III were isolated from the water

448   environment. One possibility is that these SNPs are important to the adaptation to

449   non-structured environments such as water. Some of the genes associated to

450   these SNPs presented signals of recombination and selection. One of the

451   functional enrichment GO terms within these genes was the GO:0006814, which is

452   involved in sodium transport; some of the genes annotated within this category

453   were the bacterial Na+/H+ antiporter B (NhaB) that has been suggested to play a

454   role in the adaptation of halophilic and haloalkaliphilic proteobacteria to marine

455   habitats [62]. This gene has also been found to play a role in homeostasis in *Vibrio*

456   spp. [63]. Our data suggest that there is a selective pressure over some clades

457   regarding the water environment.

458       When we analyzed unique genes for each clade disregarding the isolation

459   environment, in the case of Clade I we found the term GO:0030153 enriched,

460   which is related to bacteriocin immunity. However, antibiotic resistance associated

461   genes did not show particular signals of selection, suggesting that overall there is

462   no ongoing selective pressure for defense. In the large generalist Clade II, we

463   found three GO terms enriched, two of them related to cell wall structure while the

464   third is related to siderophore transport, a group of genes that were rare in the

465   previous metagenomic analysis of the same site [35]. In the case of Clade III, the

16

466 enriched GO term is related to lipopolysaccharide biosynthesis. Meanwhile, in

467 Clade IV, we identified six enriched GO terms, where most of them were related to

468 transport and signal transduction. Finally, for Clade V, we identified four terms

469 enriched mostly related to transport. These results suggest that distinct clades are

470 indeed responding to their environment in different ways reinforcing the idea of

471 genetic isolation as a way to preserve local adaptation (Additional file 2: Table 8).

472 **Perspectives and conclusions.** At CCB, most of the environments present an

473 extremely low phosphorus concentration, a factor that acts as an effective

474 migration barrier maintaining conditions of the ancient sea as well as ancestral

475 microbial diversity [27]. However, due to natural perturbation, we had the

476 opportunity to observe in Pozas Rojas what happens when that nutrimental barrier

477 is lifted temporarily. Apparently, rare biosphere strains that normally had a hard

478 time surviving low P conditions can follow a feast-famine cycles and have

479 population expansion when the P availability is less limiting.

480     In order to understand the other dimensions of local adaptation, further

481 sampling of *Vibrio* spp. in CCB is needed. Unfortunately, this extraordinary oasis is

482 disappearing, given the loss of more than 95% of CCB wetlands due to

483 groundwater overexploitation by agriculture [27, 47,51, 64].

484 **Methods.**

485 **Site description.** We analyzed bacterial isolates from sediment and water of nine

486 ponds in the Pozas Rojas area of CCB (Figure 1). This site is composed of several

487 small ponds (locally called *pozas*) that surround a larger pond in the system of Los

488 Hundidos [30, 35]. These small ponds become hypersaline in summer [30], and

489 used to have the highest stoichiometric unbalance (i.e., lowest P concentration)

490 reported in CCB (C:N:P 15820:157:1)[35]. The ponds have seasonal high

491 fluctuations in temperature (around 1 °C in winter to up to 60 °C in some summer

492 moments in some cases)[35] and are small but permanent, separated from each

493 other by ca. 9 meters or more, along an arch around the larger pond. However, the

494 Pozas Rojas were flooded by hurricane Alex during summer 2010, merging most of

495 the small ponds into a single large pond, until autumn 2011, when the water

496 receded, leaving the moon shaped array of small red ponds at the same place

497 (Figure 1).

498 **Sample collection and strains isolation.** We collected water and sediment

499 samples in duplicate from nine ponds located in Pozas Rojas, Los Hundidos, CCB,

500 during March 2013 and stored them at 4 °C until processing. Sediment was

501 collected for nutrient analysis in 50 ml Falcon tubes and covered with aluminum foil

502 before storage. Water was collected for nutrient quantification in 1 liter volumes

503 and stored in the dark at 4 °C. Chemical analyses were performed at the Instituto

504 de Investigaciones en Ecosistemas y Sustentabilidad, UNAM, in Morelia, Mexico.

505 Cultivable strains from both sediment and water were isolated in PIA

506 (*Pseudomonas* isolation agar) and TCBS (Thiosulfate Citrate Bile Sucrose Agar)

507 as previously described [52, 65], obtaining a total of 174 isolates, being 88 isolates

508 from sediment and 86 from water.

509 **Environmental variables measurement.** For nutrient quantification, sediment

510 samples were dried, and water samples were filtered through a Millipore 0.42 μm

511 filter. Total carbon (TC) and inorganic carbon (IC) were determined by combustion

512 and colorimetric detection [66] using a total carbon analyzer (UIC model CM5012,

513 Chicago, USA). Total organic carbon (TOC) was calculated as the difference

514 between TC and IC.  For total N (TN) and total P (TP) determination, samples were

515 acid digested with $H_2SO_4$, $H_2O_2$, $K_2SO_4$ and $CuSO_4$ at 360°C. Soil N was

516 determined by the macro-Kjeldahl method [67], while P was determined by the

517 molybdate colorimetric method following ascorbic acid reduction [68]. The N and P

518 forms analyzed were determined colorimetrically in a Bran-Luebbe Auto analyzer 3

519 (Norderstedt, Germany).

520 **DNA Extraction and PCR Amplification of 16S rRNA.** For the 174 isolates

521 obtained, DNA extraction was performed as described by Aljanabi and Martinez

522 (1997) [69]. 16S rRNA genes were amplified using universal primers 27F (5′-AGA

523 GTT TGA TCC TGG CTC AG-3′) and 1492R (5′-GGT TAC CTT GTT ACG ACT T-

524 3′) [70]. All reactions were carried out in an Applied Biosystems Veriti 96 Well

525 Thermal cycler (California, USA) using an Amplificasa DNA polymerase

526    (BioTecMol, Mexico) with the following program: 94°C for 5 min, followed by 30

527    cycles consisting of 94°C for 1 min, 50°C for 30 s, 72°C for 1 min and 72°C for 5

528    min. Polymerase chain reaction (PCR) amplification products were

529    electrophoresed on 1% agarose gels. Sanger sequencing was performed at the

530    University of Washington High-Throughput Genomics Center.

531    **Phylogenetic analysis of 16S rRNA sequences.** The first 700 bps of the 16S

532    rRNA gene, were aligned with Clustalw [71] and quality control was performed with

533    Mothur [72]. Genera level identification was made using the classifier tool [73] from

534    the Ribosomal Database Project (RDP) Release 11.4 [74] (Additional file 1: Table

535    3). Blastn searches were performed against Refseq database from NCBI to select

536    reference sequences. A total of 101 sequences were identified as members of the

537    Vibrionaceae family, 41 were isolates from water and 60 from sediment. These

538    isolates were used in subsequent analyses. A maximum likelihood phylogenetic

539    reconstruction was obtained with PhyML version 3.0 [75], using the HKY+I+G

540    substitution model estimated with jModelTest 2 [76]. The degree of support for the

541    branches was determined with 1,000 bootstrap iterations.

542    **Environmental association of phylogroups.** To test whether the community of

543    cultivable strains was structured based on its isolation environment (i.e., water or

544    sediment), we performed an AdaptML analysis [37], including our 101 isolates

545    belonging to Vibrionaceae and an *Halomonas* spp*.* strain as an out-group. Three

546    categorical environmental variables were tested, including pond of isolation, high

547    and low nutrient concentrations, and the two sampled environments (water or

548    sediment).

549    **Genome sequencing, assembly, and annotation.** For whole-genome

550    sequencing, we selected from the AdaptML analysis 39 *Vibrio* spp. isolates, 23

551    isolated from sediment and 16 from water, plus 3 isolates of *Photobacterium* spp*.*

552    (a lineage closely related to the *Vibrio* spp. genus) isolated from sediment. DNA

553    extractions were performed with the DNeasy Blood and Tissue kit (Qiagen).

554    Sequencing was performed with Illumina MiSeq 2x250 technology, with

555    insert libraries of 650 bps and an expected coverage of ca.10x per genome. At

556 first, we planned an assembly strategy using a genome reference; for this reason,

557 the strain V15_P4S5T153 had a second library that was designed using the Jr 454

558 Roche technology, in order to reduce sequencing bias and get higher coverage.

559 However, due to divergence among genomes, we performed *de novo* assemblies

560 for all genomes. All sequencing was performed at the Laboratorio Nacional de

561 Genómica para la Biodiversidad (LANGEBIO), México.

562 The quality of raw reads was analyzed using FASTQC software

563 (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/). A minimum quality

564 value of 25 was set, and low-quality sequences were removed with

565 fastq_quality_filter from the FASTX-Toolkit

566 (http://hannonlab.cshl.edu/fastx_toolkit/index.html). Adapter sequences were

567 identified, removed and paired-end reads were merged using SeqPrep

568 (https://github.com/jstjohn/SeqPrep). *De novo* assemblies were performed with

569 Newbler (Roche/ 454 Life Sciences) using both single-end and merged reads.

570 For scaffolding process, we used SSPACE [77], gaps were closed using

571 GapFiller [78] and final error correction was performed with iCORN [79] (Additional

572 file 2: Table 10). Coding sequences were inferred with Prodigal 2.0 [80]

573 implemented in PROKKA software [81]. InterProScan 5 allowed annotation [82]

574 with the databases enabled by default. Genome completeness was assessed with

575 BUSCO using the Gamma-proteobacteria database [38].

576 **Pan-genome analyses.** The 42 genomes from CCB where compared with

577 genomes of 5 reference *Vibrio* spp. strains: *Vibrio alginolyticus NBRC 15630 =*

578 *ATCC 17749*, *V. anguillarum 775*, *V. furnissii NCTC 11218*, *V. parahaemolyticus*

579 *BB22OP* and *V. metschnikovii CIP 69 14* (Additional file 1: Tables 4, 5; Additional

580 file 2: Tables 10). Ortholog gene families were predicted from all 47 genomes using

581 the DeNoGAP comparative genomics pipeline [83]. To minimize false positive

582 prediction of orthologs, we assigned *Photobacterium* spp. genomes as outgroup.

583 The completely sequenced genome of *V. anguillarum* strain 775 was used as seed

584 reference.

20

585       We estimated the core genome based on presence and absence of gene

586    families across the genomes. If the genes were present in all strains, the orthologs

587    were classified as *core*, while genes were classified as *accessory* when present in

588    more than one strain but not in all of them, and *unique* genes when it was present

589    only in a single strain. Since most of the genomes in our dataset are not completely

590    sequenced, we designated core ortholog families as those present in at least 95%

591    of the genomes, to avoid the impact of missing genes due to sequencing or

592    assembly artifacts.

593       The package Micropan [84] within R v.3.4 (R Core Team) [85] was used to

594    infer the open or closed nature of each pan-genome dataset, following the heaps

595    law proposed by Tettelin *et al.* [42]. The Heaps law model is fitted to the number of

596    new gene clusters observed when genomes are ordered randomly. The model has

597    two parameters: an intercept, and a decay parameter called alpha. If alpha is

598    higher than 1.0 the pan-genome is considered closed, if alpha is lower than 1.0 it is

599    considered open. Additionally, a random sub-sampling for each clade was made,

600    taking three genomes and calculating the alpha value for each group of three

601    genomes. A total of 1,000 independent sub-sampling events were made for each

602    clade.

603       Core proteins were aligned using Kalign [86] to infer the phylogenetic

604    relationship between the samples. The resulting alignments of individual ortholog

605    families were concatenated using a custom Perl script. With these concatenated

606    core genes, a maximum likelihood phylogenetic tree was constructed using the

607    FastTree program [87].

608    **Recombination analyses.** Of the total ortholog families in the *Vibrio* spp. pan-

609    genome, we only used the ortholog families found in at least three genomes for the

610    recombination analyses. Genetic recombination was examined on each CDS

611    alignment by using inference of pairwise recombination sites, obtained with

612    GENECONV [88] and by the identification of putative recombinant sequences

613    through breakpoints using GARD [89].

614 Based on the number of recombination events, we estimated the events
615 shared among isolates of the same pond and environment, among isolates of the
616 different pond and environment, among isolates of the same pond and different
617 environment and among isolates of different pond and environment. For this, we
618 normalized the data by pan-genome size, number of strains and branch length.
619 Given that the large generalist Clade II presented a clear sub-structure, we did a
620 separated analysis for the shorter branches within Clade II (Additional file 1: Figure
621 4).

622 To assess the impact of homologous recombination, we analyzed the
623 substitution pattern using two different algorithms, Gubbins [90] and
624 ClonalFrameML [43]. A whole-genome alignment for the 47 analyzed genomes
625 was performed with MAUVE [91]. The resulting alignment was used as input for
626 Gubbins [90] using RAxML [92] and default parameters. Additionally, whole
627 genome alignments were performed for each clade, excluding references, with the
628 progressive MAUVE algorithm [91]. We calculated the R/theta ratio, nu and delta
629 [43] for each sample and for 100 bootstrapped replicates.

630 **Genetic structure of Clade II.** Recombination analyses showed that in Clade II
631 there are internal groups with higher internal recombination, so we decided to
632 further investigate the structure within Clade II. For clustering analyses, we used
633 Nei's genetic distance [93] and neighbor joining. Genomes with distance less than
634 0.001 were grouped and tested with a discriminant analysis of principal
635 components of the genetic variation, using the adegenet library in R [94]. For this
636 study, we used 20 principal components and 3 discriminant functions.

637

638 **Selection analyses.** We used FUBAR [95] to identify signatures of positive
639 selection among ortholog gene families found in at least three genomes. We
640 accounted for recombination breakpoints in the ortholog families, while calculating
641 positively selected sites based on GARD results [89]. We considered any site to be
642 positively selected if it showed P-value $\leq$ 0.05. We also conducted a Gene

643 Ontology (GO) enrichment analysis using topGO [45] to find overrepresented

644 biological functions in this set of genes.

645 **Effective population size estimation.** We followed a simulation approach to

646 estimate the posterior distribution of the effective population size ($N_e$) of each of

647 the six clades. According to the previous clustering and recombination analysis, for

648 Clades I, III, IV, V and VI we simulated a single population, while for Clade II we

649 simulated three sub-populations that diverged from an ancestral population.

650 Simulations were performed using Fastsimcoal2 [44, 96]. For each clade,

651 we simulated DNA sequences having a similar length equal to the number of

652 nucleotides in the given clade, as well as a sample size equal to the number of

653 sequences sampled for each clade. We assumed no recombination within the

654 genome, and used the *Escherichia coli* mutation rate of $2.2 \times 10^{-10}$ mutations per

655 nucleotide per generation [97]. We ran between two and four simulations for each

656 clade. For the initial runs, we generated 100,000 replicates extracting $N_e$ values

657 from a prior log-uniform distribution that ranged from 100,000 to 20,000,000

658 individuals. For Clade II, we also estimated the age of divergence of each Sub-

659 clade, by setting the prior distribution of time ranging from 1,000-4,000,000

660 generations. After a first run, we narrowed the prior ranges based on those

661 simulations that had similar summary statistics compared to the observed data and

662 performed another 100,000 simulations using the narrowed priors.

663 To compare the previously simulated and observed data based on summary

664 statistics, we used the ape [98] and pegas [99] libraries in R to estimate the

665 number of polymorphic sites and the Tajima's *D* based on the entire genomes.

666 Tajima's *D* is commonly used to estimate demographic changes in populations

667 [100, 101]. Also, we obtained 1,000 sliding windows frames to estimate the

668 Tajima's *D* along the genomes, as well as the mean and standard deviation of

669 Tajima's *D*. Tajima's *D,* π, and Watterson's theta (θw) were estimated for each

670 clade as well as for Sub-clades A, B and G. Since clades I and VI had three

671 sequences and it was not possible to obtain Tajima's *D*, we did 1,000 replicates in

672 which we subsampled with replacement 10 sequences. For each replicate, we

23

673    calculated Tajima's $D$ and we obtained as the proximate value the median

674    estimated across the 1,000 replicates.

675    Based on the summary statistics, we used the abc function in the ABC

676    package [102] in R to calculate the distribution of the $N_e$ parameter based on a

677    0.05 % threshold distance between the simulated and observed data. For each

678    clade, we report the median and the 95% interval confidence of $N_e$. For Clade II,

679    we further reported the average and 95% interval confidence of the number of

680    generations since each Sub-clade diverged from an ancestral clade.

681    **Association between genotypes and environmental variables**. We evaluated

682    whether the genetic variation within the Vibrionaceae genomes could be explained

683    by particular adaptations to the environment (water or sediment). We used

684    progressiveMauve [91] to perform a global multiple alignment between the

685    assembled genomes. We extracted the variant sites within the alignment and

686    exported them as SNPs using snp-sites [103].

687    We obtained 38,533 SNPs, which we used to search for private alleles using

688    Poppr [104]. Afterwards, we obtained a subset of 25,892 SNPs by filtering biallelic

689    sites with minor allele frequencies > 0.05. We used PLINK [105] to perform a

690    GWAS to detect possible associations between our SNP set and either the water

691    or sediment environments. We conducted Fisher exact tests and regarded as

692    significant all SNPs whose associations had $p$-values < 0.01 after Bonferroni

693    corrections. These analyses may be informative even considering these sampling

694    differences [106, 107].

695    To test whether these associations could be explained by convergent

696    evolution rather than by common ancestry, we compared an UPGMA tree

697    reconstructed from the total set of SNPs from an UPGMA tree using only the SNPs

698    that were significantly associated to the environment. We analyzed the distribution

699    of the SNPs within the genomes to find the genes associated to those SNPs.

700    We mapped the SNPs positions in the genome alignment moving by 1 Kb

701    windows; this window size was selected considering the average bacterial gene

24

702 size and retrieved all the associated genes. We conducted a Gene Ontology (GO)

703 enrichment analysis using topGO [45] to find overrepresented biological functions

704 in this set of genes.

**Availability of data and materials**

706 The datasets generated and analysed during the current study are available in the

707 genome assembly project BioProject: PRJNA361510; PRJNA361511. The

708 resulting InterProScan annotation files, CDS fasta files and the predicted protein

709 fasta files for all taxa are available at Dryad. As by the politics of Dryad, the data

710 will be available once the manuscript is accepted.

**Competing interests**

712 The authors declare that they have no competing interests

**Funding**

**Author Contributions.**

723 MV-R-L design the sampling, obtained the biological material, analyzed the data,

724 prepared figures and tables, and wrote the paper. GYP-S analyzed the data and

725 participated in all stages of writing. JA-L, ST, ES, and JB-R analyzed the data. EI-L

726 analyzed the data and provided computing facilities. DS-G provided computing

727 facilities and contributed substantially to the analysis and discussion of the data.

728 LEE made contributions for the design, analysis, discussion of the data and writing.

729 V-S conceived, designed the study and the analyses, managed the obtaining

730 financial resources and participated in all stages of writing.

**Acknowledgments.**

**References.**

737 1.  Lilburn TG, Gu J, Cai H, Wang Y. Comparative genomics of the family

738 vibrionaceae reveals the wide distribution of genes encoding virulence-associated

739 proteins. BMC Genomics. 2010;11:369. doi:10.1186/1471-2164-11-369.

740 2.  Moriel DG, Tan L, Goh KGK, Phan M-D, Ipe DS, Lo AW, et al. A novel

741 protective vaccine antigen from the core Escherichia coli genome. mSphere.

742 2016;1. doi:10.1128/msphere.00326-16.

743 3.  Sanglas A, Albarral V, Farfán M, Lorén JG, Fusté MC. Evolutionary roots and

744 diversification of the genus Aeromonas. Frontiers in Microbiology. 2017;8.

745 doi:10.3389/fmicb.2017.00127.

746 4.  Lapierre P, Gogarten JP. Estimating the size of the bacterial pan-genome.

747 Trends in Genetics. 2009;25:107–10. doi:10.1016/j.tig.2008.12.004.

748 5.  Collins RE, Higgs PG. Testing the infinitely many genes model for the evolution

749 of the bacterial core genome and pangenome. Molecular Biology and Evolution.

750 2012;29:3413–25. doi:10.1093/molbev/mss163.

751 6.  Gordienko EN, Kazanov MD, Gelfand MS. Evolution of pan-genomes of

752 Escherichia coli, Shigella spp., and Salmonella enterica. Journal of Bacteriology.

753 2013;195:2786–92. doi:10.1128/jb.02285-12.

754 7.  Valdivia-Anistro JA, Eguiarte-Fruns LE, Delgado-Sapién G, Gasca-Pineda PM-

755 ZJ, Learned J, Elser JJ, et al. Variability of rRNA operon copy number and growth

756  rate dynamics of bacillus isolated from an extremely oligotrophic aquatic

757  ecosystem. Frontiers in Microbiology. 2016;6. doi:10.3389/fmicb.2015.01486

758  8.   Zhi X-Y, Jiang Z, Yang L-L, Huang Y. The underlying mechanisms of genetic

759  innovation and speciation in the family corynebacteriaceae : A phylogenomics

760  approach. Molecular Phylogenetics and Evolution. 2017;107:246–55.

761  doi:10.1016/j.ympev.2016.11.009.

762  9.   Hou Y, Lin S. Distinct gene number-genome size relationships for eukaryotes

763  and non-eukaryotes: Gene content estimation for dinoflagellate genomes. PLoS

764  ONE. 2009;4:e6978. doi:10.1371/journal.pone.0006978.

765  10.   McInerney JO, McNally A, O MJ. Why prokaryotes have pangenomes. Nature

766  Microbiology. 2017;2. doi:10.1038/nmicrobiol.2017.40.

767  11.   Kuo C-H, Ochman H. Deletional bias across the three domains of life.

768  Genome Biology and Evolution. 2009;1:145–52. doi:10.1093/gbe/evp016.

769  12.   Morris JJ, Lenski RE, Zinser ER. The black queen hypothesis: Evolution of

770  dependencies through adaptive gene loss. mBio. 2012;3. doi:10.1128/mbio.00036

771  13.   Mas A, Jamshidi S, Lagadeuc Y, Eveillard D, Vandenkoornhuyse P. Beyond

772  the black queen hypothesis. The ISME Journal. 2016;10:2085–91.

773  doi:10.1038/ismej.2016.22.

774  14.   Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, et al.

775  Genome analysis of multiple pathogenic isolates of Streptococcus agalactiae:

776  Implications for the microbial "pan-genome". Proceedings of the National Academy

777  of Sciences. 2005;102:13950–5. doi:10.1073/pnas.0506758102.

778  15.   Andreani NA, Hesse E, Vos M. Prokaryote genome fluidity is dependent on

779  effective population size. The ISME Journal. 2017;11:1719–21.

780  doi:10.1038/ismej.2017.36.

781  16.   Smith JM, Smith NH, Spratt MOBG. How clonal are bacteria? Proceedings of

782  the National Academy of Sciences. 1993;90:4384–8. doi:10.1073/pnas.90.10.4384.

783    17.   Souza V, Eguiarte LE. Bacteria gone native vs. bacteria gone awry?:

784    Plasmidic transfer and bacterial evolution. Proceedings of the National Academy of

785    Sciences. 1997;94:5501–3. doi:10.1073/pnas.94.11.5501.

786    18.   Lawrence JG, Ochman H. Molecular archaeology of the Escherichia coli

787    genome. Proceedings of the National Academy of Sciences. 1998;95:9413–7.

788    doi:10.1073/pnas.95.16.9413.

789    19.   Ochman H, Lawrence JG, Groisman EA. Lateral gene transfer and the nature

790    of bacterial innovation. Nature. 2000;405:299–304. doi:10.1038/35012500.

791    20.   Fournier GP, Gogarten JP. Evolution of acetoclastic methanogenesis in

792    methanosarcina via horizontal gene transfer from cellulolytic clostridia. Journal of

793    Bacteriology. 2007;190:1124–7. doi:10.1128/jb.01382-07.

794    21.   Soucy SM, Fullmer MS, Papke RT, Gogarten JP. Inteins as indicators of gene

795    flow in the halobacteria. Frontiers in Microbiology. 2014;5.

796    doi:10.3389/fmicb.2014.00299.

797    22.   Roze D, Barton NH. The hill Robertson effect and the evolution of

798    recombination. Genetics. 2006;173:1793–811. doi:10.1534/genetics.106.058586.

799    23.   Comeron JM, Williford A, Kliman RM. The hill Robertson effect: Evolutionary

800    consequences of weak selection and linkage in finite populations. Heredity.

801    2007;100:19–31. doi:10.1038/sj.hdy.6801059.

802    24.   Souza V, Nguyen TT, Hudson RR, Pinero D, Lenski RE. Hierarchical analysis

803    of linkage disequilibrium in Rhizobium populations: Evidence for sex? Proceedings

804    of the National Academy of Sciences. 1992;89:8389–93.

805    doi:10.1073/pnas.89.17.8389.

806    25.   Bobay L-M, Ochman H. Factors driving effective population size and pan-

807    genome evolution in bacteria. BMC Evolutionary Biology. 2018;18.

808    doi:10.1186/s12862-018-1272-4.

809    26.   Cohan FM. Bacterial species and speciation. Systematic Biology.

810    2001;50:513–24. doi:10.1080/10635150118398.

811   27.   Souza V, Moreno-Letelier A, Travisano M, Alcaraz LD, Olmedo G, Eguiarte

812   LE. The lost world of Cuatro Ciénegas basin, a relictual bacterial niche in a desert

813   oasis. eLife. 2018;7. doi:10.7554/elife.38278.

814   28.   Escalante AE, Eguiarte LE, Espinosa-Asuar L, Forney LJ, Noguez AM,

815   Saldivar VS. Diversity of aquatic prokaryotic communities in the Cuatro Cienegas

816   basin. FEMS Microbiology Ecology. 2008;65:50–60. doi:10.1111/j.1574-

817   6941.2008.00496.x.

818   29.   Rebollar EA, Avitia M, Eguiarte LE, González-González A, Mora L, Bonilla-

819   Rosso G, et al. Water-sediment niche differentiation in ancient marine lineages of

820   Exiguobacterium endemic to the Cuatro Cienegas basin. Environmental

821   Microbiology. 2012;14:2323–33. doi:10.1111/j.1462-2920.2012.02784.x.

822   30.   Avitia M, Escalante AE, Rebollar EA, Moreno-Letelier A, Eguiarte LE, Souza

823   V. Population expansions shared among coexisting bacterial lineages are revealed

824   by genetic evidence. PeerJ. 2014;2:e696. doi:10.7717/peerj.696.

825   31.   Souza V, Eguiarte LE, Siefert J, Elser JJ. Microbial endemism: Does

826   phosphorus limitation enhance speciation? Nature Reviews Microbiology.

827   2008;6:559–64. doi:10.1038/nrmicro1917.

828   32.   Souza V, Eguiarte LE, Travisano M, Elser JJ, Rooks C, Siefert JL. Travel,

829   sex, and food: Whats speciation got to do with it? Astrobiology. 2012;12:634–40.

830   doi:10.1089/ast.2011.0768.

831   33.   Vos M, Didelot X. A comparison of homologous recombination rates in

832   bacteria and archaea. The ISME Journal. 2008;3:199–208.

833   doi:10.1038/ismej.2008.93.

834   34.   Cui Y, Yang X, Didelot X, Guo C, Li D, Yan Y, et al. Epidemic clones, oceanic

835   gene pools, and eco-LD in the free living marine pathogen Vibrio

836   parahaemolyticus. Molecular Biology and Evolution. 2015;32:1396–410.

837   doi:10.1093/molbev/msv009.

838   35.   Peimbert M, Alcaraz LD, Bonilla-Rosso G, Olmedo-Alvarez G, Garc-Oliva F,

839   Segovia L, et al. Comparative metagenomics of two microbial mats at Cuatro

840   Ciénegas basin II: Ancient lessons on how to cope with an environment under

841   severe nutrient stress. Astrobiology. 2012;12:648–58. doi:10.1089/ast.2011.0694.

842   36.   Redfield AC. James Johnstone Memorial Volume. Daniel RJ, ed. Liverpool

843   Univ. Press;1934. p. 176-92

844   37.   Hunt DE, David LA, Gevers D, Preheim SP, Alm EJ, Polz MF. Resource

845   partitioning and sympatric differentiation among closely related bacterioplankton.

846   Science. 2008;320:1081–5. doi:10.1126/science.1157890.

847   38.   Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM.

848   BUSCO: Assessing genome assembly and annotation completeness with single-

849   copy orthologs. Bioinformatics. 2015;31:3210–2.

850   doi:10.1093/bioinformatics/btv351.

851   39.   Lux TM, Lee R, Love J. Complete genome sequence of a free-living Vibrio

852   furnissii sp. nov. strain (NCTC 11218). Journal of Bacteriology. 2011;193:1487–8.

853   doi:10.1128/jb.01512-10.

854   40.   Naka H, Dias GM, Thompson CC, Dubay C, Thompson FL, Crosa JH.

855   Complete genome sequence of the marine fish pathogen Vibrio anguillarum

856   harboring the pJM1 virulence plasmid and genomic comparison with other virulent

857   strains of V. anguillarum and V. ordalii. Infection and Immunity. 2011;79:2889–900.

858   doi:10.1128/iai.05138-11.

859   41.   Xu F, Ilyas S, Hall JA, Jones SH, Cooper VS, Whistler CA. Genetic

860   characterization of clinical and environmental Vibrio parahaemolyticus from the

861   northeast USA reveals emerging resident and non-indigenous pathogen lineages.

862   Frontiers in Microbiology. 2015;6. doi:10.3389/fmicb.2015.00272.

863   42.   Tettelin H, Riley D, Cattuto C, Medini D. Comparative genomics: The bacterial

864   pan-genome. Current Opinion in Microbiology. 2008;11:472–7.

865   doi:10.1016/j.mib.2008.09.006.

866  43.  Didelot X, Wilson DJ. ClonalFrameML: Efficient inference of recombination in

867  whole bacterial genomes. PLOS Computational Biology. 2015;11:e1004041.

868  doi:10.1371/journal.pcbi.1004041.

869  44.  Excoffier L, Foll M. Fastsimcoal: A continuous-time coalescent simulator of

870  genomic diversity under arbitrarily complex evolutionary scenarios. Bioinformatics.

871  2011;27:1332–4. doi:10.1093/bioinformatics/btr124.

872  45.  Alexa A, Rahnenfuhrer J, Lengauer T. Improved scoring of functional groups

873  from gene expression data by decorrelating GO graph structure. Bioinformatics.

874  2006;22:1600–7. doi:10.1093/bioinformatics/btl140.

875  46.  Taboada B, Isa P, Gutiérrez-Escolano AL, del Ángel RM, Ludert JE, Vázquez

876  N, et al. The geographic structure of viruses in the Cuatro Ciénegas basin, a

877  unique oasis in northern Mexico, reveals a highly diverse population on a small

878  geographic scale. Applied and Environmental Microbiology. 2018;84.

879  doi:10.1128/aem.00465-18.

880  47.  Souza V, Espinosa-Asuar L, Escalante AE, Eguiarte LE, Farmer J, Forney L,

881  et al. An endangered oasis of aquatic microbial biodiversity in the Chihuahuan

882  desert. Proceedings of the National Academy of Sciences. 2006;103:6565–70.

883  doi:10.1073/pnas.0601434103.

884  48.  Montiel-González C, Bautista F, Delgado C, García-Oliva F. The Climate of

885  Cuatro Ciénegas Basin: Drivers and Temporal Patterns. In Souza V, Olmedo-

886  Álvarez G, Eguiarte LE, eds. Cuatro Ciénegas Ecology, Natural History and

887  Microbiology. New York, NY:Springer, Cham; 2018. p. 35-42

888  49.  Bonilla-Rosso G, Peimbert M, Alcaraz LD, Hernández I, Eguiarte LE, Olmedo-

889  Alvarez G, et al. Comparative metagenomics of two microbial mats at Cuatro

890  Ciénegas basin II: Community structure and composition in oligotrophic

891  environments. Astrobiology. 2012;12:659–73. doi:10.1089/ast.2011.0724.

892  50.  Lee ZM-P, Poret-Peterson AT, Siefert JL, Kaul D, Moustafa A, Allen AE, et al.

893  Nutrient stoichiometry shapes microbial community structure in an evaporitic

894  shallow pond. Frontiers in Microbiology. 2017;8. doi:10.3389/fmicb.2017.00949.

895    51.   Anda VD, Zapata-Peñasco I, Blaz J, Poot-Hernández AC, Contreras-Moreira

896    B, González-Laffitte M, et al. Understanding the mechanisms behind the response

897    to environmental perturbation in microbial mats: A metagenomic-network based

898    approach. Frontiers in Microbiology. 2018;9. doi:10.3389/fmicb.2018.02606.

899    52.   Ponce-Soto GY, Aguirre-von-Wobeser E, Eguiarte LE, Elser JJ, Lee ZM-P,

900    Souza V. Enrichment experiment changes microbial interactions in an ultra-

901    oligotrophic environment. Frontiers in Microbiology. 2015;6.

902    doi:10.3389/fmicb.2015.00246.

903    53.   Ghalayini M, Launay A, Bridier-Nahmias A, Clermont O, Denamur E, Lescat

904    M, et al. Evolution of a dominant natural isolate of escherichia coli in the human gut

905    over the course of a year suggests a neutral evolution with reduced effective

906    population size. Applied and Environmental Microbiology. 2018;84.

907    doi:10.1128/aem.02377-17.

908    54.   Farfan M, Minana-Galbis D, Fuste MC, Loren JG. Allelic diversity and

909    population structure in Vibrio cholerae o139 bengal based on nucleotide sequence

910    analysis. Journal of Bacteriology. 2002;184:1304–13. doi:10.1128/jb.184.5.1304-

911    1313.2002.

912    55.   Gonzalez-Escalona N, Martinez-Urtaza J, Romero J, Espejo RT, Jaykus L-A,

913    DePaola A. Determination of molecular phylogenetics of Vibrio parahaemolyticus

914    strains by multilocus sequence typing. Journal of Bacteriology. 2008;190:2831–40.

915    doi:10.1128/jb.01808-07.

916    56.   Korneliussen TS, Moltke I, Albrechtsen A, Nielsen R. Calculation of tajima's d

917    and other neutrality test statistics from low depth next-generation sequencing data.

918    BMC Bioinformatics. 2013;14. doi:10.1186/1471-2105-14-289.

919    57.   Petit N, Barbadilla A. Selection efficiency and effective population size in

920    Drosophila species. Journal of Evolutionary Biology. 2009;22:515–26.

921    doi:10.1111/j.1420-9101.2008.01672.x.

922    58.    Jensen JD, Bachtrog D. Characterizing the influence of effective population

923    size on the rate of adaptation: Gillespie's darwin domain. Genome Biology and

924    Evolution. 2011;3:687–701. doi:10.1093/gbe/evr063.

925    59.    Gossmann TI, Keightley PD, Eyre-Walker A. The effect of variation in the

926    effective population size on the rate of adaptive molecular evolution in eukaryotes.

927    Genome Biology and Evolution. 2012;4:658–67. doi:10.1093/gbe/evs027.

928    60.    Shapiro BJ, Friedman J, Cordero OX, Preheim SP, Timberlake SC, Szabó G,

929    et al. Population genomics of early events in the ecological differentiation of

930    bacteria. Science. 2012;336:48–51. doi:10.1126/science.1218198.

931    61.    Vozza NF, Abdian PL, Russo DM, Mongiardini E, Lodeiro A, Molin S, et al. A

932    Rhizobium leguminosarum CHDL- (cadherin-like-) lectin participates in assembly

933    and remodeling of the biofilm matrix. Frontiers in Microbiology. 2016;7.

934    doi:10.3389/fmicb.2016.01608.

935    62.    Kurz M, Brünig AN, Galinski EA. NhaD type sodium/proton-antiporter of

936    Halomonas elongata: a salt stress response mechanism in marine habitats? Saline

937    Systems. 2006;2:10. doi:10.1186/1746-1448-2-10

938    63.    Vimont S, Berche P. NhaA, an Na(+)/H(+) antiporter involved in environmental

939    survival of Vibrio cholerae. Journal of Bacteriology. 2000;182:2937–44.

940    doi:10.1128/jb.182.10.2937-2944.2000.

941    64.    Wolaver BD, Crossey LJ, Karlstrom KE, Banner JL, Cardenas MB, Ojeda CG,

942    et al. Identifying origins of and pathways for spring waters in a semiarid basin using

943    he, sr, and c isotopes: Cuatro Cienegas basin, Mexico. Geosphere. 2012;9:113–

944    25. doi:10.1130/ges00849.1.

945    65.    Vázquez-Rosas-Landa M, Ponce-Soto GY, Eguiarte LE, Souza V.

946    Comparative genomics of free-living gammaproteobacteria: Pathogenesis-related

947    genes or interaction-related genes? Pathogens and Disease. 2017;75.

948    doi:10.1093/femspd/ftx059.

949  66.  Huffman EW. Performance of a new automatic carbon dioxide coulometer.

950  Microchemical Journal. 1977;22:567–73. doi:10.1016/0026-265x(77)90128-x.

951  67.  Bremner JM. Total nitrogen. In: Sparks DL ed. Methods of Soil Analysis. Part

952  2 Chemical Methods. Madison, WI: Soil Science Society of America; 1996. p.

953  1085-6

954  68.  Murphy J, Riley J. A modified single solution method for the determination of

955  phosphate in natural waters. Analytica Chimica Acta. 1962;27:31–6.

956  doi:10.1016/s0003-2670(00)88444-5.

957  69.  Aljanabi S. Universal and rapid salt-extraction of high quality genomic DNA for

958  PCR- based techniques. Nucleic Acids Research. 1997;25:4692–3.

959  doi:10.1093/nar/25.22.4692.

960  70.  Lane DJ. 16S/23S rRNA Sequencing. In: Stackebrandt E, Goodfellow M, eds.

961  Nucleic Acid Techniques in Bacterial Systematic. New York: John Wiley and Sons;

962  1991. p. 115-75

963  71.  Larkin M, Blackshields G, Brown N, Chenna R, McGettigan P, McWilliam H, et

964  al. Clustal W and Clustal X version 2.0. Bioinformatics. 2007;23:2947–8.

965  doi:10.1093/bioinformatics/btm404.

966  72.  Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al.

967  Introducing mothur: Open-source, platform-independent, community-supported

968  software for describing and comparing microbial communities. Applied and

969  Environmental Microbiology. 2009;75:7537–41. doi:10.1128/aem.01541-09.

970  73.  Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive bayesian classifier for rapid

971  assignment of rRNA sequences into the new bacterial taxonomy. Applied and

972  Environmental Microbiology. 2007;73:5261–7. doi:10.1128/aem.00062-07.

973  74.  Cole JR, Wang Q, Cardenas E, Fish J, Chai B, Farris RJ, et al. The ribosomal

974  database project: Improved alignments and new tools for rRNA analysis. Nucleic

975  Acids Research. 2009;37 Database:D141–5. doi:10.1093/nar/gkn879.

976  75.  Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. New

977      algorithms and methods to estimate maximum-likelihood phylogenies: Assessing

978      the performance of PhyML 3.0. Systematic Biology. 2010;59:307–21.

979      doi:10.1093/sysbio/syq010.

980  76.  Darriba D, Taboada GL, Doallo R, Posada D. jModelTest 2: More models,

981      new heuristics and parallel computing. Nature Methods. 2012;9:772–2.

982      doi:10.1038/nmeth.2109.

983  77.  Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. Scaffolding pre-

984      assembled contigs using SSPACE. Bioinformatics. 2010;27:578–9.

985      doi:10.1093/bioinformatics/btq683.

986  78.  Nadalin F, Vezzi F, Policriti A. GapFiller: A de novo assembly approach to fill

987      the gap within paired reads. BMC Bioinformatics. 2012;13. doi:10.1186/1471-2105-

988      13-s14-s8.

989  79.  Otto TD, Sanders M, Berriman M, Newbold C. Iterative correction of reference

990      nucleotides (iCORN) using second generation sequencing technology.

991      Bioinformatics. 2010;26:1704–7. doi:10.1093/bioinformatics/btq269.

992  80.  Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal:

993      Prokaryotic gene recognition and translation initiation site identification. BMC

994      Bioinformatics. 2010;11. doi:10.1186/1471-2105-11-119.

995  81.  Seemann T. Prokka: Rapid prokaryotic genome annotation. Bioinformatics.

996      2014;30:2068–9. doi:10.1093/bioinformatics/btu153.

997  82.  Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al.

998      InterProScan 5: Genome-scale protein function classification. Bioinformatics.

999      2014;30:1236–40. doi:10.1093/bioinformatics/btu031.

1000  83.  Thakur S, Guttman DS. A de-novo genome analysis pipeline (DeNoGAP) for

1001      large-scale comparative prokaryotic genomics studies. BMC Bioinformatics.

1002      2016;17. doi:10.1186/s12859-016-1142-2.

1003    84.   Snipen L, Liland KH. Micropan: An R-package for microbial pan-genomics.

1004    BMC Bioinformatics. 2015;16. doi:10.1186/s12859-015-0517-0.

1005    85.   R Core Team (2017). R: A language and environment for statistical

1006    computing. R Foundation for Statistical Computing, Vienna, Austria.

1007    (https://www.R-project.org/)

1008    86.   Lassmann T, Frings O, Sonnhammer ELL. Kalign2: High-performance

1009    multiple alignment of protein and nucleotide sequences allowing external features.

1010    Nucleic Acids Research. 2008;37:858–65. doi:10.1093/nar/gkn1006.

1011    87.   Price MN, Dehal PS, Arkin AP. FastTree 2 approximately maximum-likelihood

1012    trees for large alignments. PLoS ONE. 2010;5:e9490.

1013    doi:10.1371/journal.pone.0009490.

1014    88.   Sawyer S. Statistical tests for detecting gene conversion. Molecular Biology

1015    and Evolution. 1989;6:526-38. doi;10.1093/oxfordjournals.molbev.a040567

1016    89.   Pond SLK, Posada D, Gravenor MB, Woelk CH, Frost SD. GARD: A genetic

1017    algorithm for recombination detection. Bioinformatics. 2006;22:3096–8.

1018    doi:10.1093/bioinformatics/btl474.

1019    90.   Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al.

1020    Rapid phylogenetic analysis of large samples of recombinant bacterial whole

1021    genome sequences using gubbins. Nucleic Acids Research. 2014;43:e15–5.

1022    doi:10.1093/nar/gku1196.

1023    91.   Darling AE, Mau B, Perna NT. progressiveMauve: Multiple genome alignment

1024    with gene gain, loss and rearrangement. PLoS ONE. 2010;5:e11147.

1025    doi:10.1371/journal.pone.0011147.

1026    92.   Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-

1027    analysis of large phylogenies. Bioinformatics. 2014;30:1312–3.

1028    doi:10.1093/bioinformatics/btu033.

1029    93.   Nei, M. (1978). Estimation of average heterozygosity and genetic distance

1030    from a small number of individuals. Genetics, 89(3):583-590.

1031   94.   Jombart T, Ahmed I. Adegenet 1.3-1: New tools for the analysis of genome-

1032   wide SNP data. Bioinformatics. 2011;27:3070–1.

1033   doi:10.1093/bioinformatics/btr521.

1034   95.   Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Pond SLK, et al.

1035   FUBAR: A Fast, Unconstrained BAyesian AppRoximation for inferring selection.

1036   Molecular Biology and Evolution. 2013;30:1196–205. doi:10.1093/molbev/mst030.

1037   96.   Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. Robust

1038   demographic inference from genomic and SNP data. PLoS Genetics.

1039   2013;9:e1003905. doi:10.1371/journal.pgen.1003905.

1040   97.   Lee H, Popodi E, Tang H, Foster PL. Rate and molecular spectrum of

1041   spontaneous mutations in the bacterium escherichia coli as determined by whole-

1042   genome sequencing. Proceedings of the National Academy of Sciences.

1043   2012;109:E2774–83. doi:10.1073/pnas.1210309109.

1044   98.   Paradis E, Schliep K. Ape 5.0: An environment for modern phylogenetics and

1045   evolutionary analyses in R. Bioinformatics. 2018;35:526–8.

1046   doi:10.1093/bioinformatics/bty633.

1047   99.   Paradis E. Pegas: An R package for population genetics with an integrated-

1048   modular approach. Bioinformatics. 2010;26:419–20.

1049   doi:10.1093/bioinformatics/btp696.

1050   100.   Eckshtain-Levi N, Weisberg AJ, Vinatzer BA. The population genetic test

1051   Tajima's D identifies genes encoding pathogen-associated molecular patterns and

1052   other virulence-related genes in Ralstonia solanacearum. Molecular Plant

1053   Pathology. 2018;19:2187–92. doi:10.1111/mpp.12688.

1054   101.   Shen H-M, Chen S-B, Cui Y-B, Xu B, Kassegne K, Abe EM, et al. Whole-

1055   genome sequencing and analysis of Plasmodium falciparum isolates from China-

1056   Myanmar border area. Infectious Diseases of Poverty. 2018;7.

1057   doi:10.1186/s40249-018-0493-5.

102.   Csillery K, François O, Blum MGB. Abc: An R package for approximate bayesian computation (ABC). Methods in Ecology and Evolution. 2012;3:475–9. doi:10.1111/j.2041-210x.2011.00179.x.

103.   Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, et al. SNP-sites: Rapid efficient extraction of SNPs from multi-FASTA alignments. Microbial Genomics. 2016;2. doi:10.1099/mgen.0.000056.

104.   Kamvar ZN, Tabima JF, Grünwald NJ. Poppr: An r package for genetic analysis of populations with clonal, partially clonal, and/or sexual reproduction. PeerJ. 2014;2:e281. doi:10.7717/peerj.281.

105.   Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: A tool set for whole-genome association and population-based linkage analyses. The American Journal of Human Genetics. 2007;81:559–75. doi:10.1086/519795.

106.   Li N, Stephens M. Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. Genetics. 2003;165:2213-33

107.   Marchini J, Howie B. Genotype imputation for genome-wide association studies. Nature Reviews Genetics. 2010;11:499–511. doi:10.1038/nrg2796

108. Sung, W., Ackerman, M. S., Miller, S. F., Doak, T. G., & Lynch, M. Drift-barrier hypothesis and mutation-rate evolution. Proceedings of the National Academy of Sciences 109.45 (2012): 18488-18492. doi.org/10.1073/pnas.1216223109

109. Sivasundar A, Hey J. Population genetics of Caenorhabditis elegans: the paradox of low polymorphism in a widespread species. Genetics. 2003;163:147-57

1084 **Tables**

1085 Table 1. Pan-genome metrics of each Vibrionaceae clades isolated from Poza

1086 Rojas, CCB.

| Group Clade | Number of CCB genomes included in each clade | Pan-genome metrics | | | | Heaps law parameters | |
|---|---|---|---|---|---|---|---|
| | | Core | Flexible | Unique | Total number of genes | Intercept value | Alpha |
| Clade I | 3 | 3617 | 346 | 603 | 4566 | 692.8508 | 1.1293 |
| Clade II | 22 | 1746 | 5770 | 1745 | 9261 | 244.2096 | 0.7913 |
| Clade III | 5 | 2672 | 718 | 324 | 3714 | 658.0634 | 1.6625 |
| Clade IV | 5 | 2055 | 1445 | 180 | 3680 | 2726.7580 | 2.0000 |
| Clade V | 4 | 2853 | 1660 | 1332 | 5845 | 1196.2571 | 1.3109 |
| Clade VI | 3 | 2448 | 3476 | 1028 | 4992 | 3295.5770 | 2.0000 |
| Vibrionaceae all Clades | 47 | 1254 | 14072 | 4795 | 20121 | 2263.7472 | 0.6621 |

1087 The first column shows the Clade ID, next is the number of genomes used for the

1088 analysis regarding each clade, followed by the general metrics of pan-genome, and

1089 last columns show the heaps values obtained. If alpha >1.0 the pan-genome is

1090 considered closed if alpha <1.0 it is considered open.

1091

1092

1093

1094

1095

1096

39

1097    Table 2. Genetic diversity statistics.

| Clade | | Number of individuals | Number of segregating sites | π | θw | Tajima's *D* | P-value of Tajima's *D* |
|---|---|---|---|---|---|---|---|
| Clade I | | 3 | 100971 | 0.0164894 | 0.0163978 | 0 | 0 |
| Clade II | All individuals | 22 | 103197 | 0.01148342 | 0.01106029 | 0.15738106 | 0.8582025 |
| | All individuals in the three larger sub-Clades | 14 | 49946 | 0.00916203 | 0.00613614 | 2.23866585 | 0.02142617 |
| | Sub-clade G | 4 | 13 | 2.54E-06 | 2.77E-06 | -0.84306779 | 0.77323024 |
| | Sub-clade D | 6 | 42 | 5.47E-06 | 7.19E-06 | -1.52560731 | 0.02458297 |
| | Sub-clade A | 4 | 82 | 1.61E-05 | 1.75E-05 | -0.83190864 | 0.8020116 |
| Clade III | | 5 | 40593 | 0.0051088 | 0.0061293 | -1.27467187 | 0.01772241 |
| Clade IV | | 5 | 209 | 2.86E-05 | 3.46E-05 | -1.31696234 | 0 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Clade V | | 4 | 34843 | 0.00398639 | 0.00434715 | -0.87361739 | 0.56601856 |
| Clade VI | | 3 | 204388 | 0.04622002 | 0.04621538 | 0 | 0 |

1098    From left to right are displayed the values for segregation sites, nucleotide diversity

1099    (π) Watterson's theta (θw), Tajima's *D* and Tajima's D *p*-value.  The values were

1100    estimated for all six Clades and Sub-clades with 3 or more individuals.

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118    Table 3. Recombination vs. mutation estimates.

| Group Clade | Recombination vs mutation estimates | |
|---|---|---|
| | rho/theta | $r/m$ |
| Clade I | 0.1036 | 2.7249 |
| Clade II | 0.1171 | 0.5299 |
| Clade III | 0.1498 | 1.1163 |
| Clade IV | 0.1437 | 0.9090 |
| Clade V | 0.0278 | 0.2825 |
| Clade VI | 0.0074 | 0.0052 |
| | | |
| *P. leiognathi* | 0.0064 | 0.2261 |
| *V. anguillarum* | 0.2889 | 4.0014 |
| *V. ordalii* | 0.0667 | 0.5659 |
| *V. parahaemolyticus* | 0.0025 | 0.1246 |

1119    First column shows the names of the CCB Clades and reference strains used for

1120    the calculus. Second and third columns shows the Rho/theta and *r/m* estimates

1121    [43].

1122

1123

1124

1125

1126

42

1127    Table 4. Estimates of effective population sizes obtained through simulations with

1128    Fastsimcoal2 [44, 96].

| Group Clade | Sample size | Median Value | Range | | Environment | Reference |
|---|---|---|---|---|---|---|
| | | | Lower value | Larger value | | |
| Clade I | 3 | 12,822,270 | 10,110,043 | 16,231,765 | Sediment | This work |
| Clade II | | | | | | |
| Sub-clade A | 4 | 55,938 | 34,079 | 392,104 | Sediment | This work |
| Sub-clade D | 6 | 20,849 | 2,795 | 218,603 | Water-Sediment | This work |
| Sub-clade G | 4 | 29,791 | 6,174 | 226,658 | Water-Sediment | This work |
| Clade III | 4 | 15,018,880 | 8,970,283 | 22,432,331 | Water-Sediment | This work |
| Clade IV | 4 | 383,067 | 345,564 | 427,557 | Sediment | This work |
| Clade V | 4 | 9,594,874 | 5,894,074 | 12,914,770 | Sediment | This work |
| Clade VI | 3 | 4,141,870 | 2,582,483 | 10,645,019 | Sediment | This work |
| *H. pylori* | | 39,665,437 | - | - | - | [108] |
| *S. enterica* | | 348,991,354 | - | - | - | [108] |
| *E. coli* | | 179,600,000 | - | - | - | [108] |

| | | | | | | |
|---|---|---|---|---|---|---|
| *H. sapiens* | | 20,348 | - | - | - | [108] |
| *A. thaliana* | | 266,769 | - | - | - | [59] |
| *C. elegans* | | 3,998,701 | - | - | - | [109] |
| *T. brucei* | | 5,332,244 | - | - | - | [108] |

1129 Summary table of effective population sizes of CCB Clades and prokaryotic and

1130 eukaryotic references. First column shows the names of the CCB Clades and

1131 reference strains used for the calculus, second column represents the number of

1132 strains within each group, followed by the median Ne value estimated and the

1133 range. Last two columns display the isolation environment and the reference.

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145  Table 5. GO terms enriched estimated with TopGO [45], regarding the gene

1146  families with signals of positive selection.

| GO.ID | Term | Annotated | Significant | Expected | Fisher test with Bonferroni |
|---|---|---|---|---|---|
| GO:0000902 | cell morphogenesis | 398 | 67 | 26.93 | 0.00020748 |
| GO:0009234 | menaquinone biosynthetic process | 240 | 38 | 16.24 | 0.00150024 |
| GO:0009245 | lipid A biosynthetic process | 240 | 38 | 16.24 | 0.00150024 |
| GO:0008360 | regulation of cell shape | 244 | 37 | 16.51 | 0.0059052 |
| GO:0007156 | homophilic cell adhesion via plasma membrane adhesion molecules | 13 | 7 | 0.88 | 0.0122892 |
| GO:0006304 | DNA modification | 295 | 41 | 19.96 | 0.01596 |
| GO:0009058 | biosynthetic process | 26775 | 1675 | 1811.62 | 0.017556 |

1147  First two columns show the enriched GO IDs and its name, third column the

1148  number of annotated genes, fourth and fifth column the number of significant

1149  genes and the expected, last column shows the significance corrected with

1150  Bonferroni.

1151

1152

1153

1154

1155

45

1156 Table 6. GO terms enriched in the genes found to have an association with the

1157 isolation environment (water or sediment).

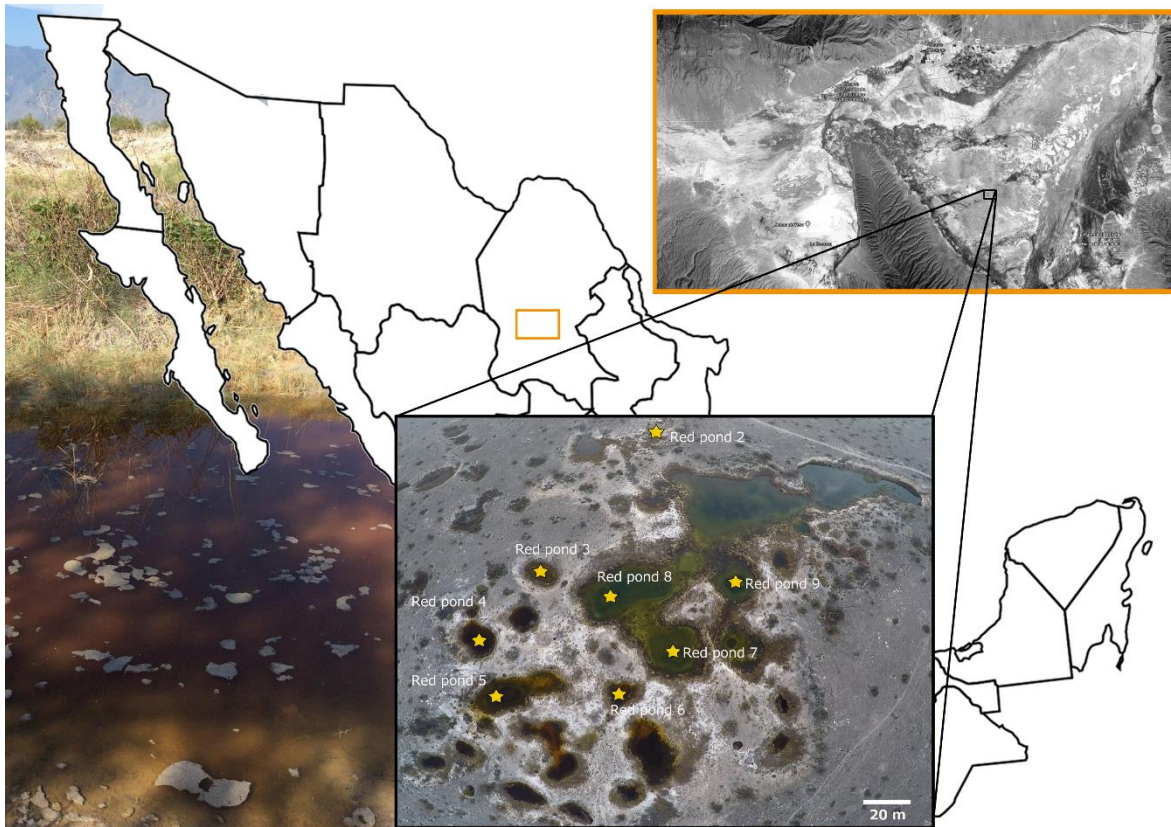| Genes with signals of recombination or selection | GO.ID | Term | Annotated | Significant | Expected | Fisher test with Bonferroni |
|---|---|---|---|---|---|---|
| Recombination | GO:0006066 | alcohol metabolic process | 446 | 8 | 0.55 | 0.000146 |
| Recombination | GO:0006429 | leucyl-tRNA aminoacylation | 41 | 4 | 0.05 | 0.000338 |
| Recombination | GO:0006419 | alanyl-tRNA aminoacylation | 48 | 4 | 0.06 | 0.000643 |
| Recombination | GO:0006265 | DNA topological change | 339 | 6 | 0.42 | 0.006914 |
| Selection | GO:0006814 | sodium ion transport | 685 | 9 | 1.47 | 0.03216 |

1158 First two columns show the enriched GO IDs and its name, with signals of
1159 recombination or selection. Third column the number of annotated genes, fourth
1160 and fifth column the number of significant genes and the expected, last column
1161 shows the significance corrected with Bonferroni.

1162

1163

1164

1165 **Figures.**



1166

1167 **Figure 1.** Study site, Pozas Rojas in Los Hundidos within Cuatro Ciénegas Basin,

1168 Mexico. Sampling sites are signaled in yellow. Cuatro Ciénegas location is also

1169 shown in a map (Pozas Rojas photos were provided by David Jaramillo, a map

1170 showing the location of Cuatro Ciénegas Valley was obtained from Google Earth,
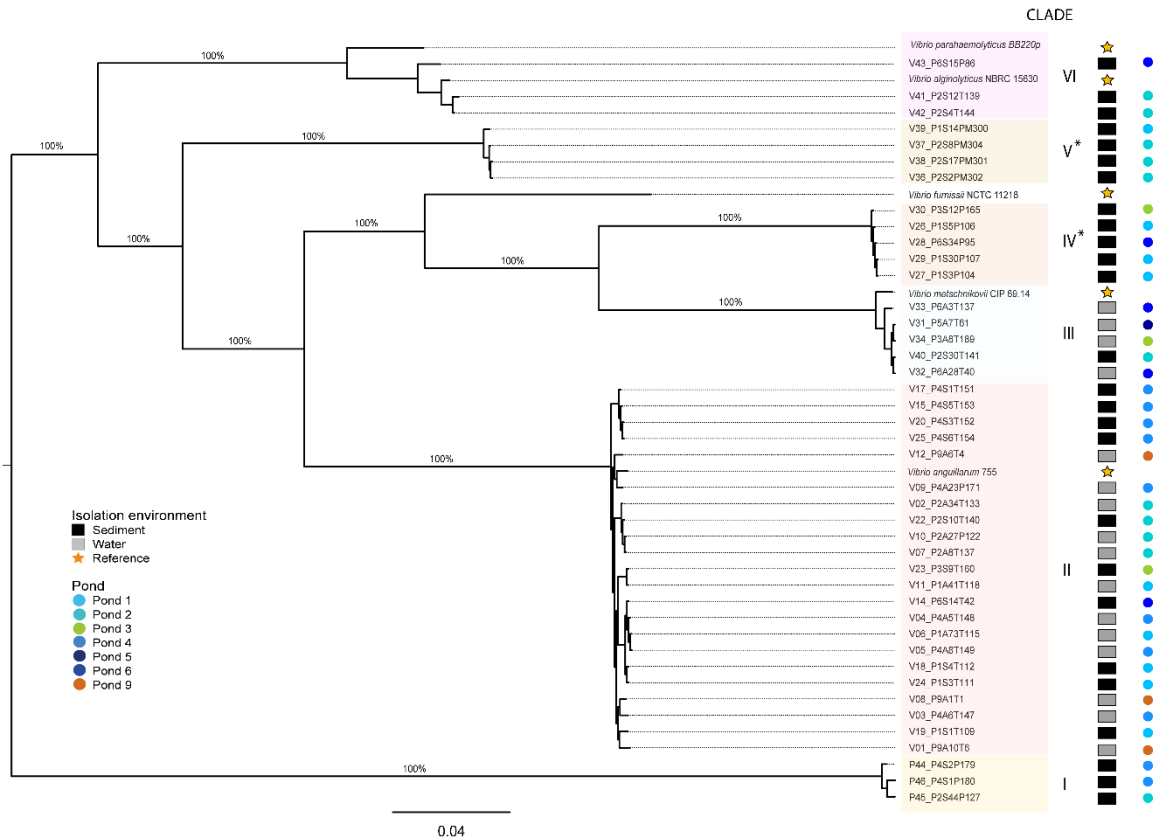
1171 earth.google.com/web/).

1172

1173

1174

1175

1176

1177

1178

1179

1180 **Figure 2.** Core gene phylogeny of the 1,254 orthologs. Maximum-likelihood

1181 phylogenetic reconstruction of core genes, supporting branch values are shown.

1182 Each square represents the isolation environment, water or sediment, while yellow

1183 stars indicate reference strains. Circles indicate isolation pond. Clades are

1184 distinguished with colors. Clades IV and V which are likely to be exclusive to CCB

1185 are highlight with an asterisk.

1186
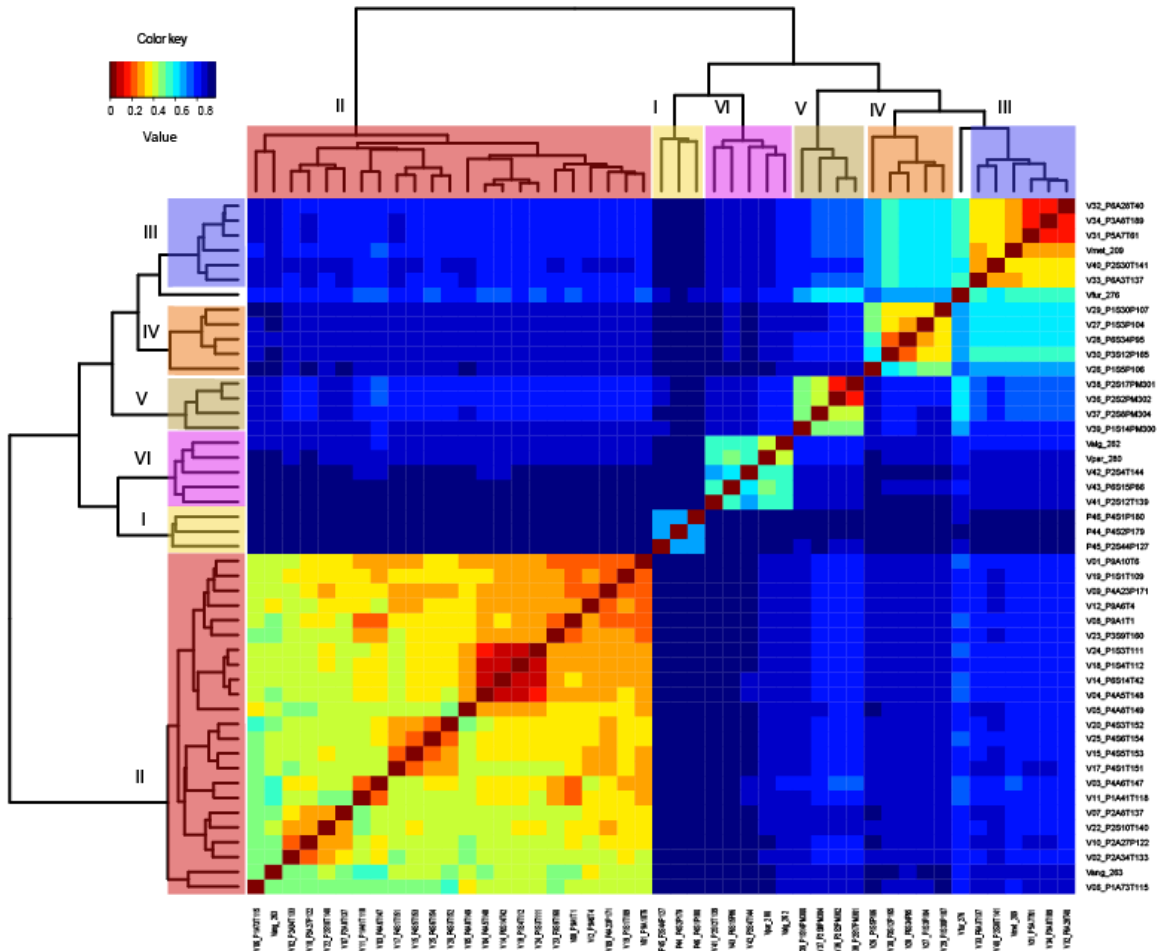
1187

1188

1189

1190

1191

48

1192

**Figure 3.** Patterns of recombination events among isolated strains. Heatmap of the frequency of recombination events among different strains; red colors indicates more recombination events within strains while blue events indicate few recombination events. Distances were estimated with the Jaccard dissimilarity index.
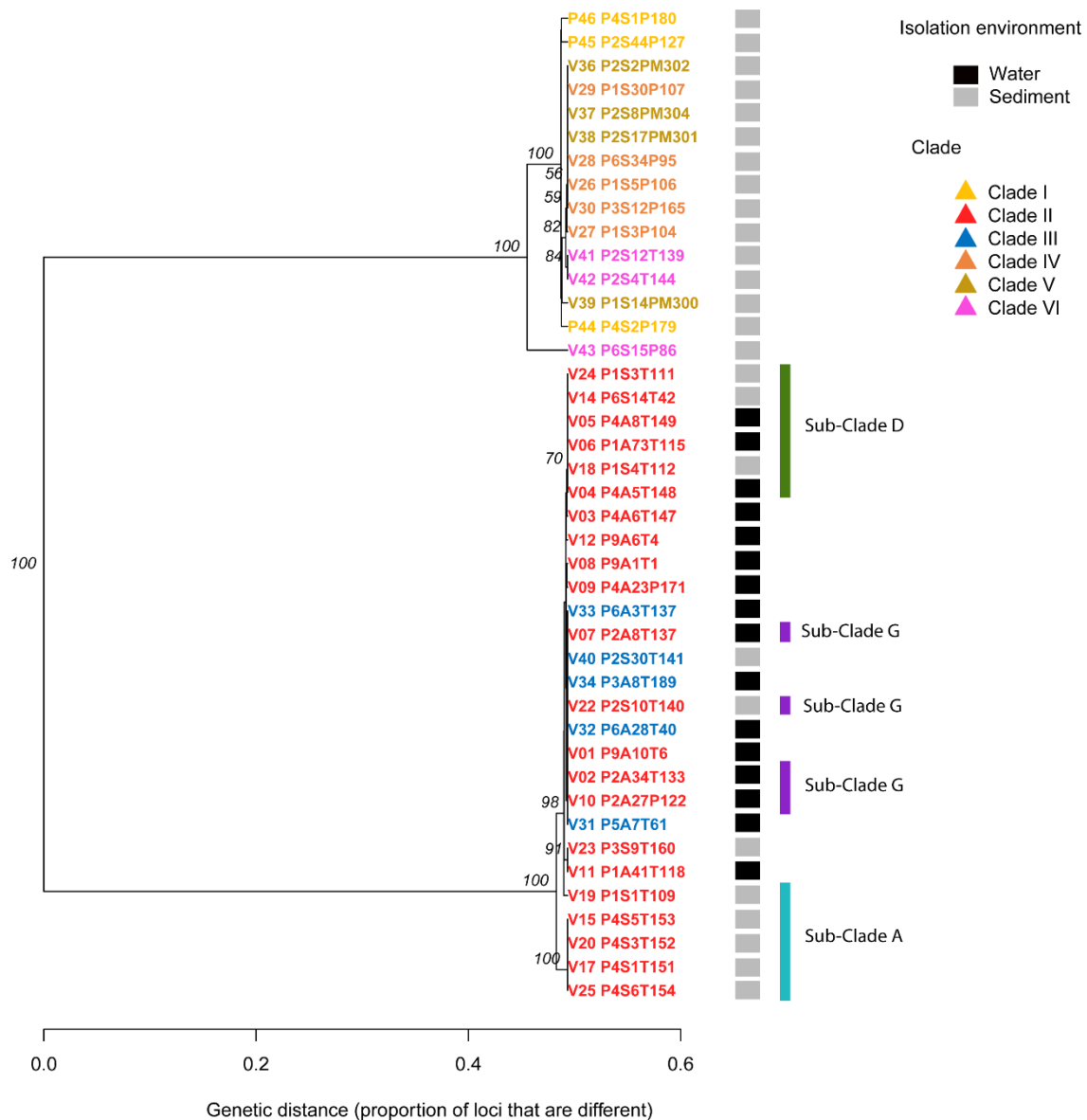
1198

1199

1200

1201

1202

1203

**Figure 4.** UPGMA of the 598 SNPs associated with the isolation environment. Tip colors represent clade membership, for Clade II, Sub-clades are also indicated. Squares represent the isolation environment. Distances were calculated with the bitwise distance function of poppr v2.8.1.