1    **Allopatric plant pathogen population divergence following disease emergence**

2

3    Andreina I. Castillo [1], Isabel Bojanini [1], Hongyu Chen[2,3], Prem P. Kandel[2,4], Leonardo De La

4    Fuente [2], and Rodrigo P.P. Almeida [1]

5    [1] Department of Environmental Science, Policy and Management, University of California,

6    Berkeley, CA, USA.

7    [2] Department of Entomology and Plant Pathology, Auburn University, Auburn, Alabama, USA.

8    [3] Current address: Department of Entomology and Plant Pathology, North Carolina State

9    University, Raleigh, North Carolina, USA.

10    [4] Current address: Department of Plant Pathology and Environmental Microbiology, PennState

11    University, University Park, Pennsylvania, USA.

12

13

14    Corresponding author: rodrigoalmeida@berkeley.edu

15

16    **ABSTRACT**

17    Within the landscape of globally distributed pathogens, populations differentiate via both

18    adaptive and non-adaptive forces. Individual populations are likely to show unique trends of

19    genetic diversity, host-pathogen interaction, and ecological adaptation. In plant pathogens,

20    allopatric divergence may occur particularly rapidly within simplified agricultural monoculture

21    landscapes. As such, the study of plant pathogen populations in monocultures can highlight the

22    distinct evolutionary mechanisms that lead to local genetic differentiation. *Xylella fastidiosa* is a

23    plant pathogen known to infect and damage multiple monocultures worldwide. One subspecies,

24    *Xylella fastidiosa* subsp. *fastidiosa* was first introduced to the USA ~150 years ago, where it was

25    found to infect and cause disease in grapevines (Pierce's disease of grapevines, PD). Here, we

26    studied PD-causing subsp. *fastidiosa* populations, with an emphasis on those found in the USA.

27    Our study shows that following its establishment in the USA, PD-causing strains likely split into

28    populations in the East and West Coast. This diversification has occurred via both changes in

29    gene content (gene gain/loss events) and variations in nucleotide sequence (mutation and

30    recombination). In addition, we reinforce the notion that PD-causing populations within the USA

31    acted as the source for subsequent subsp. *fastidiosa* outbreaks in Europe and Asia.

32

33    **IMPORTANCE**

34    Compared to natural environments, the reduced diversity of monoculture agricultural landscapes

35    can lead bacterial plant pathogens to quickly adapt to local biological and ecological conditions.

36    Because of this, accidental introductions of microbial pathogens into naïve regions represents a

37    significant economic and environmental threat. *Xylella fastidiosa* is a plant pathogen with an

38    expanding host and geographic range due to multiple intra- and inter-continental introductions.

39   *X. fastidiosa* subsp. *fastidiosa*, infects and causes disease in grapevines (Pierce's disease of

40   grapevines; PD). This study focused on PD-causing *X. fastidiosa* populations, particularly those

41   found in the USA but also invasions into Taiwan and Spain. The analysis shows that PD-causing

42   *X. fastidiosa* has diversified via multiple co-occurring evolutionary forces acting at an intra- and

43   inter-population level. This analysis enables a better understating of the mechanisms leading to

44   the local adaptation of *X. fastidiosa*, and how a plant pathogen diverges allopatrically after

45   multiple and sequential introduction events.

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

## INTRODUCTION

The worldwide distribution of microbial plant pathogens is constantly shifting. Global trade and movement of infected plant material enables pathogen introductions from native and endemic aeras to naïve regions (1, 2). Likewise, the intentional introduction of non-native plant species of agronomic and ornamental value to novel environments facilitates the host range expansion of endemic pathogens (3, 4). One crucial factor in the formation of novel plant-pathogen associations is the amount of genetic diversity on which natural selection can act, in other words, the adaptive potential (5). Differences in adaptive potential between host and microbial populations can have a significant role in determining the host and geographic range of a pathogen. For instance, in the case of plant pathogens, higher genetic diversity in effector proteins and virulence genes has a positive effect on host range (6–9). Alternatively, multiple studies have highlighted how reduced genetic diversity in plant hosts can enhance the spread of pathogens within a population (10–12).

Factors that influence genetic diversity, whether via the action of distinct evolutionary mechanisms (13, 14) or as a product of ecological and evolutionary history, affect adaptive potential (15). In plant pathogens, geographical and ecological specialization have been frequently described (16, 17). This is partly explained by plant pathogen differentiation and specialization occurring rapidly within agricultural systems (14, 18, 19). Overall, it is expected that in the absence of gene flow, plant pathogens of agricultural crops will rapidly adapt to local environmental, ecological, and biological conditions (20, 21). Therefore, understanding the mechanisms leading to pathogen adaptation, either to a new crop or environmental condition, has great relevance in developing effective management and control strategies (22). This is particularly pertinent in plant pathogens with a proven capacity to adapt to multiple crops as well

85  as having an expanding geographic and host range. This is the case of the emerging pathogen

86  *Xylella fastidiosa* (23).

87      The bacterial species *X. fastidiosa* has been reported to infect 563 plant species from 82

88  distinct botanical families (23). However, the host range of *X. fastidiosa* varies among and within

89  described subspecies and phylogenetic clades (24). The geographic distribution of the three main

90  *X. fastidiosa* subspecies is also unique, with most of them having experienced one or several

91  dispersal and establishment events at the continental scale. For this reason, efficient

92  identification and tracking of *X. fastidiosa* subspecies has important implications for the

93  development of adequate disease control and mitigation strategies (25, 26). Three *X. fastidiosa*

94  subspecies have an ancestrally allopatric range that has recently expanded: *X. fastidiosa* subsp.

95  *multiplex* is native to temperate and subtropical North America (27, 28), and has been introduced

96  multiple times into Europe (29); *X. fastidiosa* subsp. *pauca* is native to South America (28) but

97  has been recently reported in the Apulian region in Italy and in Costa Rica (30, 31); finally, *X.*

98  *fastidiosa* subsp. *fastidiosa* is native to Central America (32, 33), and was introduced to the USA

99  (24, 34), and subsequently to Europe (35) and Taiwan (36). Other non-monophyletic but

100  proposed subspecies include *X. fastidiosa* subsp. *sandyi*, found in Southern regions of the United

101  States (37, 38) and also introduced into Europe (39); and *X. fastidiosa* subsp. *morus*, only found

102  in regions were subsp. *multiplex* and subsp. *fastidiosa* co-occur (24, 40).

103      The hypothesis that subsp. *fastidiosa* was introduced once to the United States (USA)

104  ~150 years ago leading to the emergence of Pierce's Disease of grapevines (PD) is well

105  supported (24, 33, 41). PD is a grapevine malady that results in significant economic losses to

106  the wine industry in California (42) and the Southeast USA (43). Current knowledge of the

107  evolution of subsp. *fastidiosa* suggests that the ability to infect grapevines was acquired after its

108    introduction to the USA (33). Furthermore, there is evidence that local adaptation to

109    environmental factors has occurred in grape-infecting isolates across a latitudinal gradient in

110    California (34). Finally, available genomic and MLSA-E data suggest that PD-causing isolates in

111    the West and East Coast of the USA are phylogenetically distinct (34, 44)

112         These studies are indicative that after its introduction and establishment in the USA, the

113    subsp. *fastidiosa* clade causing disease in grapevines dispersed to different geographic regions

114    and diversified genetically to adapt to a range of biotic and abiotic conditions. To better

115    understand how *X. fastidiosa* evolved with the emergence of a novel plant disease (PD) and

116    diversified in allopatry in different regions of the USA, we studied populations of the pathogen

117    from the USA and abroad. We evaluated the evolutionary relationship between both USA

118    populations and their relationship with recent introduction events derived from them (i.e.

119    introductions to Spain and Taiwan associated with the emergence of PD in those regions). In

120    addition, we identified the evolutionary mechanisms facilitating population diversification by

121    defining intra-population patterns of gene gain/loss, intra-subspecific recombination, and

122    nucleotide diversity.

123

124    **RESULTS**

125         **PD isolates are split into regional clades within the USA, with Europe and Asia**

126    **introductions originating from these regions.** We arbitrarily split grapevine isolates into 3

127    phylogenetically supported clades, PD-I to -III (Fig.1). These phylogenetically supported clades

128    were also observed in the non-recombinant phylogenetic tree (Fig. S1). PD-I only included

129    isolates from the Southeast USA; PD-II and PD-III were dominated by California isolates, but at

130    the base of those clades there was one isolate from Texas (PD-II) and a sister clade from Georgia

131   (PD-III). No isolates from California grapevines clustered within the PD-I clade. From the data

132   available alone, it is not possible to infer the dispersal history of the non-California isolates in

133   PD-II and PD-III (i.e. basal sister clades or introductions to California). Isolates from Taiwan

134   were phylogenetically placed within the PD-I clade, while those from Spain were nested in the

135   PD-II clade. These represent two distinct introductions, originating from different regions in the

136   USA. Isolates from the same geographic region tended to cluster together within each major

137   clade. For instance, in the PD-I clade, most Georgia isolates from Site1 (i.e. 14B1, 14B4, 14B6,

138   16B2, 15B2, 14B3, and 16B4) and Site2 (i.e. 16M5, 16M6, 16M7, 16M8, and 16M9) clustered

139   together. Isolates from each site formed separate subclades within this group (Fig. 2). Other

140   Georgia isolates from Site1 (i.e. 14B2, 14B5, 14B7, 16B1, 16B3, 16B5, and 16B6) were more

141   closely related to those from Florida and North Carolina. In a similar manner, isolates from the

142   West Coast (i.e. California) tended to group geographically. Specifically, isolates obtained from

143   Southern California (i.e. Je81, Je104, Je112, Je110, etc.) were ancestral to those from Northern

144   California (i.e. Hopland, Stag Leap, Conn-Creek, CV17-3, Je65, Je73, etc.) in the PD-III clade.

145       A total of 141 different haplotypes named using roman numerals (I-CXIV) (Fig. 1a) were

146   found in the PD-causing core genome alignment. Haplotypes were structured by geographic

147   location and largely matched the evolutionary relationships observed in phylogenetic analyses

148   (Fig. 1b,2). Overall, haplotypes were grouped similarly to the phylogenetic analyses. Isolates

149   originating from the West and East Coast were split by 979 mutations. California had the largest

150   number of haplotypes (106) as well as haplotypes with the highest frequency: XXVIII (7), XLIV

151   (6), XXXIV (5), LVIII (4), LXVI (4), LXXIII (3), and XCIV (3). On the other hand, Southeast

152   USA haplotypes (31) were generally found in low frequency (i.e. one or two isolates). In

153   addition, Southeast isolates in PD-III formed a distinct group separated from the California group

154 by 243 mutations. Likewise, GB514 (Texas, PD-II) was closely connected to California isolates,

155 from which it differentiated by 159 mutations. Isolates originating from recent introduction

156 events (i.e. Spain and Taiwan) had unique haplotypes. Spanish associated haplotypes were linked

157 to a haplotype originating from California (PD-II) and were differentiated by 61 mutations.

158 Similarly, the Taiwan haplotypes were closely linked to the haplotype group originating from

159 Southeast USA (PD-I) and differentiated by 13 mutations.

160 **Gene gain and loss events occur following subsp. *fastidiosa* introduction events.**

161 Estimated rates of gene gain/loss were highest in branches leading to the introduction of subsp.

162 *fastidiosa* from Central America. Furthermore, a total of 35 core genes were absent in the PD-

163 causing population compared to the Costa Rican isolates, while 49 core genes were present in the

164 PD-causing population but absent in the Costa Rican isolates. In addition, gene gain/loss events

165 also occurred within the USA populations. In California (Fig. S2a) gene gain/loss rates were

166 highest in the branches leading to each cluster than within clusters, but PD-III had higher gene

167 gain/loss rates compared to PD-II. Likewise, two clades were observed within the Southeast

168 USA population (Fig. S2b). The first clade was formed by isolates 16M2, 16M3, XF51_CCPM1

169 (from Georgia, clustering with PD-III), and GB514 (from Texas, clustering with PD-II); and the

170 second by the remaining Southeast USA isolates (PD-I).

171 Some unique genes were identified through estimating gene gain/loss rates within each

172 population. We found that, when considering geographical origins of isolates alone, gene

173 presence/absence was similar in PD-II and PD-III isolates regardless of geographical origin (Fig.

174 3a). In the case of PD-I, PD-II, and PD-III isolates from Southeast USA, three genes were

175 uniquely found in PD-I and nine in PD-III (Table S3). When gene gain/loss was compared

176 between PD-II and PD-III isolates from California and PD-I, three genes coding for hypothetical

177    proteins were found in PD-II and PD-III isolates from California but absent in PD-I. In addition,

178    two genes were absent in isolates from Spain but present in PD-II and PD-III isolates from

179    California. On the other hand, two genes were found in PD-I but absent in PD-II and PD-III

180    isolates from California (Fig. 3b). In addition, five genes were absent in isolates from Taiwan,

181    which was considered as the descendant population of Southeast USA (Table 1).

182        These unique genes were annotated using eggNOG-mapper and searched in the GenBank

183    and Pfam databases, using both BLAST and interproscan5 (Table 1, Table S4). Two hypothetical

184    proteins and a gene coding for the HTH-type transcriptional regulator (*prtR*) were found for PD-I

185    (Table S3); while nine were hypothetical proteins, the protein coding genes *traC_2* (DNA

186    primase), and *higB_2* (endoribonuclease) were found for the PD-III Southeast USA isolates. Two

187    of the three genes found in PD-II and PD-III isolates from California but absent in PD-I coded

188    for hypothetical proteins and one coded for an alpha/beta fold hydrolase. For the two genes

189    absent in Spain, one of them was listed as glutamate 5-kinase, and another had a conserved

190    LacZ, Beta-galactosidase/beta-glucuronidase domain. For the two genes found in PD-I but

191    absent in PD-II and PD-III isolates from California, one was annotated as a hypothetical protein

192    and the other one as a phage head morphogenesis protein. For the five genes absent in isolates

193    from Taiwan, two were annotated as site-specific DNA-methyltransferase; another two were

194    annotated as peptidoglycan DD-metalloendopeptidase family protein and hypothetical protein,

195    respectively; the last one could be a pseudo gene with unknown function.

196        **Intra-subspecific recombination events are pervasive in both the West and East**

197    **Coast.** Intra-subspecific recombination was pervasive in both populations (Fig. 4 and Fig. S3-4).

198    The r/m estimate (recombination to mutation rates) for the California/Spain core genome

199    alignment was 3.29, while the same estimate for Southeast USA/Taiwan core genome alignment

200    was 5.65. In the Southeast USA (Fig. 4b), recombination events were more frequently observed

201    in isolates from the PD-II/PD-III group (recipient) than in isolates from the PD-I group (donor).

202    Within the PD-II/PD-III group, the Texas isolate GB514 (PD-II) was the most frequent

203    recombinant recipient. Donor sequences for the Texas isolate originated from both PD-I and

204    from an 'unknown' donor (representing genetic variability present in the population but not

205    characterized in the original sampling). A total of 188 core genes were entirely contained within

206    recombinant regions in the Southeast USA population; out of this group, 101 genes were

207    classified as hypothetical proteins. The remaining recombinant core genes belonged to a variety

208    of functions (Table S5). These functions were grouped by their COG class resulting in 12 genes

209    belonging to the 'Cellular Processes and Signaling' class, 5 genes associated with the

210    'Information Storage and Processing' class, 41 genes from the 'Metabolism' class, and 7 genes

211    belonging to two or more functional classes ('Multiple Categories'). Based on gene annotation,

212    some CDs functions are related to virulence and/or host adaptation. These include vitamin $B_{12}$

213    import (*btuD*), ferric uptake regulation protein (*fur*), response regulator (*gacA*), virulence protein

214    (PD_1332 in Temecula1 assembly AE009442.1, COG0346), polygalacturonase (*pglA*), export

215    protein (*secB*) and ABC transporter (*uup*).

216        Likewise, sequence exchange occurred between isolates from the PD-III and the PD-II

217    clusters in California. Recombination events were observed among isolates from the same

218    geographic regions (Fig. 4a). Specifically, recombination was frequent between sequences

219    originating from the Temecula Valley in Southern California (Fig. S3). Sequences in both groups

220    acted as donors and recipients. In addition, Northern California isolates were recipients of

221    recombinant segments from Southern California. This group was also a recipient of 'unknown'

222    sequence fragments. A total of 180 core genes were exclusively contained within these

223    recombinant regions (Table S5). Eighty-five genes were described as hypothetical proteins. The

224    remaining genes were classified by their COG as: 'Cellular Processes and Signaling' (19 genes),

225    'Information Storage and Processing' (6 genes), 'Metabolism' (38 genes), and 'Multiple

226    Categories' (6 genes). From these genes, those with annotated function related to host

227    adaptation/virulence include: biofilm growth-associated repressor (*bigR*), periplasmic serine

228    endoprotease (*degP*) (*htrA* in Temecula1 assembly AE009442), putative TonB-dependent

229    receptor (*phuR* in Temecula1 assembly AE009442, COG1629), virulence protein (PD_1332 in

230    Temecula1 assembly AE009442.1, COG0346), sec-independent translocase protein (*tatA-D*),

231    and PhoH-like protein (*ybeZ*).

232        Based on the used genome annotations, a total of 13 recombinant genes were shared in

233    both populations. These genes were: *glk_1* and *glk_2* (glucokinases), *glmM_2* (a

234    phosphoglucosamine mutase), *glmS_1* and *glmS_2* (glutamine--fructose-6-phosphate

235    aminotransferases [isomerizing]), *grpE* (a GrpE protein), *grxD* (a glutaredoxin 4), *gshB* (a

236    glutathione synthetase), *gtaB* (a UTP--glucose-1-phosphate uridylyltransferase), *pepQ* (a Xaa-

237    Pro dipeptidase), *petA* (an Ubiquinol-cytochrome c reductase iron-sulfur subunit), *petC* (an

238    ammonia monooxygenase gamma subunit), an unnamed PKHD-type hydroxylase (COG3128),

239    and a unnamed Virulence protein (COG0346).

240        **Grapevine-infecting populations in the East and West USA are largely genetically**

241    **isolated.** Nucleotide diversity ($\pi$) varied within and among populations (Table 2). Overall,

242    nucleotide diversity was higher within the Southeast USA (947 SNPs, $\pi=1.36\times10^{-05}$) compared

243    to California (458 SNPs, $\pi=3.22\times10^{-06}$). When compared to their corresponding source

244    populations, nucleotide diversity was lower within Spain (2 SNPs, $\pi=1.38\times10^{-07}$) and Taiwan (6

245    SNPs, $\pi=4.15\times10^{-07}$). When diversity in phylogenetically distinct clusters was evaluated, PD-I

Castillo  11

246 (93 SNPs, $\pi$=7.58x10e$^{-07}$) and PD-II (114 SNPs, $\pi$=9.65x10e$^{-07}$) had lower nucleotide diversity

247 than PD-III (509 SNPs, $\pi$=3.25x10e$^{-06}$).

248        The frequency of polymorphism present in the population in regard to expectations under

249 neutrality was calculated using a Tajima's D. Briefly, negative Tajima's D values indicate an

250 excess of rare polymorphisms than expected under neutrality, which can be caused by a selective

251 sweep or a recent population expansion. Positive Tajima's D values indicate excess of

252 intermediate frequency polymorphism than expected under neutrality, which could suggest

253 balancing selection or a recent population contraction. Tajima's D in California and the

254 Southeast USA was negative (Table 2); however, the magnitude of the statistic in California was

255 roughly twice that of the Southeast USA (-1.448 and -0.658, respectively). Due to the reduced

256 sample size, it was not possible to estimate Tajima's D in Spain or Taiwan. When populations

257 were divided phylogenetically, PD-I isolates had a lower Tajima's D (-2.060) compared to PD-II

258 (-1.781) and PD-III (-1.743). On the other hand, Watterson's $\theta$ estimates the population mutation

259 rate from the observed nucleotide diversity. This estimator decreases with increased sample size

260 or with recombination rate. Watterson's $\theta$ estimated a higher mutation rate in the Southeast USA

261 ($\theta$=1.64x10e$^{-05}$) compared to California ($\theta$=5.75x10e$^{-06}$). When populations were divided based

262 on phylogeny, mutation rate was higher in PD-III ($\theta$=6.72x10e$^{-06}$) than PD-I ($\theta$=1.64x10e$^{-06}$) or

263 PD-II ($\theta$=1.87x10e$^{-06}$).

264        In addition, Fst values were used to measure population differentiation across geographic

265 and phylogenetic groups. Briefly, Fst values compare the amount of genetic variability within

266 and between populations, values of 1 indicate complete population structuring while values of 0

267 indicate complete panmixia. Pairwise Fst values (Table S6) for California vs. Southeast USA

268 (Fst = 0.814) and California vs. Taiwan (Fst = 0.964) were higher than California vs. Spain (Fst

269    = 0.566). This was also the case for comparisons involving Southeast USA vs. Spain (Fst =

270    0.847) and Southeast USA vs. Taiwan (Fst = 0.114). Taiwan vs. Spain also showed strong

271    differentiation (Fst = 0.994). Once populations were divided phylogenetically, PD-I was more

272    differentiated from PD-II (Fst=0.987) and PD-III (Fst=0.960), than PD-II and PD-III from each

273    other (Fst=0.541).

274           An MKT was used to estimate the rate of synonymous and non-synonymous

275    polymorphism vs. the rate of synonymous and non-synonymous fixed differences across

276    geographic populations and phylogenetic groups. Under neutrality, it is expected that both rates

277    will be the same (NI=1). Therefore, departures of neutrality (NI≠1) will indicate either the action

278    of balancing selection (e.g. maintenance of population polymorphisms; NI>1) or the action of

279    positive selection (e.g. accumulation of fixed differences between populations; NI<1). The

280    Neutrality Index (NI) was larger than 1 in all comparisons except for Spain vs. Taiwan. NI was

281    significant only for California vs. Taiwan (p-value=$9.87 \times 10^{-05}$) (Table S6). Many

282    polymorphisms were observed in Southeast USA and California, while few were observed

283    within Spain or Taiwan. The largest number of fixed differences were observed for Taiwan vs.

284    California. When populations were divided phylogenetically, the NI values were larger than 1

285    only in comparison between PD-I with PD-II and PD-III. In this instance, the only significant NI

286    was observed for PD-I vs. PD-III (p-value=$6.26 \times 10^{-05}$). The number of polymorphisms was

287    larger in PD-III compared to PD-I and PD-II. The number of fixed differences were similar

288    between PD-I vs. PD-II and PD-III, but smaller in PD-II vs. PD-III.

289           Selective sweep signatures were pervasive in both California and Southeast USA (Fig.

290    5a), thought the magnitude of the sweep was larger in California. Alternatively, CLR peaks were

291    smaller and scattered in Spain and Taiwan. When the populations were split phylogenetically,

Castillo  13

292    CLR peaks were more numerous and prominent in PD-III, followed by PD-II, and finally PD-I

293    (Fig. 5b). Regardless if the populations were subdivided geographically or phylogenetically,

294    some CLR peaks co-located across populations, while others were group specific.

295

296    **DISCUSSION**

297         Our analyses show that after its introduction from Central America (33, 41), PD-causing

298    subsp. *fastidiosa* split into two populations: one in the East Coast (31 haplotypes) and one in the

299    West Coast (106 haplotypes). Apart from PD-II/PD-III isolates from the Southeast USA, each

300    population formed a sister monophyletic clade with long basal branch lengths. This indicates that

301    the populations split shortly after introduction to the USA. Moreover, isolates from the same

302    location clustered together, suggesting stronger sequence similarity within than between

303    locations. With the current information available, it is not possible to know if the clustering of

304    PD-II/PD-III isolates from the Southeast USA with the California clades instead of Southeast

305    USA (PD-I) reflects a recent introduction to California or if there is a higher diversity within

306    Southeast USA isolates than currently represented. Alternatively, it is feasible the East and West

307    Coast populations originated via independent introduction events. Previous studies have pointed

308    out the large genetic diversity of subsp. *fastidiosa* within Central America (33) and the

309    importation of plant material from this region into the USA (67). Our data do not exclude the

310    possibility that additional subsp. *fastidiosa* strains circulate within Central America and could

311    have been introduced to the USA in relatively simultaneous events. This is a hypothesis that

312    should be evaluated as additional whole genomic data from both native and introduced

313    populations of subsp. *fastidiosa* becomes available. However, previously published MLST data

314    (41, 67) and results based on whole genome sequence analysis (monophyly of the PD-causing

Castillo   14

315    population, age and diversity of PD-causing clades, and their evolutionary relationship with the

316    native subsp. *fastidiosa* population) are indicative of a single introduction event.

317         Pathogen introductions into Spain and Taiwan were closely related to isolates from

318    California and the Southeast USA, respectively. Though closely related to their source

319    populations, both Spain and Taiwan had their unique core haplotypes which could be indicative

320    of early local adaptation. Nonetheless, we cannot discard the possibility that differences in

321    unique core haplotypes might also be the result of a founder effect. Small sample size in both

322    populations does not allow to test between these two possibilities; however, this should be

323    addressed once additional genomic data becomes available.

324         **Gene gain/loss events are common between and within populations.** Bacterial gene

325    content is in constant flux (68); in bacteria, evolution via gene gain and loss often precedes

326    evolution at the sequence level (i.e. nucleotide substitutions and indels) (69). Therefore

327    variations in gene content can act as a source for adaptive differentiation (70). Gene gain and

328    loss rates were highest following the introduction to the USA (e.g. 35 genes gained and 49 loss

329    vs. 4 genes gained and 5 loss between the East Coast and Taiwan); however, gene content

330    changes were also detected within each geographic population. The higher number of gene

331    gain/loss events observed in basal tree branches can be explained by a founder event. However,

332    they could also be the result of accumulated gene gain/loss events over longer evolutionary time.

333    It is likely that both factors contribute to gene gain/loss between the native and ancestral subsp.

334    *fastidiosa* populations. The highest intra-population gene gain/loss rates were localized in

335    branches following clade splits. Within California, intra-population splits were associated with

336    locations along a latitudinal gradient (PD-II in Southern California vs. PD-III in Southern and

337    Northern California). In other organisms, selection driven gene gain/loss has been described in

Castillo   15

338   genes involved in environmental interactions (69, 71, 72). Likewise, previous studies have found

339   evidence of local adaptation to environmental conditions within California (34). Thus, it is

340   possible that changes in gene content might be adaptive to the local environment. This is further

341   supported by PD-III, which encompasses a larger latitudinal gradient, having four times higher

342   gene gain/loss rates compared to PD-II.

343        Alternatively, while gene gain/loss rates were higher in PD-II/PD-III Southeast samples

344   compared to PD-I, the difference was not as pronounced as that seen in California. Sampling of

345   PD-causing isolates has been more extensive in California; therefore, detecting environmentally

346   linked gene gain/loss might require further sampling in the Southeast USA. Based on the current

347   annotation, it is difficult to interpret the possible benefit or disadvantage of unique genes found

348   in specific *X. fastidiosa* populations. Functional analysis of these genes will be needed to

349   understand their biological role. Still, a small number of genes involved in transcription

350   regulation (*prtR* and *higB_2*) and DNA replication (*traC_2*) were exclusive to PD-I and PD-III

351   Southeast USA isolates. These functions are linked to changes in bacterial transcription and

352   replication in response to environmental cues (73, 74).

353        However, it should be noted that gene gain/loss events can also be a product of non-

354   adaptive evolution. In bacteria, genetic drift promotes genome reduction and neutral gene losses

355   are favored by small population size (75, 76). In addition, homologous recombination facilitates

356   core genome homogenization but might not affect accessory genes, leading to gene content

357   divergence and pangenome expansion (69). As such, these gene gain/loss events might not be

358   linked to the adaptive potential of each population. Likewise, this could also be the case of more

359   recent introduction events and smaller population sizes (i.e. Spain and Taiwan populations).

360

361      **Unequal recombination frequencies drive inter- and intra-population**

362 **differentiation.** r/m estimates showed that recombination contributes more than mutation to

363 genetic diversity. The r/m values for California/Spain (r/m=3.29) and Southeast USA/Taiwan

364 (r/m=5.65) were higher than previous reports on subsp. *fastidiosa* (r/m = 2.074, (33)). However,

365 both values were lower than reports focused specifically for a California population (r/m=6.797,

366 (34)). Location-specific core genomes analyses can detect nucleotide changes uniquely to a

367 geographic region. Therefore, the high r/m found here is likely due to location-specific SNPs.

368      The number of genes located within intra-subspecific recombination was similar across

369 functional classes showing that there were no specific gene functions more prone to

370 recombination. These results are like those found in a previous analysis (33). On the other hand,

371 the frequency of recombination varied among phylogenetic clusters. PD-II/PD-III Southeast

372 isolates were recipients to sequence fragments from both PD-I and an 'unknown' group.

373 Similarly, recombination occurred among geographically close isolates from the PD-II and PD-

374 III clusters in California. These results show that genetic exchange is actively occurring within

375 the West and East Coast. Variations in recombination frequency across isolates have been

376 reported in native subsp. *fastidiosa* populations (33). Furthermore, recombinant genotypes form

377 distinct phylogenetic groups in subsp. *multiplex* (77). Also, *in vitro* analyses have shown that the

378 natural competency in both subsp. *fastidiosa* and subsp. *multiplex* is strain dependent (78, 79).

379 Taken together this shows that intra- and inter-subspecific recombination does not equally affect

380 all strains and that different gene functions, at least within the core genome, are not differentially

381 prone to recombination.

382      Recombination events also contribute to the differentiation between the East and West

383 Coast, as well as between PD-I, PD-II, and PD-III. Previous studies have shown allele exchange

384    between co-occurring subsp. *multiplex* and subsp. *fastidiosa* isolates in the Southeast USA, but

385    not in California (40). Therefore, the presence of multiple *X. fastidiosa* subspecies within the

386    same geographic regions can enable divergence of recombinant prone isolates or clades.

387    Moreover, highly recombinant clades also experienced higher gene gain/loss in the East and

388    West Coast. Homologous recombination can aid in maintaining core genome cohesiveness while

389    allowing extensive gene gain/loss in the accessory genome (69) and variations in gene content

390    can enable ecological divergence (80). Therefore, intra-subspecific recombination can act as a

391    source of differentiation in PD-causing isolates, not only by mediating allelic exchange but also

392    by facilitating gene gain/loss.

393         From the genes found to recombine in the Southeast USA and Californian populations

394    with putative function as host adaptation and/or virulence, most have been already identified as

395    recombinants among *X. fastidiosa* populations (79). Genes with the same annotation found in

396    both studies include *btuD, secB, uup, tatD,* and *ybeZ.* In other cases the identified genes were not

397    exactly the same, but genes with similar functions were found in both studies, including genes

398    related to iron acquisition (*fur* in the current study), biofilm-associated-repressor (*bigR* in the

399    current study) (81–83) and sulfide sensor (84), other members of the *sec* pathway (*tatA-D*) (85,

400    86), and other serine proteases (*degP* here) (79, 85, 86). Interestingly the vitamin $B_{12}$ transporter

401    BtuD was the single annotated gene with highest recombination inter- and intra-subspecific

402    identified in a previous study (79), and has been described in other bacteria as regulating gene

403    expression, abundance of microorganisms and virulence (87, 88), although no functionality has

404    been attributed yet to *X. fastidiosa*. Genes like *fur* and *gacA* have been identified as

405    transcriptionally regulated by calcium (89), an abundant element inside xylem vessels. Other

406    genes like the putative TonB-dependent receptor (*phuR* in the Temecula1 assembly AE009442,

407    COG1629), are involved in twitching motility and biofilm formation (90); and PhoH-like protein

408    (*ybeZ*), is putatively linked to detection and response to changes of phosphate concentration (91).

409    **West and East Coast populations show unique trends of genetic diversity and**

410    **mutation rate.** At a first glance, isolates originating from the Southeast USA population were

411    more genetically diverse than those originating from California. However, this trend was less

412    clear when populations were assigned phylogenetically. PD-II (California + 1 Texas isolate) had

413    slightly higher than PD-I (exclusively Southeast USA), and PD-III (California + 3 Georgia

414    isolates) had higher genetic diversity than either PD-I or PD-II.

415    The negative Tajima's D values indicate an excess of rare polymorphisms, which can be

416    caused by a selective sweep or a recent population expansion. In the case of subsp. *fastidiosa*, a

417    population expansion could have occurred following a founder effect. This result, in addition to

418    previously published data (24, 33, 67), supports the hypothesis that subsp. *fastidiosa* was

419    introduced to the USA. Furthermore, they show that limitation on genetic diversity caused by a

420    founder effect can be long lasting. Tajima's D values were markedly reduced in PD-I compared

421    to the geographic Southeast USA population (PD-I+PD-II(Texas)/PD-III(Georgia)). This is

422    indicative that there is more than one phylogenetic cluster circulating in the East Coast.

423    Similarly, Tajima's D was smaller in PD-II and PD-III compared to California, further

424    supporting the idea of ongoing latitudinal distinction within the West Coast.

425    Watterson's θ estimates were also affected by grouping criteria. In the case of Southeast

426    USA compared to PD-I, the Watterson estimator remained roughly unchanged suggesting that

427    mutation rate in the region is captured by current sampling. Watterson's θ was larger in

428    California compared to either PD-II or PD-III, and lower in PD-III compared to PD-II. The

429    values were comparable to previous reports in California (34). This could be indicative that

430  mutation rate within the West Coast is, to a certain point, location dependent and that mutation

431  itself contributes less to population differentiation than other evolutionary forces.

432       **PD-causing strains have differentiated phylogenetically and geographically.** The Fst

433  values for different groups of PD-causing isolates were higher than those reported for other

434  global bacterial pathogens(92). It is possible that these values reflect rapid differentiation of PD-

435  causing populations. Pairwise Fst values between PD-I (Southeast only) vs. PD-II (California + 1

436  Texas) and PD-III (California + 3 Georgia) were higher than between Southeast USA and

437  California. These results further support the phylogenetic and geographic separation of the East

438  and West Coast, and the more recent differentiation within California. How much this

439  differentiation can be linked to the Southeast USA PD-II/PD-III group, needs to be further

440  analyzed. Our Fst analyses indicate a complex phylogeographic history between USA

441  populations, yet, the effects of sample size in these calculations should not be ignored. For

442  example, recently introduced populations (e.g. Spain and Taiwan) showed even higher

443  population differentiation than comparisons involving their source populations. Whether this

444  suggests higher differentiation as a product of a founder effect remains to be determined.

445       In general, genetic diversity has a high impact on adaptive potential (93); however, some

446  genetic variants might be considered neutral and can be estimated based on the number of

447  synonymous polymorphisms (94). Variables associated with local adaptation are linked to non-

448  synonymous polymorphisms. There were more non-synonymous than synonymous

449  polymorphism in both the East and West Coast. This suggests that, though the number of

450  polymorphisms might be limited due to a recent introduction event, each population maintains a

451  certain level of genetic variation (as evidenced by NI > 1) which could be a source for local

452  adaptation (95).

453    When populations were divided according to their phylogenetic relationships, a

454    significant NI > 1 was only observed between PD-I and PD-III. Polymorphism largely

455    accumulated in PD-III compared to PD-I. However, the number of fixed differences was

456    comparable between PD-I vs. PD-II and PD-III. This shows that a significant number of intra-

457    clade polymorphisms in PD-III have not yet been fixed. Instead, fixed differences seem to mostly

458    reside between the PD-I compared to PD-II and PD-III. This further supports the idea that East

459    and West Coast populations split early following introduction to the USA, with local population

460    differentiation within a latitudinal gradient in the West Coast.

461    **Selective sweeps have occurred following the introduction of *X. fastidiosa* to the**

462    **USA.** Many CLR peaks were co-localized in the same region while others were group exclusive.

463    The localized nature of CLR peaks in the core genome alignment suggests that selective sweeps

464    can only be detected on certain genes. The location and intensity of selective sweeps are the

465    product of evolutionary and ecological variables. A founder effect can result in reduced selection

466    strength, but it might not affect recombination potential, particularly in *X. fastidiosa* (96).

467    Therefore, the CLR patterns observed here likely reflect genes undergoing strong selection,

468    either following subsp. *fastidiosa* introduction from Central America (co-localized CLR peaks)

469    or via selective pressures associated to a specific environment (group specific CLR peaks).

470    Strong CLR signals in both PD-II and PD-III are indicative that selective sweeps have been more

471    prevalent within the West Coast. Some genes located in CLR peak include: outer membrane

472    protein assembly factors (*BamA-B*), a beta-barrel assembly-enhancing protease (*bepA_4*), a

473    ubiquinol cytochrome C oxidoreductase (*fbcH*), a glycine cleavage system transcriptional

474    repressor (*gcvR*), a glutamine--fructose-6-phosphate aminotransferase (*glmS_2*), a

475    proton/glutamate-aspartate symporter (*gltP*), and a sensor histidine kinase (*rcsC*). Branch-site

476 analyzes aimed to detect signals of positive selection should be performed to further evaluate

477 these results.

478

479 **CONCLUSIONS**

480 We identified a series of evolutionary mechanisms that led to the diversification of PD-causing

481 subsp. *fastidiosa* populations. Diversification has occurred in core genome sequences via

482 mutation and recombination, and in gene content via gain/loss events. These differences have the

483 potential of facilitating local adaptation to environmental conditions, and in the absence of gene

484 flow, lead to pathogen specialization. The host range and geographic distribution of *X. fastidiosa*

485 is expanding and each new introduction can result in significant economic and ecological

486 damage. Understanding the mechanisms and speed of local adaptation in *X. fastidiosa* is

487 important to manage emerging *X. fastidiosa* diseases and hopefully limit the number of novel

488 epidemics.

489

490 **MATERIALS AND METHODS**

491 **Sampling, culturing, and isolation.** The following study encompasses 175 *X. fastidiosa*

492 subsp. *fastidiosa* isolates obtained from infected PD-symptomatic grapevines from diverse

493 geographic regions. The number of isolates from each region were: California (N=140),

494 Southeast USA (N=31), Spain (N=2), and Taiwan (N=2). In addition, three non-grapevine-

495 infecting *X. fastidiosa* subsp. *fastidiosa* isolates from Costa Rica were used as an outgroup (33).

496 New subsp. *fastidiosa* isolates were obtained from infected grapevines in the Southeast USA

497 during 2014-2016; these isolates were cultured from symptomatic leaves as previously described

498 (44). Colonies growing after ~1-2 weeks under 28°C incubation were re-streaked, cloned, and

499 had identity confirmed with *X. fastidiosa* specific PCR primer sets (45). Isolates were obtained

500 from different grapevine varieties. Specifically, the varieties found in Site1 were: Merlot (N=5,

501 years 2014-2016), Mourvedre (N=1, year 2014), Cabernet Sauvignon (N=1, year 2014),

502 Chardonnay (N=5, years 2014-2016), Viognier (N=2, years 2014-2015), Sangiovese (N=1, year

503 2014), and Touriga (N=1, year 2014). The varieties found in Site2 were: Montaluce (N=1, year

504 2015), Merlot (N=3, year 2016), Pinot grigio (N=3, year 2016), and Vidal (N=3, year 2016).

505 Except for Site1 (N=16) and Site2 (N=8), all data included in the following study has been

506 previously made publicly available. Detailed metadata on each assembly has been compiled in

507 Table S1; assembly statistics for new whole genome sequences are provided in Table S2.

508 **Sequencing, assembly, and annotation of *X. fastidiosa* subsp. *fastidiosa* isolates.** All

509 isolates were sequenced using Illumina HiSeq2000. Samples were sequenced at the University of

510 California, Berkeley Vincent J. Coates Genomics Sequencing Laboratory (California Institute for

511 Quantitative Biosciences; QB3), and the Center for Genomic Sciences, Allegheny Singer

512 Research Institute, Pittsburgh, PA. All raw reads and information regarding each newly

513 sequenced strain can be accessed under the NCBI BioProject accession PRJNA655351. The

514 quality of raw paired FASTQ reads was evaluated using FastQC (46) and visualized using

515 MultiQC (47). Low quality reads and adapter sequences were removed from all paired raw reads

516 using seqtk v1.2 (https://github.com/lh3/seqtk) and cutadapt v1.14 (48) with default parameters.

517 After pre-processing, isolates were assembled *de novo* with SPAdes v3.13 (49, 50) using the -

518 *careful* parameter and -k of 21, 33, 55, and 77. Assembled contigs were reordered using Mauve's

519 contig mover function (51) with the complete publicly available Temecula1 assembly

520 (GCA_000007245.1) used as reference. Assembled and reordered genomes were then

Castillo   23

521  individually annotated using the Prokka pipeline (52). In addition, published genome sequences

522  were also re-annotated with Prokka.

523  **Core genome alignments, construction of Maximum Likelihood trees, and haplotype**

524  **network.** Roary v3.11.2 (53) was used to calculate the number of genes in the core (genes shared

525  between 99-100% strains), soft-core (genes shared between 95-99% strains), shell (genes shared

526  between 15-95% strains), and cloud (genes shared between 0-15% strains) genomes of PD-

527  causing isolates (N=175).A core genome alignment of PD-causing isolates plus the three Costa

528  Rica isolates (non-PD) was created using the -e (codon aware multisequence alignment of core

529  genes) and -n (fast nucleotide alignment) flags in Roary. This core genome alignment was used

530  to build a Maximum Likelihood (ML) tree with RAxML (54). The GTRCAT substitution model

531  was used on tree construction, while tree topology and branch support were assessed with 1000

532  bootstrap replicates. In addition, a non-recombinant tree was constructed by removing detected

533  recombinant segments from the core genome alignment (see later methods). The ML non-

534  recombinant tree was constructed using the same parameters as the recombinant tree. Finally, a

535  haplotype network for PD-causing isolates was built following removal of the outgroup

536  sequences (non-PD). Core genome haplotypes were calculated based on the number of mutations

537  among the analyzed strains, and the haplotype network was built using the HaploNet function in

538  the R package 'pegas' (55). Haplotypes were then color-coded by geographic location.

539  **Estimation of recombinant segments and gene gain/loss rates within populations.**

540  Isolates were divided based on their geographical origin: California, Southeast USA, Spain, and

541  Taiwan. California and the Southeast USA were the source population for Spain and Taiwan,

542  respectively. Source and descendant relationships between populations were phylogenetically

543  determined (see results). A core genome alignment was created for California/Spain (N=142),

Castillo  24

544    and Southeast USA/Taiwan (N=33). The alignment was used to estimate the frequency and

545    location of recombinant events. FastGEAR (56) was used with default parameters to identify

546    lineage-specific recombinant segments (ancestral) and strain-specific recombinant segments

547    (recent). The size and location of recombinant segments across the length of the core genome

548    alignment were mapped within California/Spain and Southeast USA/Taiwan using the R package

549    'circlize' (57). Donor/recipient recombinant regions were visualized using fastGEAR's

550    plotRecombinations script. In addition, the number of substitutions introduced by recombination

551    vs. random point mutation (r/m) (58) was estimated for the California/Spain and Southeast

552    USA/Taiwan core genome alignments using ClonalFrameML (59). It should be noted that

553    fastGEAR was designed to test recombination in individual gene alignments instead of core

554    genome alignments; a previous study found that fastGEAR was more conservative than other

555    more appropriate recombination detection methods such as ClonalFrameML (34). Future

556    research should perform an empirical comparison of recombination detection methods for *X.*

557    *fastidiosa.*

558         Additionally, the stochastic probability of gene gain/loss per tree branch was estimated

559    with GLOOME using default parameters (60). Briefly, RAxML was used to build a ML

560    phylogenetic tree for the California/Spain and Southeast USA/Taiwan core genome alignments.

561    The parameters used were the same that for the PD-causing ML tree. Roary v3.11.2 was used to

562    calculate a binary gene presence (1)/absence (0) matrix within the California/Spain and the

563    Southeast USA/Taiwan populations. A binary accessory genome matrix was created by

564    removing core genome genes from the dataset. Subsequently, the binary accessory

565    presence/absence matrix was transposed and converted into FASTA format. The binary

566    accessory genome and the ML trees were used as inputs to the GLOOME analysis. Unique genes

567    were identified through estimating gene gain/loss rates within each population. These genes were

568    annotated by eggNOG-mapper v1.0.3 (https://github.com/eggnogdb/eggnog-mapper) and

569    searched in the Genebank and Pfam databases using BLAST and interproscan v5.47

570    (https://github.com/ebi-pf-team/interproscan).

571        **Population genomics analyses.** Global measures of genetic diversity, population

572    differentiation, and selective sweeps were estimated for the PD-causing dataset using the R

573    package 'PopGenome' (61). The dataset was subdivided in two ways: a) based on isolates'

574    geographical origin (i.e. California, Southeast USA, Spain, and Taiwan), and b) based on

575    isolates' phylogenetic relationships (i.e. PD-I, PD-II, and PD-III; see results). All calculations

576    described below were performed for both the (a) geographic and (b) phylogenetic subdivisions.

577        Genetic diversity was estimated by computing nucleotide diversity ($\pi$), Tajima's D (62),

578    and the Watterson's estimator ($\theta$) (63). Briefly, nucleotide diversity ($\pi$) measures the average

579    number of nucleotide differences per site in pairwise comparisons among DNA sequences.

580    Tajima's D evaluates the frequency of polymorphism present in a population and compares that

581    value to the expectation under neutrality. The Watterson $\theta$ estimator measures the mutation rate

582    of a population. Population differentiation was estimated by calculating the Fixation Index (Fst)

583    (64) within (a) geographic and (b) phylogenetic groups. In addition, the McDonald-Kreitman

584    Test (MKT) (65) was used to estimate the rate of synonymous (syn-P) and non-synonymous

585    (nonsyn-P) polymorphism, against the rate of fixed synonymous (syn-F) and non-synonymous

586    (nonsyn-F) differences. In each instance, the Neutrality Index (NI) was calculated. NI > 1

587    suggests an excess of preserved polymorphism maintained via balancing selection. Alternatively,

588    NI < 1 suggests population divergence via positive selection. Finally, the location and magnitude

589    of selective sweeps was calculated using Nielsen's composite-likelihood-ratio (CLR) (66). This

590  test identifies regions with aberrant allele frequency spectra and estimates if the aberrant allele

591  distribution fits the expectations of a selective sweep. The test was performed on a 1500bp

592  sliding window across the length of the PD-causing core genome alignment.

593      **Data Availability.** The raw sequence data files for the newly published isolates were

594  submitted to the NCBI Sequence Read Archive under accession number SAMN15732826

595  through SAMN15732849. All other used data has been previously published. All accession

596  numbers are listed in Table S1.

597

598

599    **ACKNOWLEDGMENTS**

606

Castillo   28

## References

607

608 1. Pimentel D, Lach L, Zuniga R, Mprrison D. 2000. Environmental and Economic Costs of
609 Nonindigenous Species in the United States. Bioscience 50:53.

610 2. Fletcher J, Bender C, Budowle B, Cobb WT, Gold SE, Ishimaru CA, Luster D, Melcher
611 U, Murch R, Scherm H, Seem RC, Sherwood JL, Sobral BW, Tolin SA. 2006. Plant
612 Pathogen Forensics: Capabilities, Needs, and Recommendations. Microbiol Mol Biol Rev
613 70:450–471.

614 3. Pyšek P, Jarošík V, Pergl J. 2011. Alien plants introduced by different pathways differ in
615 invasion success: Unintentional introductions as a threat to natural areas. PLoS One 6.

616 4. Early R, Bradley BA, Dukes JS, Lawler JJ, Olden JD, Blumenthal DM, Gonzalez P,
617 Grosholz ED, Ibañez I, Miller LP, Sorte CJB, Tatem AJ. 2016. Global threats from
618 invasive alien species in the twenty-first century and national response capacities. Nat
619 Commun 7.

620 5. Mable BK. 2019. Conservation of adaptive potential and functional diversity: integrating
621 old and new approaches. Conserv Genet 20:89–100.

622 6. Baltrus DA, Nishimura MT, Dougherty KM, Biswas S, Mukhtar MS, Vicente J, Holub
623 EB, Dangl JL. 2012. The molecular basis of host specialization in bean pathovars of
624 *Pseudomonas syringae*. Mol Plant-Microbe Interact 25:877–888.

625 7. Karasov TL, Horton MW, Bergelson J. 2014. Genomic variability as a driver of plant–
626 pathogen coevolution? Curr Opin Plant Biol 18:24–30.

627 8. Plissonneau C, Benevenuto J, Mohd-Assaad N, Fouché S, Hartmann FE, Croll D. 2017.
628 Using population and comparative genomics to understand the genetic basis of effector-
629 driven fungal pathogen evolution. Front Plant Sci 8:1–15.

630 9. Mokryakov M V., Abdeev IA, Piruzyan ES, Schaad NW, Ignatov AN. 2010. Diversity of
631 effector genes in plant pathogenic bacteria of genus *Xanthomonas*. Microbiology 79:58–
632 65.

633 10. Rowntree JK, Cameron DD, Preziosi RF. 2011. Genetic variation changes the interactions
634 between the parasitic plant-ecosystem engineer *Rhinanthus* and its hosts. Philos Trans R
635 Soc B Biol Sci 366:1380–1388.

636 11. González R, Butković A, Elena SF. 2019. Role of host genetic diversity for susceptibility-
637 to-infection in the evolution of virulence of a plant virus. Virus Evol 5:1–12.

638 12. Zhu Y, Chen H, Fan J, Wang Y, Li Y, Chen J, Fan JX, Yang S, Hu L, Leung H, Mew TW,
639 Teng PS, Wang Z, Mundt CC. 2000. Genetic diversity and disease control in rice. Nature
640 406:718–722.

641 13. Escriu F. 2012. Diversity of Plant Virus Populations: A Valuable Tool Diversity of Plant
642 Virus Populations: A Valuable Tool for Epidemiological Studies, p. 13. *In* IntechOpen.

643 14. Brown JKM. 2015. Durable Resistance of Crops to Disease: A Darwinian Perspective.
644 Annu Rev Phytopathol 53:513–539.

645 15. Zhan J. 2016. Population Genetics of Plant Pathogens. eLS 1–7.

646 16. Giraud T, Gladieux P, Gavrilets S. 2010. Linking emergence of fungal plant diseases and
647 ecological speciation. Trends Ecol Evol 25:387–395.

648 17. Mhedbi-Hajri N, Hajri A, Boureau T, Darrasse A, Durand K, Brin C, Saux MF Le,
649 Manceau C, Poussier S, Pruvost O, Lemaire C, Jacques MA. 2013. Evolutionary History
650 of the Plant Pathogenic Bacterium *Xanthomonas axonopodis*. PLoS One 8.

651 18. Zhan A, Hu J, Hu X, Zhou Z, Hui M, Wang S, Peng W, Wang M, Bao Z. 2009. Fine-scale

Castillo 29

652      population genetic structure of zhikong scallop (chlamys farreri): Do local marine currents
653      drive geographical differentiation? Mar Biotechnol 11:223–235.
654 19. McDonald BA, Linde C. 2002. The population genetics of plant pathogens and breeding
655      strategies for durable resistance. Euphytica 124:163–180.
656 20. Slatkin M. 1985. Gene flow in natural populations. Annu Rev Ecol Syst Vol 16 393–430.
657 21. McDermott JM, McDonald BA. 1993. Gene flow in plant pathosystems. Annu Rev
658      Phytopathol 31:353–373.
659 22. Pruvost O, Boyer K, Ravigné V, Richard D, Vernière C. 2019. Deciphering how plant
660      pathogenic bacteria disperse and meet: Molecular epidemiology of *Xanthomonas citri*
661      pv. *citri* at microgeographic scales in a tropical area of Asiatic citrus canker endemicity.
662      Evol Appl 12:1523–1538.
663 23. EFSA. 2018. Update of the *Xylella* spp. host plant database. EFSA J 16:1–87.
664 24. Vanhove M, Retchless AC, Sicard A, Rieux A, Coletta-filho HD, Fuente LD La, Stenger
665      DC, Almeida PP. 2019. Genomic Diversity and Recombination among *Xylella fastidiosa*
666      Subspecies. Appl Environ Microbiol 85:1–17.
667 25. Bragard C, Dehnen-Schmutz K, Di Serio F, Gonthier P, Jacques MA, Jaques Miret JA,
668      Justesen AF, MacLeod A, Magnusson CS, Milonas P, Navas-Cortés JA, Potting R,
669      Reignault PL, Thulke HH, Van der Werf W, Vicent Civera A, Yuen J, Zappalà L,
670      Makowski D, Delbianco A, Maiorano A, Muñoz Guajardo I, Stancanelli G, Guzzo M,
671      Parnell S. 2019. Effectiveness of in planta control measures for *Xylella fastidiosa*. EFSA J
672      17.
673 26. Almeida RPP, De La Fuente L, Koebnik R, Lopes JRS, Parnell S, Scherm H. 2019.
674      Addressing the New Global Threat of *Xylella fastidiosa*. Phytopathology 109:172–174.
675 27. Nunney L, Hopkins DL, Morano LD, Russell SE, Stouthamer R. 2014. Intersubspecific
676      Recombination in *Xylella fastidiosa* Strains Native to the United States: Infection of Novel
677      Hosts Associated with an Unsuccessful Invasion. Appl Environ Microbiol 80:1159–1169.
678 28. Nunney L, Yuan X, Bromley RE, Stouthamer R. 2012. Detecting Genetic Introgression:
679      High Levels of Intersubspecific Recombination Found in *Xylella fastidiosa* in Brazil. Appl
680      Environ Microbiol 78:4702–4714.
681 29. Landa BB, Castillo AI, Giampetruzzi A, Kahn A, Román-Écija M, Velasco-Amo MP,
682      Navas-Cortés JA, Marco-Noales E, Barbé S, Moralejo E, Coletta-Filho HD, Saldarelli P,
683      Saponari M, Almeida RPP. 2020. Emergence of a plant pathogen in europe associated
684      with multiple intercontinental introductions. Appl Environ Microbiol 86:1–15.
685 30. Giampetruzzi A, Saponari M, Loconsole G, Boscia D, Savino VN, Almeida RPP, Zicca S,
686      Landa BB, Chacón-Diaz C, Saldarelli P. 2017. Genome-Wide Analysis Provides Evidence
687      on the Genetic Relatedness of the Emergent *Xylella fastidiosa* Genotype in Italy to
688      Isolates from Central America. Phytopathology 107:816–827.
689 31. Saponari M, Giampetruzzi A, Loconsole G, Boscia D, Saldarelli P. 2018. *Xylella*
690      *fastidiosa* in Olive in Apulia: Where We Stand. Phytopathology 109:175–186.
691 32. Nunney L, Azad H, Stouthamer R. 2019. An Experimental Test of the Host-Plant Range
692      of Nonrecombinant Strains of North American *Xylella fastidiosa* subsp. *multiplex*.
693      Phytopathology 109:294–300.
694 33. Castillo AI, Chacón-díaz C, Rodríguez-murillo N, Coletta- HD, Almeida RPP, Rica C.
695      2020. Impacts of local population history and ecology on the evolution of a globally
696      dispersed pathogen . BMC Genomics 21:1–51.
697 34. Vanhove M, Sicard A, Ezennia J, Leviten N, Almeida RPP. 2020. Population structure

Castillo 30

698 and adaptation of a bacterial pathogen in California grapevines. Env Microbiol.

699 35. Gomila M, Moralejo E, Busquets A, Segui G, Olmo D, Nieto A, Juan A, Lalucat J. 2018.
700 Draft Genome Resources of Two Strains of *Xylella fastidiosa* XYL1732/17 and
701 XYL2055/17 Isolated from Mallorca Vineyards. Phytopathology 109:222–224.

702 36. Castillo AI, Tuan S-J, Retchless AC, Hu F-T, Chang H-Y, Almeidaa RPP. 2019. Draft
703 Whole-Genome Sequences of *Xylella fastidiosa* subsp. *fastidiosa* Strains TPD3 and TPD4,
704 Isolated from Grapevines in Hou-li, Taiwan. Microbiology 8:1–3.

705 37. Schuenzel EL, Scally M, Stouthamer R, Nunney L. 2005. A Multigene Phylogenetic
706 Study of Clonal Diversity and Divergence in North American Strains of the Plant
707 Pathogen *Xylella fastidiosa*. Appl Environ Microbiol 71:3832–3839.

708 38. Yuan X, Morano L, Bromley R, Spring-pearson S, Stouthamer R, Nunney L. 2010.
709 Multilocus Sequence Typing of *Xylella fastidiosa* Causing Pierce's disease and Oleander
710 Leaf Scorch in the United States. Ecol Epidemiol 100:601–611.

711 39. Cella E, Angeletti S, Fogolari M, Bazzardi R, De L, Ciccozzi M, Cella E, Angeletti S,
712 Fogolari M, Bazzardi R. 2018. Two different *Xylella fastidiosa* strains circulating in Italy:
713 phylogenetic and evolutionary analyses. J Plant Interact 13:428–432.

714 40. Nunney L, Schuenzel EL, Scally M, Bromley RE, Stouthamerc R. 2014. Large-scale
715 intersubspecific recombination in the plant-pathogenic bacterium *Xylella fastidiosa* is
716 associated with the host shift to mulberry. Appl Environ Microbiol 80:3025–3033.

717 41. Nunney L, Ortiz B, Russell SA, Sánchez RR, Stouthamer R. 2014. The complex
718 biogeography of the plant pathogen *Xylella fastidiosa*: Genetic evidence of introductions
719 and subspecific introgression in Central America. PLoS One 9.

720 42. Kabir P. Tumber JMA and KBF. 2014. Pierce's disease costs California $104 million per
721 year. Calif Agric 68:20–29.

722 43. Hickey C. 2019. Pierce's Disease of Grape: Identification and Management. UGA Coop
723 Ext Bull 1514:1–6.

724 44. Parker JK, Havird JC, De La Fuente L. 2012. Differentiation of *Xylella fastidiosa* strains
725 via multilocus sequence analysis of environmentally mediated genes (MLSA-E). Appl
726 Environ Microbiol 78:1385–1396.

727 45. Francis M, Lin H, Rosa JC La, Doddapaneni H, Civerolo EL. 2006. Genome-based PCR
728 primers for specific and sensitive detection and quantification of *Xylella fastidiosa*. Eur J
729 Plant Pathol 115:203–213.

730 46. Andrews S, Wingett SW, Hamilton RS. 2018. FastQ Screen : A tool for multi-genome
731 mapping and quality control [ version 2 ; referees : 4 approved ] Referee Status :
732 F10000research 1–13.

733 47. Ewels P, Lundin S, Max K. 2016. Data and text mining MultiQC : summarize analysis
734 results for multiple tools and samples in a single report. Bioinformatics 32:3047–3048.

735 48. Marcel M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing
736 reads. EMB.netJournal 17:5–7.

737 49. Bankevich A. 2012. SPAdes: A New Genome Assembly Algorithm and Its Applications
738 to Single-Cell Sequencing. J Comput Biol 19:455–477.

739 50. Nurk S, Bankevich A, Antipov D, Gurevich A, Korobeynikov A, Lapidus A, Prjibelski A,
740 Pyshkin A, Sirotkin A, Sirotkin Y, Stepanauskas R, Clingenpeel S, Woyke T, McLean J,
741 Lasken R, Tesler G, Alekseyev M, Pevzner P. 2013. Assembly single-cell genomes and
742 mini-metagenomes from chimeric MDA products. J Comput Biol 20:714–737.

743 51. Rissman AI, Mau B, Biehl BS, Darling AE, Glasner JD, Perna NT. 2009. Reordering

744        contigs of draft genomes using the Mauve Aligner. Bioinformatics 25:2071–2073.

745  52.    Seemann T. 2014. Prokka: Rapid prokaryotic genome annotation. Bioinformatics
746        30:2068–2069.

747  53.    Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, Fookes M, Falush
748        D, Keane JA, Parkhill J. 2015. Roary : rapid large-scale prokaryote pan genome analysis.
749        Bioinformatics 31:3691–3693.

750  54.    Stamatakis A. 2014. RAxML version 8 : a tool for phylogenetic analysis and post-analysis
751        of large phylogenies. Bioinformatics 30:1312–1313.

752  55.    Paradis E. 2010. Pegas: An R package for population genetics with an integrated-modular
753        approach. Bioinformatics 26:419–420.

754  56.    Mostowy R, Croucher NJ, Andam CP, Corander J, Hanage WP, Marttinen P. 2017.
755        Efficient Inference of Recent and Ancestral Recombination within Bacterial Populations.
756        Mol Biol Evol 34:1167–1182.

757  57.    Gu Z, Gu L, Eils R, Schlesner M, Brors B. 2014. circlize implements and enhances
758        circular visualization in R. Bioinformatics 30:2811–2812.

759  58.    Guttman DS, Dykhuizen DE. 1994. Clonal divergence in Escherichia coli as a result of
760        recombination, not mutation. Science (80- ) 266:1380–1383.

761  59.    Didelot X, Wilson DJ. 2015. ClonalFrameML: Efficient Inference of Recombination in
762        Whole Bacterial Genomes. PLoS Comput Biol 11:1–18.

763  60.    Cohen O, Ashkenazy H, Belinky F, Huchon D, Pupko T. 2010. GLOOME : gain loss
764        mapping engine. Bioinformatics 26:2914–2915.

765  61.    Pfeifer B, Wittelsbu U, Ramos-onsins SE, Lercher MJ. 2014. PopGenome : An Efficient
766        Swiss Army Knife for Population Genomic Analyses in R. Mol Biol Evol 31:1929–1936.

767  62.    Tajima F. 1989. Statistical Method for Testing the Neutral Mutation Hypothesis by DNA
768        Polymorphism. Genet Soc Am 595:585–595.

769  63.    Watterson GA. 1975. On the numer of Segregating Sites in Genetical Models without
770        Recombination. Theor Popul Biol 276:256–276.

771  64.    Wrigth S. 1965. The interpretation of population structure by F-statistics with special
772        regard to systems of mating. Evolution (N Y) 19:395–420.

773  65.    McDonald JH, Kreitman M. 1991. Adaptive protein evolution at the Adh locus in
774        *Drosophila*. Nature 351:652–654.

775  66.    Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, Bustamante C. 2005. Genomic
776        scans for selective sweeps using SNP data. Genome Res 15:1566–1575.

777  67.    Nunney L, Yuan X, Bromley R, Hartung J, Montero-Astúa M, Moreira L, Ortiz B,
778        Stouthamer R. 2010. Population genomic analysis of a bacterial plant pathogen: Novel
779        insight into the origin of Pierce's disease of grapevine in the U.S. PLoS One 5.

780  68.    Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin E V. 2014. Genomes in
781        turmoil: Quantification of genome dynamics in prokaryote supergenomes. BMC Med
782        12:1–19.

783  69.    Iranzo J, Wolf YI, Koonin E V., Sela I. 2019. Gene gain and loss push prokaryotes beyond
784        the homologous recombination barrier and accelerate genome sequence divergence. Nat
785        Commun 10.

786  70.    Hartmann FE, Croll D. 2017. Distinct trajectories of massive recent gene gains and losses
787        in populations of a microbial eukaryotic pathogen. Mol Biol Evol 34:2808–2822.

788  71.    Kettler GC, Martiny AC, Huang K, Zucker J, Coleman ML, Rodrigue S, Chen F, Lapidus
789        A, Ferriera S, Johnson J, Steglich C, Church GM, Richardson P, Chisholm SW. 2007.

790      Patterns and implications of gene gain and loss in the evolution of *Prochlorococcus*. PLoS
791      Genet 3:2515–2528.

792 72. Moulana A, Anderson RE, Fortunato CS, Huber JA. 2020. Selection Is a Significant
793      Driver of Gene Gain and Loss in the Pangenome of the Bacterial Genus *Sulfurovum* in
794      Geographically Distinct Deep-Sea Hydrothermal Vents . mSystems 5:1–18.

795 73. Frick DN, Richardson CC. 2001. DNA Primases. Annu Rev Biochem 70:39–80.

796 74. Browning DF, Busby SJW. 2016. Local and global regulation of transcription initiation in
797      bacteria. Nat Rev Microbiol 14:638–650.

798 75. Kuo CH, Moran NA, Ochman H. 2009. The consequences of genetic drift for bacterial
799      genome complexity. Genome Res 19:1450–1454.

800 76. Albalat R, Cañestro C. 2016. Evolution by gene loss. Nat Rev Genet 17:379–391.

801 77. Nunney L, Vickerman DB, Bromley RE, Russell SA, Hartman JR, Morano LD,
802      Stouthamer R. 2013. Recent Evolutionary Radiation and Host Plant Specialization in the
803      *Xylella fastidiosa* Subspecies Native to the United States. Appl Environ Microbiol
804      79:2189–2200.

805 78. Kandel PP, Almeida RPP, Cobine PA, De La Fuente L. 2017. Natural Competence Rates
806      Are Variable Among *Xylella fastidiosa* Strains and Homologous Recombination Occurs In
807      Vitro Between Subspecies *fastidiosa* and *multiplex*. Mol Plant-Microbe Interact 30:589–
808      600.

809 79. Potnis N, Kandel PP, Merfa M V., Retchless AC, Parker JK, Stenger DC, Almeida RPP,
810      Bergsma-Vlami M, Westenberg M, Cobine PA, De La Fuente L. 2019. Patterns of inter-
811      and intrasubspecific homologous recombination inform eco-evolutionary dynamics of
812      *Xylella fastidiosa*. ISME J 13:2319–2333.

813 80. Schmutzer M, Barraclough TG. 2019. The role of recombination, niche-specific gene
814      pools and flexible genomes in the ecological speciation of bacteria. Ecol Evol 9:4544–
815      4556.

816 81. Barbosa RL, Rinaldi FC, Guimarães BG, Benedetti CE. 2007. Crystallization and
817      preliminary X-ray analysis of BigR, a transcription repressor from *Xylella fastidiosa*
818      involved in biofilm formation. Acta Crystallogr Sect F Struct Biol Cryst Commun
819      63:596–598.

820 82. Barbosa RL, Benedetti CE. 2007. BigR, a transcriptional repressor from plant-associated
821      bacteria, regulates an operon implicated in biofilm growth. J Bacteriol 189:6185–6194.

822 83. Guimarães BG, Barbosa RL, Soprano AS, Campos BM, De Souza TA, Tonoli CCC,
823      Leme AFP, Murakami MT, Benedetti CE. 2011. Plant pathogenic bacteria utilize biofilm
824      growth-associated repressor (BigR), a novel winged-helix redox switch, to control
825      hydrogen sulfide detoxification under hypoxia. J Biol Chem 286:26148–26157.

826 84. De Lira NPV, Pauletti BA, Marques AC, Perez CA, Caserta R, De Souza AA, Vercesi
827      AE, Paes Leme AF, Benedetti CE. 2018. BigR is a sulfide sensor that regulates a sulfur
828      transferase/dioxygenase required for aerobic respiration of plant bacteria under sulfide
829      stress. Sci Rep 8:1–13.

830 85. Federici MT, Marcondes JA, Picchi SC, Stuchi ES, Fadel AL, Laia ML, Lemos MVF,
831      Lemos EGM. 2012. *Xylella fastidiosa*: An *in vivo* system to study possible survival
832      strategies within citrus xylem vessels based on global gene expression analysis. Electron J
833      Biotechnol 15.

834 86. Da Silva Neto JF, Koide T, Gomes SL, Marques M V. 2007. The single extracytoplasmic-
835      function sigma factor of *Xylella fastidiosa* is involved in the heat shock response and

836       presents an unusual regulatory mechanism. J Bacteriol 189:551–560.

837   87.   Lee KM, Go J, Yoon MY, Park Y, Kim SC, Yong DE, Yoon SS. 2012. Vitamin B 12-
838         Mediated restoration of defective anaerobic growth leads to reduced biofilm formation in
839         *Pseudomonas aeruginosa*. Infect Immun 80:1639–1649.

840   88.   Cordonnier C, Le Bihan G, Emond-Rheault JG, Garrivier A, Harel J, Jubelin G. 2016.
841         Vitamin B12 uptake by the gut commensal bacteria bacteroides thetaiotaomicron limits
842         the production of shiga toxin by enterohemorrhagic *Escherichia coli*. Toxins (Basel) 8.

843   89.   Chen H, De La Fuente L. 2020. Calcium transcriptionally regulates movement,
844         recombination and other functions of *Xylella fastidiosa* under constant flow inside
845         microfluidic chambers. Microb Biotechnol 13:548–561.

846   90.   Cursino L, Li Y, Zaini PA, De La Fuente L, Hoch HC, Burr TJ. 2009. Twitching motility
847         and biofilm formation are associated with tonB1 in *Xylella fastidiosa*. FEMS Microbiol
848         Lett 299:193–199.

849   91.   Santos-Beneit F. 2015. The Pho regulon: A huge regulatory network in bacteria. Front
850         Microbiol 6:1–13.

851   92.   Singh J, Khan A. 2019. Distinct patterns of natural selection determine sub-population
852         structure in the fire blight pathogen, Erwinia amylovora. Sci Rep 9:1–13.

853   93.   Ørsted M, Hoffmann AA, Sverrisdóttir E, Nielsen KL, Kristensen TN. 2019. Genomic
854         variation predicts adaptive evolutionary responses better than population bottleneck
855         history. PLoS Genet 15:1–18.

856   94.   Holderegger R, Kamm U, Gugerli F. 2006. Adaptive vs. neutral genetic diversity:
857         Implications for landscape genetics. Landsc Ecol 21:797–807.

858   95.   Moutinho AF, Bataillon T, Dutheil JY. 2019. Variation of the adaptive substitution rate
859         between species and within genomes. Evol Ecol.

860   96.   Kung SH, Almeida RPP. 2011. Natural competence and recombination in the plant
861         pathogen *Xylella fastidiosa*. Appl Environ Microbiol 77:5278–5284.

862
863
864
865
866
867
868
869
870
871
872
873
874
875

876

877

878

Castillo   34

879 **FIGURES LEGENDS**

880 **Figure 1. Maximum Likelihood (ML) tree and haplotype network showing phylogenetic**

881 **and geographic diversification of worldwide PD-causing subsp.** *fastidiosa* **isolates.** Color

882 represents isolates from the same geographical location: California (Red), Texas (Pink), Georgia

883 (Green), North Carolina (Dark green), Florida (Yellow), Spain (Light blue), and Taiwan (Dark

884 blue). PD-causing strains have been divided into three phylogenetically supported clades (PD-I,

885 PD-II, PD-III). **a.** Haplotype network of PD-causing subsp. *fastidiosa* isolates. Haplotypes

886 belonging to each PD-causing clade are shown within black circles. Roman numbers identify

887 detected haplotypes (I-CXLI). Size of the circle indicates the number of isolates belonging to

888 each haplotype. **b.** Maximum likelihood (ML) tree of PD-causing subsp. *fastidiosa* isolates. Tree

889 was built using the core genome alignment without removing recombinant segments. Bootstrap

890 values mark branch support. Arrows point towards the base of PD-causing clades (-I to -III).

891

892 **Figure 2. Phylogeographic analysis showing diversification of PD-causing isolates within**

893 **the contiguous USA.** Color represents isolates from the same geographical location: California

894 (Red), Texas (Pink), Georgia (Green), North Carolina (Dark green), and Florida (Yellow).

895 Coordinates were recorded during field sampling. In absence of this information, coordinates

896 referred to the city or vineyard closest to the sample site were used. Florida coordinates were not

897 available, the location shown in the map represents central Florida. Isolates from Southern and

898 Northern California are shown within pale red circles. PD-causing strains were divided into three

899 phylogenetically supported clades: PD-I (Southeast USA isolates exclusively), PD-II (Southern

900 California and Texas isolates), and PD-III (both Southern and Northern California isolates, and

901 three Georgia isolates). Tree was built using the core genome alignment without removing

902 recombinant segments. Bootstrap values mark branch support.

Castillo 35

903 **Figure 3. Venn diagram and maps showing population linked gene gain/loss events among**

904 **PD-causing isolates.** Color represents isolates from the same geographical location: California

905 (Red), Texas (Pink), Georgia (Green), North Carolina (Dark green), Florida (Yellow), Spain

906 (Light blue), and Taiwan (Dark blue). **a.** Venn diagram shows both the number of genes shared

907 between geographic PD-causing populations and genes unique to each population. Size of the

908 oval represents sample size. **b.** Estimated number of genes gained and lost between geographical

909 locations and following introduction events. Arrows point from the source population to its

910 descendant following introduction events. California isolates belong to the phylogenetically

911 distinct clades PD-II and PD-III. Included Southeast isolates belong to the phylogenetically

912 distinct PD-I and PD-III clades. All maps were publicly available from Wikimedia commons.

913

914 **Figure 4. Frequency and location of recombination events in fastGEAR identified lineages.**

915 Analysis shows results for: **a.** the California/Spain population and **b.** the Southeast USA/Taiwan

916 population. FastGEAR's recombination plots show two distinct lineages on each population (red,

917 PD-III in California/Spain and PD-II/PD-III in Southeast USA/Taiwan; blue, PD-II in

918 California/Spain and PD-I in Southeast USA/Taiwan). The recombination events are shown

919 across the length of the core genome alignment. Larger areas represent recipient sequences while

920 shorter segments of different color within those areas represent donor sequences from another

921 lineage. Recombinant segments from unidentified lineages are shown in black. Maximum

922 Likelihood (ML) trees showing the phylogenetic relationship of isolates within each intra-

923 population cluster identified by fastGEAR are also included. Trees were built using the core

924 genome alignment without removing recombinant segments for the California/Spain and

925 Southeast USA/Taiwan populations. Bootstrap values mark branch support.

Castillo   36

926  **Figure 5. Line plot showing variations in Nielsen's composite likelihood ratio (CLR) across**

927  **the length of the core genome alignment (1500 bp window size).** The CLR identifies regions

928  with aberrant allele frequency and determines if their distribution matches those expected from a

929  selective sweep. Peaks represent higher CLR values at that position, which is indicative of a

930  putative selective sweep. Color represents isolates from the same geographical location or

931  phylogenetic cluster. **a.** Lines indicate distinct geographic population: California (Red),

932  Southeast USA (Green), Spain (Light blue), and Taiwan (Dark blue). **b.** Lines indicated distinct

933  phylogenetic clusters: PD-I (Teal), PD-II (Yellow), and PD-III (Purple).

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

Castillo   37

949 **Table 1.** List of genes gained/lost among geographic and phylogenetic PD-causing groups.

| Annotation |
| --- |
| **Genes absent in the Taiwan population (Found in the contiguous USA)** |
| Site-specific DNA-methyltransferase (QIS25725.1); ko:K00571,ko:K00590,ko:K07319 (adenine-specific DNA-methyltransferase); PF01555 |
| Hypothetical protein (QIS26419.1) |
| Peptidoglycan DD-metalloendopeptidase family protein (QIS26766.1); PF06594, PF00353 (RTX calcium-binding nonapeptide repeat) |
| Site-specific DNA-methyltransferase (QIS25737.1); ko:K00571,ko:K00590,ko:K07319 (adenine-specific DNA-methyltransferase); PF01555 |
| Pseudogene |
| **Genes absent in the Spanish population (Found in the contiguous USA)** |
| CPD*: LacZ, Beta-galactosidase/beta-glucuronidase; ko:K01192(beta-mannosidase) |
| Glutamate 5-kinase (AAO28181.1); ko:K00931; PF00696 |
| **PD-II and PD-III exclusive genes in the California population** |
| Alpha/beta fold hydrolase (QIS25057.1); ko:K02170,ko:K07002 (pimeloyl-[acyl-carrier protein] methyl ester esterase) |
| Hypothetical protein (QJP55224.1); PF04014 (Antidote-toxin recognition MazE, bacterial antitoxin) |
| Hypothetical protein (AAO28982.1) |
| **PD-I and PD-III exclusive genes in the Southeast population, excluding the Taiwanese clade** |
| Hypothetical protein (QIS26118.1) |
| Phage head morphogenesis protein (QIS26295.1); PF04233 |
| **Genes gained after introduction into Taiwan** |
| CPD*: OM_channels Superfamily, Porin superfamily |
| CPD*: DUF769 Superfamily; ko:K15125 (filamentous hemagglutinin) |
| CPD*: entero_EhxA Superfamily |
| Hypothetical protein (QIS25070.1); RTX toxin (QIS25071.1) |
| **Genes gained after introduction into Spain** |
| Pseudogene: Glycoside hydrolase family 125 protein; ko:K09704 (uncharacterized protein); PF06824 (Metal-independent alpha-mannosidase) |
| DUF596 domain-containing protein (QIS26773.1); PF04591 |
| Hypothetical protein (QID15519.1) Hemagglutinin (QID15518.1); ko:K02014 (iron complex outer membrane receptor protein) |

950 **ko: KEGG orthology; PF: Pfam database entry ID**

951

952

953

954

Castillo 38

955 **Table 2.** Diversity and neutrality statistics of PD-causing isolates. a) Geographically divided

956 populations, and b) Phylogenetically divided populations.

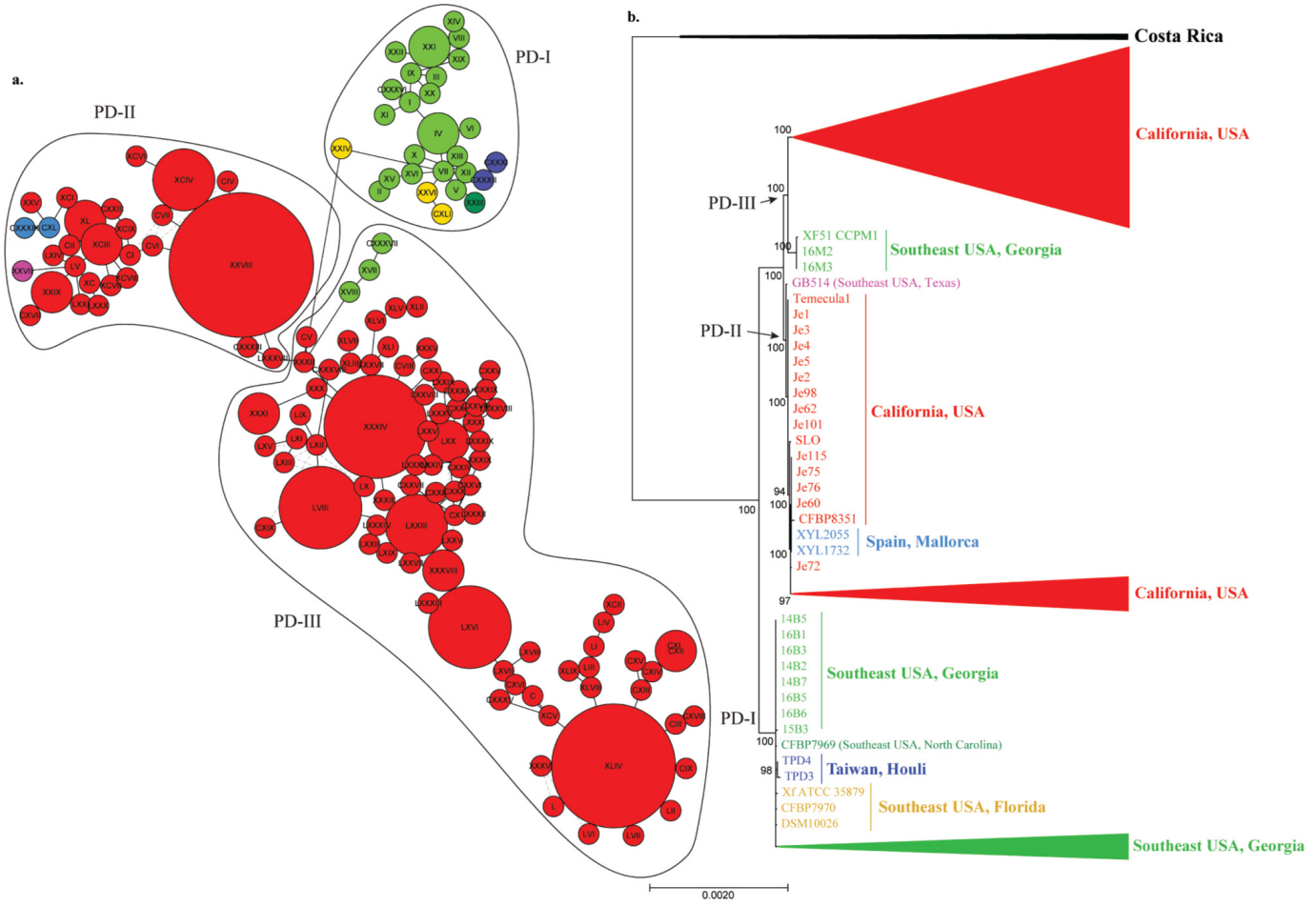| Population (a) | Core (nt) | SNPs | $\pi$ | $\theta$ | Tajima's D |
|---|---|---|---|---|---|
| California (140) | | 458 | $3.22 \times 10^{-06}$ | $1.64 \times 10^{-05}$ | -1.448 |
| Southeast USA (31) | 14,446,213 | 947 | $1.36 \times 10^{-05}$ | $5.75 \times 10^{-06}$ | -0.658 |
| Spain (2) | | 2 | $1.38 \times 10^{-07}$ | $1.38 \times 10^{-07}$ | * |
| Taiwan (2) | | 6 | $4.15 \times 10^{-07}$ | $4.15 \times 10^{-07}$ | * |

\* Spain isolates were not included.
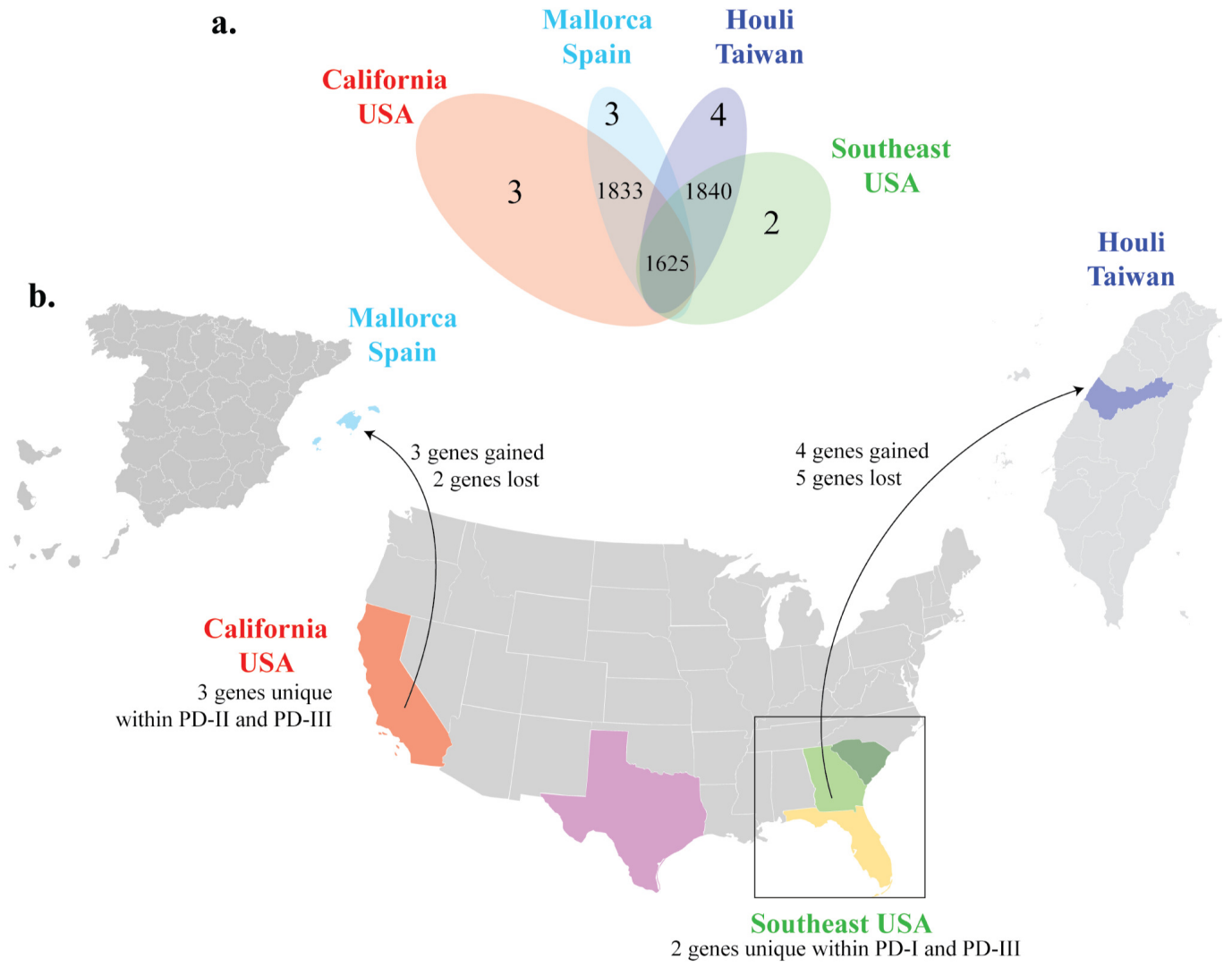\* Taiwan isolates were not included.

957

| Population (b) | Core (nt) | SNPs | $\pi$ | $\theta$ | Tajima's D |
|---|---|---|---|---|---|
| PD-I (29) | | 93 | $7.58 \times 10^{-07}$ | $1.64 \times 10^{-06}$ | -2.0604 |
| PD-II (40) | 14,446,213 | 114 | $9.65 \times 10^{-07}$ | $1.87 \times 10^{-06}$ | -1.7813 |
| PD-III (106) | | 509 | $3.25 \times 10^{-06}$ | $6.72 \times 10^{-06}$ | -1.7425 |

958

**a.**

**Mallorca Spain**

**Houli Taiwan**

**California USA**

3

4

**Southeast USA**

3

1833

1840

2

1625

**b.**

**Mallorca Spain**

**Houli Taiwan**

3 genes gained
2 genes lost

4 genes gained
5 genes lost

**California USA**
3 genes unique
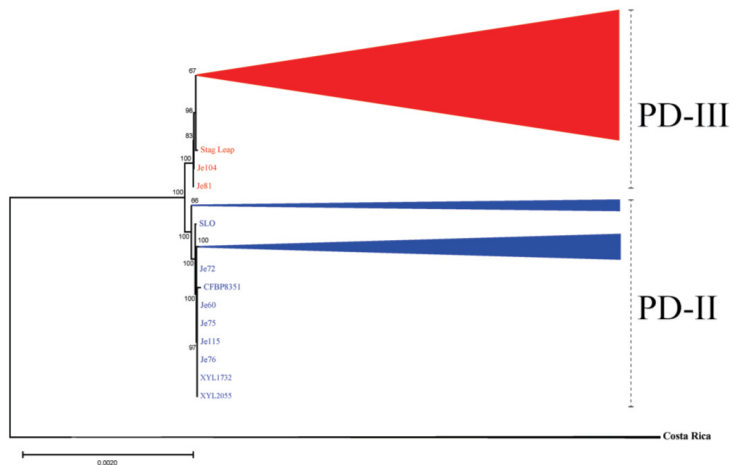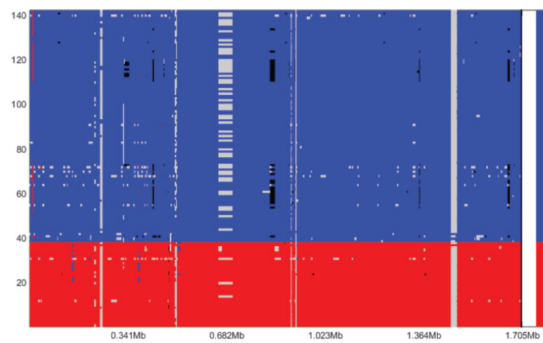within PD-II and PD-III

**Southeast USA**
2 genes unique within PD-I and PD-III

**a.** California/Spain

**b.** Southeast USA/Taiwan

PD-III

PD-II

PD-I

PD-II/PD-III

**a.**



CLR

Position in core genome alignment

—— California, USA   —— Southeast, USA   —— Mallorca, Spain   —— Hou-li, Taiwan

**b.**



CLR

Position in core genome alignment

—— PD-I   —— PD-II   —— PD-III