1 **Genomic acquisitions in emerging populations of *Xanthomonas vasicola* pv.**

2 ***vasculorum* infecting corn in the U.S. and Argentina**

3 Alvaro L Perez-Quintero[1], Mary Ortiz-Castro[1], Guangxi Wu[1], Jillian M. Lang[1], Sanzhen Liu[2],

4 Toni A Chapman[3], Christine Chang[4], Janet Ziegle[4], Zhao Peng[5], Frank F. White[5], Maria

5 Cristina Plazas[6], Jan E. Leach[1], Kirk Broders.[1,7]*

6 1. Bioagricultural Sciences and Pest Management, Colorado State University, Fort Collins,

7 CO

8 2. Department of Plant Pathology, Kansas State University, Manhattan, KS

9 3. Biosecurity and Food Safety, NSW Department of Primary Industries, Elizabeth Macarthur

10 Agricultural Institute, Menangle, NSW, Australia

11 4. Pacific Biosciences, Menlo Park, CA

12 5. Department of Plant Pathology, University of Florida, Gainesville, FL

13 6. Laboratorio de Fitopatología y Microbiología, Universidad Católica de Córdoba, Ob. Trejo

14 323 Córdoba, Argentina.

15 7. Smithsonian Tropical Research Institute, Apartado 0843-03092, Balboa, Ancon,

16 Republic of Panamá.

17 *Corresponding author

18 BrodersK@si.edu

19 **Abstract**

20 *Xanthomonas vasicola* pv. *vasculorum (Xvv)* is an emerging bacterial plant pathogen that

21 causes bacterial leaf streak on corn. First described in South Africa in 1949, reports of this

22 bacteria have greatly increased in the past years in South America and in the U.S., where it is

23 now present in most of the corn producing states. Phenotypic characterization showed that

24 the emerging U.S. and South American X*vv* populations may have increased virulence in corn

25 compared to older strains. To understand the genetic mechanisms behind the increased

26 virulence in this group, we used comparative genomics to identify gene acquisitions in *Xvv*

27 genomes from the U.S. and Argentina. We sequenced 41 genomes of *Xvv* and the related

28 sorghum-infecting *X. vasicola* pv. *holcicola* (Xvh). A comparison of all available *X. vasicola*

29 genomes showed the phylogenetic relationships in the group and identified clusters of genes

30 associated with the emerging *Xvv* populations. The newly acquired gene clusters showed

31 evidence of horizontal transfer to *Xvv* and included candidate virulence factors. One cluster, in

32 particular, corresponded to a prophage transferred from *Xvh* to all *Xvv* from Argentina and the

33 U.S. The prophage contains putative secreted proteins, which represent candidates for

34 virulence determinants in these populations and await further molecular characterization.

35 **Key words***: Xanthomonas vasicola pv. vasculorum (Xvv), Xanthomonas vasicola pv.*

36 *holcicola          (Xvv),          corn,          horizontal          gene          transfer.*

## Introduction

In the U.S., bacterial leaf streak of corn (BLS) was first observed in Nebraska in 2014 and became widespread by 2016 (Korus et al. 2017). The disease is present in dent corn and popcorn producing regions of Colorado, Kansas and Nebraska, with several fields reporting disease incidence levels above 90% and disease severity reaching greater than 50% of leaf area infected (Broders 2017). The disease has now been found in most of the corn growing region of the U.S. including Illinois, Iowa and Nebraska, which are the top three corn producing states in the U.S. (Korus et al. 2017; USDA-NASS 2017). Given the large number of acres and economic importance of corn production in the U.S., there are important implications to the emergence and spread of this new disease. Thus, understanding how this disease originated and what favors its spread is crucial to prevent future losses.

Caused by *X. vasicola* pv. *vasculorum (Xvv),* BLS was first described in 1949 on corn in South Africa (Dyer 1949), but prior to 2016 it had not been documented in any other country. It is unknown how this organism was introduced to the U.S. or if it was already present in a latent state. The only other report of BLS of corn outside of South Africa and the U.S. was in Argentina and Brazil in 2017 and 2018 (Plazas et al. 2017, Leite et al. 2019). While the official report of the disease in Argentina is relatively recent, the symptoms of BLS were first observed in 2010 in the Cordoba province and have since spread to nine other corn-producing provinces, including provinces that border the corn growing regions in Brazil and Paraguay (Leite et al. 2019; Plazas et al. 2017). It is still unclear why *Xvv* continues to spread in the Americas or how severe future epidemics may become. However, it does appear that *Xvv* may have been present in Argentina prior to reports of this pathogen in the U.S.

59    A significant amount of confusion existed around the taxonomic classification of this

60    bacterium. The nomenclature had gone through several changes, from *X. campestris* pv. *zeae*

61    to *X. vasicola* pv. *zeae* to its current designation as *X. vasicola* pv. *vasculorum* (Lang et al.

62    2017; Coutinho and Wallis 1991; Sanko et al. 2018; Bradbury 1986; Qhobela, Claflin, and

63    Nowell 1990). The *X. vasicola* species is now divided into five groups including three named

64    pathovars: 1) *Xvv* infecting corn and sugarcane, 2) *X. vasicola* pv. *holcicola (Xvh)*, commonly

65    infecting sorghum, 3) *X. vasicola* pv. *musacearum (Xvm)* infecting enset and banana, 4) a

66    group of strains isolated from *Tripsacum laxum*, and, 5) strains isolated from areca nut

67    (previously *X. campestris* pv. *arecae*) (Lang et al. 2017) (Studholme et al., this issue).

68    The term pathovar refers to a strain or set of strains with the same or similar characteristics,

69    differentiated at the infrasubspecific level from other strains of the same species or

70    subspecies on the basis of distinctive pathogenicity to one or more plant hosts (Young et al.

71    1991). While the named pathovars of *X. vasicola* seem to have defined host preferences,

72    their host ranges may be broader than initially claimed. *Xvh* and *Xvv,* in particular, may have

73    overlapping host ranges. Under laboratory conditions, isolates of *Xvv* from corn and

74    sugarcane caused disease on corn, sugarcane and sorghum, but were most virulent on corn

75    and sugarcane (Lang et al. 2017). Similarly, when infiltrated into leaves, *Xvh* infected corn,

76    sorghum and sugarcane, but caused more disease on sorghum (Lang et al. 2017). *Xvv* has

77    not been isolated from sorghum, while *Xvh* has occasionally been isolated from corn in the

78    field (Moffett 1983; Péros et al. 1994). Upon inoculation in the greenhouse, *Xvv* isolates from

79    the U.S. can infect 16 hosts, mostly monocots such as rice, oats and big blue stem and the

80    dicot yellow nutsedge (Hartman et al. this issue). Field studies confirmed these results for big

81    blue stem and bristly foxtail as hosts in a natural inoculum system (Hartman et al. this issue).

4

82    At least two host jumps have been hypothesized for *X. vasicola*, i.e. from grasses to banana

83    (Tushemereirwe et al. 2004) and from sugarcane to *Eucalyptus* spp, a dicot (Coutinho et al.

84    2015), suggesting a remarkable adaptive ability for the species.

85    The appearance and spread of *Xvv* in the U.S. and Argentina was rapid. How and why the

86    populations expanded so quickly in two countries located on the opposite side of the equator

87    at approximately the same time, while the disease remained rarely documented in South

88    Africa during the same time period, is intriguing. Reasons for these disparate observations

89    could include the occurrence of more favorable environmental conditions and susceptible

90    corn germplasm in the Americas versus South Africa, and/or, as we hypothesize here, the

91    acquisition of genetic features that favored infection or spread or virulence. In this study, we

92    employed a comparative genomics approach to identify genetic changes associated to these

93    emerging populations.

## Results

**U.S. *Xvv* strains are closely related to each other and to *Xvv* strains from Argentina.**

Draft genome assemblies were generated for fifteen strains collected in 2016 in the states of Colorado, Iowa, Kansas and Nebraska. Draft genome sequences were also prepared for *Xvv* isolates from South Africa (2 strains from corn) and Argentina (7 strains), and for *Xvh* isolates from the U.S. (1 strain, from sorghum) and Australia (8 strains from sorghum, 3 from corn). Additionally, fully assembled genomes were derived for three U.S. *Xvv* isolates isolated from corn, one sugarcane *Xvv* isolate from Zimbabwe, and one sorghum *Xvh* isolate from Mexico, totaling 41 new genomes (Supplementary Table 1).

A phylogenetic maximum-likelihood tree was determined on the basis of the pan-genome SNPs (Gardner, Slezak, and Hall 2015), using the newly sequenced genomes as well as all available *Xanthomonas vasicola* genomes, adding *Xanthomonas oryzae* pv. *oryzae* PXO99A as an outgroup (94 genomes in total) (Supplementary Table 1). The tree reveals division of the four main groups/pathovars in the species: *holcicola (Xvh), musacearum (Xvm), vasculorum (Xvv)* and an unnamed group of strains collected from *Tripsacum laxum* (here referred to as simply *Xv*) (Figure 1 A). Most *Xvv* strains from corn formed a closely-related group, separated from strains isolated from sugarcane, with the exception of strain NCPPB206, a weakly virulent isolate from South Africa collected in 1948 (Lang et al. 2017; Dyer 1949). *Xvv* strains from the U.S. tended to group together (one exception NE-7 from Nebraska), and seemed to be more closely related to strains from Argentina than to strains from South Africa (Figure 1 A).

Similar groupings were seen in phylogenetic trees based on core-genome SNPs (Supplementary Figure 1 A), multi-locus sequence alignments (MLSA) from house-keeping

6

117 genes or core-genome genes (Supplementary Figure 1), and average nucleotide identity

118 (Supplementary Figure 2). Additionally, minimum spanning trees (MSTs) based on core-

119 genome SNPs showed the same main groups but also revealed possible relations between

120 the different corn *Xvv* groups. Unlike phylogenetic trees, MSTs only assume identity based

121 only upon state (SNP) similarity and not upon relationships to hypothetical ancestors

122 (Salipante and Hall 2011). By MST, Argentinian *Xvv* strains form a central cluster connected

123 to U.S. and South African *Xvv* strains but also to *Xvh*, which indicates a possible transmission

124 path whereby the Argentinian *Xvv* population is a link between the U.S. and South African

125 ones, and possibly also received genetic material from *Xvh* (Figure 1B).

126 Disease phenotyping showed that contemporary *Xvv* isolates from the US and Argentina tend

127 to cause more severe symptoms on corn (hybrid P1151) than older *Xvv* isolates form South

128 Africa and contemporary *Xvh*, although with high variation (Figure 2). Suggesting the recent

129 epidemic may be associated to a gain of virulence in the American populations. These results

130 however may be dependent on the host genotype since more severe symptoms have been

131 reported for *Xvh* and South African *Xvv* in another corn variety (Lang et al. 2017).

132 **U.S. *Xvv* genomes contain clusters of genes often absent in other *X. vasicola***

133 **genomes.**

134 To identify genes associated to the emerging U.S. *Xvv* population, we identified ortholog

135 groups among all proteins. Overall, 3616 genes were present in all the *X. vasicola* groups

136 (Supplementary Figure 3), and the core genome (orthologs present in all strains) was 2755

137 genes (Supplementary Figure 3). Orthology distribution analysis in the different groups found

138 that 44 genes were exclusive to corn *Xvv*. No genes were found to be exclusive to *Xvv* from

139  the U.S., but 19 genes were shared exclusively with Argentinian strains. Sixty-three genes

140  were shared between *Xvh* and Argentinian and U.S. *Xvv* strains (Supplementary Figure 3).

141  To uncover more genes associated with (but not necessarily exclusive to) the U.S. *Xvv*

142  population, we performed a hyper-geometric test for each ortholog group. In this test we

143  compared the presence and copy number of each gene in the U.S. *Xvv* genomes against 100

144  sets of randomly selected genomes from the other groups (Supplementary Figure 4). The test

145  identified 278 genes that were over-represented in U.S. *Xvv* and 44 genes that were

146  underrepresented (Figure 3, Supplementary Table 2). The Xvv U.S.-associated genes in the

147  reference genome Xvv strain CO-5 were often grouped together, indicating sub-genomic

148  regions are associated with the population. Five clusters of over-represented genes (named A

149  to E) were identified (Figure 3). The clusters contain 141 genes, with two large clusters

150  containing 44(C) and 57(E) genes, respectively. Cluster C was shared among a majority of

151  the corn *Xvv* strains, and some genes are also present in *Xvh*. All genes from Cluster E are

152  shared by *Xvh* and U.S. and Argentinian *Xvv* strains, while absent of most South African *Xvv*

153  strains (Figure 3).

154  The annotations and predicted functions of over-represented genes showed an enrichment in

155  proteins with nucleic binding activity and involvement in DNA metabolism, recombination and

156  transposition, suggesting mobile elements in the strains (Supplementary Figure 5). No

157  particular enrichment was found in the group of underrepresented genes.

158  **Clusters of genes in U.S. *Xvv* are genomic islands and contain putative effectors.**

159  We further analyzed these clusters in a set of eight fully sequenced genomes. Most of the

160  identified clusters (except Cluster B) are predicted to be in genomic islands, using the

161  IslandViewer4suite, which compiles parametric and phylogenetic methods for genomic island

8

162   prediction (Bertelli et al. 2017) and implies that the clusters were acquired by horizontal

163   transfer (Figure 4). Furthermore, Cluster C is particularly enriched in insertion sequences (IS)

164   transposition-associated genes (Figure 4 and Supplementary Figure 5). IS were absent in

165   Cluster E.

166   None of the over or underrepresented genes matched against known Type III effectors (T3E)

167   by blast (Altschul et al. 1997). Furthermore, no specific association was found between

168   effector presence/absence and U.S. *Xvv* or corn *Xvv,* in general, with the possible exception

169   of XopG1, a M27 zinc protease (White et al. 2009), which is absent in most *Xvv* strains and

170   present in other *X. vasicola* pathovars (Supplementary Figure 6).

171   Additionally, the suite EffectiveDB, which allows identification of putative T3Es based on

172   prediction of secretion signals, T3 chaperone binding domains, eukaryotic-like domains and

173   eukaryotic subcellular localization (Eichinger et al. 2016), predicted that ~450 proteins were

174   putative T3Es in each genome (having, at minimum, a predicted T3E secretion signal) (Figure

175   4). Some predicted T3Es proteins were found in the clusters, including five genes in Cluster C

176   and nine in Cluster E, the latter included various hypothetical proteins, an HTH transcriptional

177   regulator, and two methyltransferases (Supplementary Table 3).

178   **A gene cluster in U.S. *Xvv* is a prophage horizontally transferred from *Xvh*.**

179   Since these clusters are likely to have been horizontally transferred we attempted to find the

180   taxonomic origin of the transfer by using Kaiju (Menzel, Ng, and Krogh 2016) to find the

181   closest match for each gene from the CO-5 strains in the progenomes database (Mende et al.

182   2017), a database of representative microbial genomes that does not include *Xvv* or *Xvh*. As

183   a whole, over 93% of the CO-5 genome was effectively assigned to *Xanthomonas* sp., with

184   most genes assigned to *Xvm* (Supplementary Figure 7). In contrast, the gene clusters

185  contained sequences from different taxonomic groups as well as several unassigned

186  sequences. In Cluster A, eight genes (67%), were assigned to *Pantoea ananatis*

187  (Supplementary Figure 7), which is frequently isolated along with *Xvv* (Lang et al. 2017;

188  Coutinho et al. 2015). Within Cluster C, 40% of the genes were assigned to groups other than

189  *Xanthomonas* including species of *Sphingobium* and *Pseudomonas* (Supplementary Figure

190  7)*.

191  Meanwhile in Cluster E, 57% of the genes were assigned to *Xanthomonas* sp., but curiously,

192  8% of the genes matched phages in the Caudovirales group (Supplementary Figure 7)*. This

193  cluster was also enriched in GO terms associated to viral life cycles (Supplementary Figure

194  5), and unlike the other clusters, was not associated to insertion sequences (Figure 4)*, so its

195  acquisition could have been mediated by phage transmission.

196  The finding of phage-related sequences prompted us to scan the genomes for additional

197  phage sequences using PHASTER (Arndt et al. 2016) Indeed Cluster E corresponded to an

198  intact prophage that included a region larger than 30 kb containing a high percentage of

199  phage-related proteins. This was the only prophage identified in U.S. *Xvv* strains (Figure 5).

200  Both *Xvh* strains examined contained the prophage in a similar genomic location, but also

201  contained four other intact prophages. No prophage was identified in the South African *Xvv*

202  strains, and strains from sugarcane and *T. laxum* both contained a prophage different from

203  the one corresponding to Cluster E (Supplementary Figure 8). Many of the proteins in the

204  Cluster E prophage have similarities to proteins from the *Xanthomonas*-infecting phages Cp1

205  and Cp2 (Figure 5).

206  The Cluster E/prophage genes are found in all U.S. and Argentinian corn *Xvv*, and some

207  genes are found in one contemporary South African strain (Xvz45) (Figure 3). Older South

208 African strains do not contain prophage genes, and neither do *Xv, Xvm* or sugarcane *Xvv.* On

209 the other hand, all *Xvh* isolates contain these genes despite when isolated. A scenario that

210 could explain this distribution is that this region was acquired at some point in *Xvh* and was

211 recently horizontally transferred to an ancestor of the Argentinian and U.S. *Xvv* populations.

212 To explore this scenario we used Ranger-DTL (Bansal et al. 2018) to reconcile a whole

213 genome phylogenetic tree with gene trees for each ortholog group. In each reconciliation, the

214 more likely horizontal gene transfer (HGT) events and their direction were identified. When

215 analyzing all suitable ortholog trees (4739 in total), various possible exchanges where

216 identified, mostly within pathovars, and more abundantly within *Xvm* (Figure 6). As for genes

217 in Cluster E, the results suggested that indeed, a transfer from *Xvh* to the U.S. and Argentina

218 *Xvv* clade occurred. However, since many of the genes in this Cluster are identical across

219 strains (Supplementary Figure 9), it was not possible to establish a clear origin or destination

220 of the transfer. Similar results were obtained for Cluster E using ALE (amalgamated likelihood

221 estimation), another reconciliation technique, albeit with lower probabilities (Supplementary

222 Figure 10).

223 Overall, our results identified possible regions associated with the emerging *Xvv* population

224 infecting corn in the U.S. and Argentina. These regions potentially contain virulence

225 determinants or genes that conferred an advantage to this population for corn colonization in

226 a way that explains its rapid proliferation.

11

227 **Discussion**

228 In this work we show that the emerging populations of *Xvv* infecting corn in Argentina and the

229 U.S. are genetically related and have acquired genomic regions, specifically a prophage

230 (Cluster E) that may be associated with their spread. In phylogenetic analyses, the

231 Argentinian *Xvv* strains are closer to South African strains than U.S. strains, suggesting a

232 possible South American origin for the current epidemic. Accordingly, Argentinian *Xvv* strains

233 also appear connected to *Xvh* strains indicating the horizontal gene transfer of the prophage

234 from *Xvh* to *Xvv* could have occurred in South America (Figure 1). One U.S. strain (NE-7),

235 grouped closer to Argentinian strains, which is consistent with at least two separate

236 introductions to the U.S. Alternative scenarios are also possible since at least one

237 contemporaneous South African *Xvv* strain (Xvz45) contained some of the genes in the

238 cluster and another strain (XvzGP) grouped with U.S. strains.

239 The rapid spread of the disease in the U.S. and the possible ongoing genetic exchange

240 between these distant populations may have been facilitated by human activity. Corn

241 breeders in the U.S. accelerate product development by maintaining year-long operations and

242 have increasingly adopted the practice of using winter nurseries for breeding and seed

243 production (Butruille et al. 2015; Brewbaker 2003); many of these winter nurseries are located

244 in the southern hemisphere, including South America (Zaworski 2016; Butruille et al. 2015)

245 Additionally, corn production and export in South America has experienced considerable

246 growth in the last decades (Meade et al. 2016). Although it is still unknown whether *Xvv* is

247 transmitted by seeds, current practices likely allow enough exchange of contaminated

248 material such that an adapted population can spread quickly between continents.

249

12

250  When inoculated on corn, Argentinian and U.S. *Xvv* strains caused more severe symptoms

251  than South African strains, indicating that the emerging American populations are

252  phenotypically different than the older (1988) South African population. We were unable to

253  test contemporary South African strains because none are available in collections. Testing of

254  newer populations will be needed to determine if the current disease spread is associated

255  with an increase in virulence, since it is also possible that the tested South African strains

256  have reduced virulence due to their extended time in storage. Testing different corn

257  accessions will also be needed to confirm a gain of virulence since phenotypes seem to vary

258  between varieties (Lang et al. 2017).

259  We identified five clusters of genes that were over-represented in the U.S. *Xvv* strains. These

260  gene clusters overlapped with predicted genomic island regions, consistent with acquisition

261  through horizontal gene transfer. Several clusters may represent important genomic

262  acquisitions, if not for the current emerging population, for corn *Xvv* strains overall. Cluster C,

263  for instance, is a group of ~44 genes found in all contemporary corn *Xvv* strains (although not

264  all strains contain all genes) and some *Xvh* strains. This cluster is enriched in mobility genes:

265  transposases, insertion sequences, and various DNA binding genes, and it contains 5

266  predicted T3 secreted proteins. Taxonomic analyses revealed that a large percentage of

267  genes in this cluster match other groups of bacteria including *Pseudomonas,* S*phingobium*

268  and various *Burkholderiales.*

269  Cluster A contained eight genes found in *Pantoea ananatis,* including genes involved in

270  replication (replication proteins A and C) and conjugation (P-type conjugative transfer

271  proteins). *P. ananatis* is the causal agent of brown stalk rot of corn (Goszczynska et al. 2007)

272  but it is also a versatile organism able to infect monocotyledonous and dicotyledonous hosts,

13

273   and it is also a common epiphyte and endophyte (Coutinho and Venter 2009). *P. ananatis* was

274   documented in association with *Xvv* on *Eucalyptus* in S. Africa (Coutinho et al. 2015) and with

275   *Xvv* BLS symptomatic corn in the U.S. (Lang et al. 2017). However, the *Pantoea* strains alone

276   were unable to cause BLS symptoms in corn (Lang et al. 2017), and brown stalk rot

277   symptoms have not been reported on plants infected with BLS. The relationship between *Xvv*

278   and *P. ananatis* is intriguing and it is possible *Xvv* may have acquired important virulence

279   capacity from this association.

280   We focused on Cluster E because it is associated with the Argentina and U.S. *Xvv*

281   populations, and thus may be related to their emergence. The prophage region in Cluster E is

282   shared by *Xvh* and U.S. and Argentinian *Xvv,* and contained genes resembling elements of

283   *Xanthomonas*-infecting bacteriophages CP1, CP2 and Xp10. This prophage is absent in other

284   groups, including a recently reclassified available genome of a strain isolated from Areca nut

285   (Wicker et al, this issue) (Bradbury 1986) that is closely related to *Xvv* (absence of the

286   prophage was verified using PHASTER (Arndt et al. 2016)).

287   Prophages are temperate (non-infective or non-lytic) viruses that are integrated into bacterial

288   genomes by recombination (Varani et al. 2013). They are important vehicles for horizontal

289   gene transfer, they can promote recombination and rearrangements in the bacterial genome,

290   and they often carry additional non-essential cargo genes (morons) that may confer new

291   phenotypic properties to the bacteria (Varani et al. 2013; Brüssow, Canchaya, and Hardt

292   2004). Prophages have been known to carry virulence factors or factors that enhance

293   bacterial fitness (Brüssow, Canchaya, and Hardt 2004; Figueroa-Bossi et al. 2001), although

294   reduction of virulence has also been reported (Ahmad et al. 2014). Prophages harboring

295   elements conferring virulence activity have been found in different plant pathogenic bacteria

14

296   including *Xylella* sp*. (Varani et al. 2008)* and *Candidatus* Liberibacter asiaticus (Jain, Fleites,

297   and Gabriel 2015). And in *Xanthomonas arboricla*, strains pathogenic on walnut carry a higher

298   number (and also a different repertoire) of prophages than non-pathogenic strains (Cesbron

299   et al. 2015).

300   We hypothesize that the Cluster E prophage region contains genes that play a role in

301   virulence in *Xvh*, and when horizontally transferred to *Xvv,* it enhanced virulence or fitness to

302   the emerging *Xvv* populations. In *Xvh,* Cluster E was found in all examined strains, and was

303   one of only two shared pro-phages between two geographically and temporarily distant *Xvh*

304   strains analyzed (1961 Zimbabwe vs 2016 U.S.) (Figure 5, Supplementary Figure 8).

305   Furthermore most of the genes in this cluster were identical across all compared *Xvv* and *Xvh*

306   (Supplementary Figure 9), suggesting they are not subject to prophage decay (Brüssow,

307   Canchaya, and Hardt 2004) and may indeed play a beneficial role for the bacteria. This

308   cluster contains genes predicted to be T3-secreted, peptidases and transcription factors that

309   could have virulence activity, and various non-phage related hypothetical proteins with

310   unknown function. Further characterization of these genes, as well as of other over-

311   represented *Xvv* genes that were not assigned to clusters, is needed to establish a possible

312   role in virulence.

313   Here we have used comparative genomics to address questions about the origin of an

314   epidemic and the genetic determinants associated with pathogen population spread. Based

315   on our findings, we postulate different exciting hypotheses that will be the subject of future

316   work to understand the lifestyle and evolution of *Xvv* and related bacteria.

15

317  **Materials and Methods**

318  **Strain collection and molecular detection**

319  Isolation of *Xvv* from corn leaves was performed as in Lang et al. (2017) with minor

320  modifications. Instead of placing the tissue in distilled water, fresh tissues were dissolved in 1

321  mL of 10 mM MgCl$_2$, macerated with sterile pellet pestle and incubated for at least 1.5 hours

322  at room temperature. For bacterial isolation, one loop-fill (10 µL) of solution was spread onto

323  nutrient agar (NA). Plates were incubated at 28ºC for two days. Single characteristic bright

324  yellow colonies were selected, and re-streaked for further isolation until pure colonies were

325  obtained. Samples from the United States were collected across several fields in Colorado,

326  Iowa, Kansas and Nebraska. Samples from Argentina were collected from fields located in

327  San Luis, Córdoba, and Santa Fe states. (Supplementary Table 1).

328  South African *Xvv* strains where obtained from the L. E. Claflin collection (Qhobela, Claflin,

329  and Nowell 1990). Australian *Xvh* strains were obtained from the NSW Department of Primary

330  Industries Plant Pathology and Mycology Herbarium Culture Collection

331  (https://www.dpi.nsw.gov.au/about-us/services/collections/herbarium).

332  Molecular detection of *Xvv* was performed following one of these procedures: For the cases

333  using colony PCR, one single colony was suspended in 10 µL of sterile water and boiled at

334  95ºC for 5 min. First, colony PCR of suspected *Xvv* samples were performed using Xvv3 or

335  Xvv5 primers as described previously (Lang et al. 2017). To further confirm isolates, a second

336  method using 16S rRNA gene and a housekeeping gene, atpD (ATP synthase β chain), was

337  used to identify bacteria to species level. PCR reactions for 16S rRNA (50 µL) contained 2 µL

338  of boiled DNA template, 0.2 µM of each primer (Supplementary Table 4), 1X GoTaq reaction

339  buffer, 2 mM MgCl2, 0.2 mM dNTP, and 0.25 unit/µL GoTaq DNA polymerase enzyme

16

340   (Promega, Madison, WI). The cycling conditions were as follows: initial denaturation at 94ºC

341   for 3 min, following 35 cycles of 94ºC for 45 s, 50ºC for 1 min and 72ºC for 1:30 min, and the

342   final extension period at 72ºC for 10 min. PCR fragments were separated in a 1.5% agarose

343   gel for 45 min at 90 V, and fragments were extracted and purified using the DNA clean &

344   concentrator kit (ZYMO Research). Sequencing was performed with 5 ng/μL of each PCR

345   product at Quintara Biosciences (Fort Collins, CO) and analyzed using Geneious software

346   (version 10.0.7). Sequence identities to the genus level were determined using Blastn from

347   the NCBI database.

**Disease phenotyping on corn**

349   Corn (hybrid P1151) was grown in a 1:1 mix of Promix-BX Biofungicide + Mycorrhizae

350   (Quakertown, PA) under greenhouse conditions (30 ± 1 ºC, 16 hour day length, and 80%

351   relative humidity). Three weeks after planting, plants were inoculated with 24 *Xvv* isolates and

352   one *Xvh* isolate (Supplementary Table 1). Each bacterial strain was cultured in peptone

353   sucrose agar (PSA) for 24 hours at 28ºC and then suspended to $10^8$ CFU mL$^{-1}$ in sterile,

354   distilled water. Bacterial suspensions were infiltrated as described by Lang et al. 2017. Two

355   leaves were inoculated on at least seven individual plants. Infiltration experiments were

356   repeated three times and data was combined to perform statistical analysis. Sterile, distilled

357   water was used as a negative control in all inoculations. Quantification of the lesion length

358   was done by measuring the expansion distance beyond the infiltration site at seven days post

359   inoculation (dpi).

360   For statistical analysis a one-way ANOVA using lesion length ~ Isolate was made using the

361   aov function in R (R Core Team 2013), square root transformation of lesion length data was

362   done to satisfy ANOVA requirements. Treatment groups were obtained using a Tukey's HSD

17

363 (honestly significant difference) test, with the HSD.test in the agricolae package (de

364 Mendiburu and de Mendiburu 2019).

**Genome Sequencing, assembly and data collection**

366 Genomic DNA for the *Xanthomonas* positive samples was extracted using Easy-DNA kit

367 (Invitrogen) and PCR amplification of the *atpD* gene was carried out for further confirmation to

368 the species level. PCR reactions for *atpD* gene (40 µL) contained 25 ng/µL of DNA template,

369 0.4 µM of each primer (Supplementary Table 4), 1X GoTaq reaction buffer, 1.5 mM MgCl2, 0.2

370 mM dNTPs, and 0.1 unit/µL GoTaq DNA polymerase enzyme (Promega, Madison, WI).

371 Cycling conditions were performed as described by (Fargier, Saux, and Manceau 2011).

372 For 24 *Xvv* strains and one *Xvh* strain from the U.S., Illumina sequencing was performed by

373 BGI (www.bgi.com) using HiSeq 4000 with paired-end 100 bp reads. All Illumina reads were

374 first trimmed with Trimmomatic (PE ILLUMINACLIP:CO-2.adapters.fa LEADING:2

375 TRAILING:2 SLIDINGWINDOW:4:2 MINLEN:30) (Bolger, Lohse, and Usadel 2014) and then

376 assembled into scaffolds using SPAdes (Bankevich et al. 2012) with default settings. For 11

377 Australian *Xvh* strains sequencing was performed using Illumina Miseq and assembled using

378 the A5 pipeline (Coil, Jospin, and Darling 2015).

379 Five strains (CO-5, XV1601, NE744, Mex-1 and ZCP611) were sequenced using long read,

380 single molecule real time sequencing (SMRT Sequel, PacBio, Menlo Park, CA). SMRT read

381 sequences were assembled using HGAP v4 (Chin et al. 2013). Genomes were circularized

382 using circlator (Hunt et al. 2015). For XV1601, Illumina reads were also available and were

383 used to polish the PacBio assembly using Canu v1.3-r7616 (Koren et al. 2017). No major

384 differences were found between with the polished assembly nor with the other SMRT-

18

385 generated genomes as examined with multiple alignments with Mauve (Darling et al. 2004).

386 All generated genomes have been deposited to the NCBI (Supplementary Table 1)

387 We obtained all available assemblies of *X. vasicola* (NCBI:txid56459) and *X. campestris* pv.

388 *musacearum* (NCBI:txid454958, here referred to as *X. vasicola* pv. *musacearum*) as of

389 November 2018 (Supplementary Table 1). Assemblies with an N50 of minimum 10 kbp were

390 kept (thus excluding *Xvv* strains NCPPB895 and NCPPB890). For 10 recently published *Xvv*

391 strains form South Africa (named *Xanthomonas vasicola* pv. *zeae)* (Sanko et al. 2018) the

392 available assemblies ranged in size from 3.8 to 4.5 Mbp, significantly less than the average

393 size for other *X. vasicola* genomes (~4.9Mbp), and alignments to reference genomes

394 revealed large fragments missing from the assemblies (Supplementary Figure 11). These

395 genomes were thus reassembled from available Illumina raw reads (Biosample accessions

396 SAMN10286417-26) using Unicycler v0.4.8-beta (--mode bold), which functions as a SPAdes-

397 optimiser (Wick et al. 2017). New assemblies had expected sizes and were not missing large

398 regions and were thus kept for analysis in this work (Supplementary Figure 11).

**Genome annotation and ortholog identification**

400 All assemblies were automatically annotated using prokka v1.14-dev (--rfam) (Seemann

401 2014). Ortholog groups from prokka-annotated proteins were identified using Orthofinder v.

402 2.2.6 (default parameters) (Emms and Kelly 2015). Additionally, orthologs were also identified,

403 and core genome size was estimated, using Pan-x (Ding, Baumdicker, and Neher 2018)

404 Similar ortholog groups with similar distribution were found with both strategies (Pan-X= 6155

405 groups and 2163 unassigned genes, Orthofinder = 6084 groups and 1896 unassigned

406 genes), since Orthofinder grouped more genes together, these results were kept for further

407 analyses.

19

**Phylogenetic analyses**

Phylogenetic trees were used were obtained with various methods using whole genome data. KSNP3 (Gardner, Slezak, and Hall 2015) was used to obtain parsimony and maximum-likelihood trees based on pan-genome SNPs from identified k-mers (K-mer size =21), the parsimony tree had overall higher branch support and was kept for analysis. CSI Phylogeny 1.4 (Kaas et al. 2014) was used to obtain trees based on core genome SNPs, with *Xoo* PXO99A used as reference for SNP calling. Trees based on whole genome protein alignments obtained using the STAG method implemented in orthofinder (Emms and Kelly 2015) and RAxML+FasTree in Pan-X (Ding, Baumdicker, and Neher 2018; Price, Dehal, and Arkin 2010; Stamatakis 2014) were also analyzed.

MLSA neighbor-joining trees were obtained by identifying 31 housekeeping genes using AMPHORA v2 (Kerepesi, Bánky, and Grolmusz 2014), creating multiple alignments form their concatenated sequences using MUSCLE v3.8.31 (Edgar 2004), and generating the trees using functions of the R package phangorn (pml, optim.pml (model = "Blosum62") and bootstrap.pml(bs=100)) (Schliep 2011). Average nucleotide identity values were obtained using the ANI-matrix script from the enveomics collection (v1.3) (Rodriguez-R and Konstantinidis 2016).

Minimum spanning trees were generated with MSTGold v2.5 (Salipante and Hall 2011) using a multiple alignment of core genome SNPs identified with CSI Phylogeny (Kaas et al. 2014), a consensus tree of eight estimated different MSTs out of maximum 3000 tested was kept and bootstraped 500 times (parameters= -n 3000 -m 43200 -b 10 -t 50 -s 500).

RangerDTL was used to explore reconciliations between species trees and gene trees for all identified orthologs using the dated method (Bansal et al. 2018). Gene and species trees

431  used were generated by orthofinder (Emms and Kelly 2015); the species tree was made

432  ultrametric for this analysis using the chronos function of the ape package ((Paradis, Claude,

433  and Strimmer 2004). For each gene, 100 trees (using variable –seed from 1 to 100) were

434  reconciled with the species tree and possible horizontal gene transfer (HGT) were identified

435  with a probability corresponding to the number of trees were the event was identified. To

436  analyze multiple genes simultaneously (as for Cluster E), the probabilities for each event in

437  each tree were averaged. Amalgamated likelihood estimation, ALE v0.5 (Szöllősi Gergely J.

438  et al. 2015) was also used to perform the same analysis with the same trees.

439  **Identification of over-represented regions**

440  A hypergeometric test was designed and applied to each ortholog group identified with

441  orthofinder to look for over or under-represented genes in corn U.S. *Xvv* strain. The test was

442  applied using the function phyper(q, m, n, k, lower.tail = TRUE, log.p = FALSE) in R, where for

443  each ortholog: q= strains in the U.S. *Xvv* group that contain the gene, m=total number of

444  strains in the U.S. *Xvv* group, n=number of strains in the comparison group, k=total strains

445  that contain the gene in both groups.

446  The test was applied for each gene in both directions, for over-representation (q – 1) and

447  under-representation, and the lowest *p* value was chosen (if the lowest *p* value was for the

448  under-representation test, it was multiplied by -1 to differentiate them). The test was applied

449  100 times for each gene, each time changing the comparison group by randomly selecting a

450  group of non-U.S. *Xvv* strains of a random size between 10 and 69 (total of non-U.S. *Xvv*

451  strains) (Supplementary Figure 4). The average *p* value of the 100 tests was taken, and a

452  correction for multiple testing (p.adjust function, method BH (Benjamini and Hochberg 1995)

453  in R) was applied to the *p* values obtained for all genes.

21

454    Genes with an absolute adjusted *p* value < 0.05 were considered as over or under-

455    represented in the U.S. *Xvv* group. The position of the selected genes in the genome of the

456    strain CO-5 was then used to establish clusters. Groups of more than 10 genes over-

457    represented genes found less than 5 kb from each other were considered a Cluster and

458    assigned a letter (A-E) according to their distance to the replication origin.

459    **Annotation of genomic regions**

460    For a more thorough annotation of genes in each clusters and to assess enrichment in

461    functional categories, protein sequences of the CO-5 strain were further annotated using

462    Blast2GO v5.2.5 (Conesa et al. 2005) by combining hits against the ncbi nr- protein database

463    (blast-p fast, e-value 0.01, number of hits 10), InterPro, Gene Onthology terms (GOs), and

464    KEGG enzyme codes (default parameters). Enrichment of GO terms was assessed for the

465    different groups using a hyper-geometric test as implemented in the GoFuncR package

466    (Grote 2018).

467    Genomic Islands were predicted using the IslandViewer 4 suite (Bertelli et al. 2017) Insertion

468    Sequences (IS) were identified using ISEScan (v1.6) (default parameters) (Xie and Tang

469    2017). And possible prophage were identified using PHASTER (Arndt et al. 2016). All three

470    analyses were made using prokka-annotated files for strains with complete genomes.

471    Known Type III (T3) effectors were identified by blastp (v. 2.6.0+, results were filtered keeping

472    hits with -evalue < 0.0001, >30% identity in >40% the query length) (Altschul et al. 1997) of

473    consensus effectors sequences obtained from http://xanthomonas.org/ against the protein

474    sequences obtained using Prokka. Novel T3 effectors were predicted using effective DB

475    (default parameters + plant model for Predotar) (Eichinger et al. 2016), results were filtered to

476    keep proteins with an EffectiveT3 (signal peptide) of minimum 0.9999, plus any additional

477    predictions with other methods.

478    The web version of Kaiju (Menzel, Ng, and Krogh 2016) was used to annotate possible

479    taxonomic origin of cluster genes against the progenomes database (default parameters)

480    (Mende et al. 2017).

481    **Visualization and other analyses**

482    Most figures were generated using R (R Core Team 2013). Phylogenetic trees were

483    generated using the ggtree package (Yu et al. 2017). Circular genome and genomic region

484    visualizations were generated using ggbio (Yin, Cook, and Lawrence 2012). Heatmaps were

485    generated using pheatmap (Kolde and Kolde 2015). Upset plot was generated using UpsetR

486    (Conway, Lex, and Gehlenborg 2017). Comparisons of genomic regions were made using

487    GenomicRanges (Lawrence et al. 2013). Genomic alignments were visualized using Mauve v

488    Jan-19-2018 (Darling et al. 2004).

489 **Acknowledgments**

## References

494 Ahmad, A. A., Askora, A., Kawasaki, T., Fujie, M., and Yamada, T. 2014. The filamentous
496 phage XacF1 causes loss of virulence in Xanthomonas axonopodis pv. citri, the causative
497 agent of citrus canker disease. Front Microbiol. 5:321.

498 Altschul, S. F., Madden, T. L., Schäffer, A. A., Zhang, J., Zhang, Z., Miller, W., et al. 1997.
499 Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
500 Nucleic Acids Res. 25:3389–3402.

501 Arndt, D., Grant, J. R., Marcu, A., Sajed, T., Pon, A., Liang, Y., et al. 2016. PHASTER: a
502 better, faster version of the PHAST phage search tool. Nucleic Acids Res. 44:W16–W21.

503 Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., et al. 2012.
504 SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing.
505 Journal of Computational Biology. 19:455–477.

506 Bansal, M. S., Kellis, M., Kordi, M., and Kundu, S. 2018. RANGER-DTL 2.0: rigorous
507 reconstruction of gene-family evolution by duplication, transfer and loss. Bioinformatics.
508 34:3214–3216.

509 Benjamini, Y., and Hochberg, Y. 1995. Controlling the false discovery rate: a practical and
510 powerful approach to multiple testing. Journal of the Royal statistical society: series B
511 (Methodological). 57:289–300.

512 Bertelli, C., Laird, M. R., Williams, K. P., Lau, B. Y., Hoad, G., Winsor, G. L., et al. 2017.
513 IslandViewer 4: expanded prediction of genomic islands for larger-scale datasets. Nucleic
514 Acids Res. 45:W30–W35.

515   Bolger, A. M., Lohse, M., and Usadel, B. 2014. Trimmomatic: a flexible trimmer for Illumina

516   sequence data. Bioinformatics. 30:2114–2120.

517   Bradbury, J. F. 1986. *Guide to plant pathogenic bacteria.* CAB international.

518   Brewbaker, J. L. 2003. *Corn production in the tropics: The Hawaii experience*. University of

519   Hawaii.

520   Broders, K. 2017. Status of bacterial leaf streak of corn in the United States. In *Proceedings*

521   *of the Integrated Crop Management Conference*, Iowa State University, Digital Press.

522   Available at: https://lib.dr.iastate.edu/icm/2017/proceedings/18/ [Accessed February 25,

523   2019].

524   Brüssow, H., Canchaya, C., and Hardt, W.-D. 2004. Phages and the Evolution of Bacterial

525   Pathogens: from Genomic Rearrangements to Lysogenic Conversion. Microbiol Mol Biol Rev.

526   68:560–602.

527   Butruille, D. V., Birru, F. H., Boerboom, M. L., Cargill, E. J., Davis, D. A., Dhungana, P., et al.

528   2015. Maize Breeding in the United States: Views from Within Monsanto. In *Plant Breeding*

529   *Reviews: Volume 39*, John Wiley & Sons, Ltd, p. 199–282.

530   Cesbron, S., Briand, M., Essakhi, S., Gironde, S., Boureau, T., Manceau, C., et al. 2015.

531   Comparative Genomics of Pathogenic and Nonpathogenic Strains of Xanthomonas arboricola

532   Unveil Molecular and Evolutionary Events Linked to Pathoadaptation. Front Plant Sci. 6:1126.

533   Chin, C.-S., Alexander, D. H., Marks, P., Klammer, A. A., Drake, J., Heiner, C., et al. 2013.

534   Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data.

535   Nature Methods. 10:563–569.

536  Coil, D., Jospin, G., and Darling, A. E. 2015. A5-miseq: an updated pipeline to assemble

537  microbial genomes from Illumina MiSeq data. Bioinformatics. 31:587–589.

538  Conesa, A., Götz, S., García-Gómez, J. M., Terol, J., Talón, M., and Robles, M. 2005.

539  Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics

540  research. Bioinformatics. 21:3674–3676.

541  Conway, J. R., Lex, A., and Gehlenborg, N. 2017. UpSetR: an R package for the visualization

542  of intersecting sets and their properties. Bioinformatics. 33:2938–2940.

543  Coutinho, T. A., and Wallis, F. M. 1991. Bacterial Streak Disease of Maize (Zea mays L.) in

544  South Africa. Journal of Phytopathology. 133:112–112.

545  Coutinho, T. A., and Venter, S. N. 2009. Pantoea ananatis: an unconventional plant pathogen.

546  Molecular Plant Pathology. 10:325–335.

547  Coutinho, T. A., Westhuizen, L. van der, Roux, J., McFarlane, S. A., and Venter, S. N. 2015.

548  Significant host jump of Xanthomonas vasicola from sugarcane to a Eucalyptus grandis clone

549  in South Africa. Plant Pathology. 64:576–581.

550  Darling, A. C. E., Mau, B., Blattner, F. R., and Perna, N. T. 2004. Mauve: Multiple Alignment of

551  Conserved Genomic Sequence With Rearrangements. Genome Res. 14:1394–1403.

552  Ding, W., Baumdicker, F., and Neher, R. A. 2018. panX: pan-genome analysis and

553  exploration. Nucleic Acids Res. 46:e5–e5.

554  Dyer, R. A. 1949. Botanical surveys and control of plant diseases. Farming in South Africa.

555  24:119–121.

556  Edgar, R. C. 2004. MUSCLE: multiple sequence alignment with high accuracy and high

557  throughput. Nucleic Acids Res. 32:1792–1797.

558  Eichinger, V., Nussbaumer, T., Platzer, A., Jehl, M.-A., Arnold, R., and Rattei, T. 2016.

559  EffectiveDB—updates and novel features for a better annotation of bacterial secreted proteins

560  and Type III, IV, VI secretion systems. Nucleic Acids Res. 44:D669–D674.

561  Emms, D. M., and Kelly, S. 2015. OrthoFinder: solving fundamental biases in whole genome

562  comparisons dramatically improves orthogroup inference accuracy. Genome Biology. 16:157.

563  Fargier, E., Saux, M. F.-L., and Manceau, C. 2011. A multilocus sequence analysis of

564  Xanthomonas campestris reveals a complex structure within crucifer-attacking pathovars of

565  this species. Systematic and Applied Microbiology. 34:156–165.

566  Figueroa-Bossi, N., Uzzau, S., Maloriol, D., and Bossi, L. 2001. Variable assortment of

567  prophages provides a transferable repertoire of pathogenic determinants in Salmonella.

568  Molecular Microbiology. 39:260–272.

569  Gardner, S. N., Slezak, T., and Hall, B. G. 2015. kSNP3.0: SNP detection and phylogenetic

570  analysis of genomes without genome alignment or reference genome. Bioinformatics.

571  31:2877–2878.

572  Goszczynska, T., Botha, W. J., Venter, S. N., and Coutinho, T. A. 2007. Isolation and

573  Identification of the Causal Agent of Brown Stalk Rot, A New Disease of Maize in South Africa.

574  Plant Disease. 91:711–718.

575  Grote, S. 2018. GOfuncR: gene ontology enrichment using FUNC. R package version. 1.

576  Hunt, M., Silva, N. D., Otto, T. D., Parkhill, J., Keane, J. A., and Harris, S. R. 2015. Circlator:

577  automated circularization of genome assemblies using long sequencing reads. Genome

578  Biology. 16:294.

579  Jain, M., Fleites, L. A., and Gabriel, D. W. 2015. Prophage-Encoded Peroxidase in

580  'Candidatus Liberibacter asiaticus' Is a Secreted Effector That Suppresses Plant Defenses.

581  MPMI. 28:1330–1337.

582  Kaas, R. S., Leekitcharoenphon, P., Aarestrup, F. M., and Lund, O. 2014. Solving the Problem

583  of Comparing Whole Bacterial Genomes across Different Sequencing Platforms. PLOS ONE.

584  9:e104984.

585  Kerepesi, C., Bánky, D., and Grolmusz, V. 2014. AmphoraNet: the webserver implementation

586  of the AMPHORA2 metagenomic workflow suite. Gene. 533:538–540.

587  Kolde, R., and Kolde, M. R. 2015. Package 'pheatmap.' R Package. 1.

588  Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H., and Phillippy, A. M. 2017.

589  Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat

590  separation. Genome Res. 27:722–736.

591  Korus, K., Lang, J. M., Adesemoye, A. O., Block, C. C., Pal, N., Leach, J. E., et al. 2017. First

592  Report of *Xanthomonas vasicola* Causing Bacterial Leaf Streak on Corn in the United States.

593  Plant Disease. 101:1030.

594  Lang, J. M., DuCharme, E., Ibarra Caballero, J., Luna, E., Hartman, T., Ortiz-Castro, M., et al.

595  2017. Detection and Characterization of Xanthomonas vasicola pv. vasculorum (Cobb 1894)

596  comb. nov. Causing Bacterial Leaf Streak of Corn in the United States. Phytopathology.

597  107:1312–1321.

598    Lawrence, M., Huber, W., Pages, H., Aboyoun, P., Carlson, M., Gentleman, R., et al. 2013.

599    Software for computing and annotating genomic ranges. PLoS computational biology.

600    9:e1003118.

601    Leite, R. P., Custódio, A. a. P., Madalosso, T., Robaina, R. R., Duin, I. M., and Sugahara, V. H.

602    2019. First Report of the Occurrence of Bacterial Leaf Streak of Corn Caused by

603    Xanthomonas vasicola pv. vasculorum in Brazil. Plant Disease. 103:145–145.

604    Meade, B., Puricelli, E., McBride, W. D., Valdes, C., Hoffman, L., Foreman, L., et al. 2016.

605    *Corn and Soybean Production Costs and Export Competitiveness in Argentina, Brazil, and the*

606    *United States*. United States Department of Agriculture, Economic Research Service.

607    Available at: https://ideas.repec.org/p/ags/uersib/262143.html.

608    Mende, D. R., Letunic, I., Huerta-Cepas, J., Li, S. S., Forslund, K., Sunagawa, S., et al. 2017.

609    proGenomes: a resource for consistent functional and taxonomic annotations of prokaryotic

610    genomes. Nucleic Acids Res. 45:D529–D534.

611    de Mendiburu, F., and de Mendiburu, M. F. 2019. Package 'agricolae.' R Package, Version.

612    :1.2-1.

613    Menzel, P., Ng, K. L., and Krogh, A. 2016. Fast and sensitive taxonomic classification for

614    metagenomics with Kaiju. Nature Communications. 7:11257.

615    Moffett, M. L. 1983. Bacterial plant pathogens recorded in Australia. In *Plant Bacterial*

616    *Diseases: A Diagnostic Guide.*, Academic Press, Sydney., p. 317–336.

617    Paradis, E., Claude, J., and Strimmer, K. 2004. APE: Analyses of Phylogenetics and Evolution

618    in R language. Bioinformatics. 20:289–290.

619    Péros, J. P., Girard, J. C., Lombard, H., Janse, J. D., and Berthier, Y. 1994. Variability of
620    Xanthomonas Campestris pv. vasculorum From Sugarcane and Other Gramineae in Reunion
621    Island. Characterization of a Different Xanthomonad. Journal of Phytopathology. 142:177–
622    188.

623    Plazas, M. C., De Rossi, R. L., Brücher, E., Guerra, F. A., Vilaró, M., Guerra, G. D., et al.
624    2017. First Report of Xanthomonas vasicola pv. vasculorum Causing Bacteria Leaf Streak of
625    Maize (Zea mays) in Argentina. Plant Disease. 102:1026–1026.

626    Price, M. N., Dehal, P. S., and Arkin, A. P. 2010. FastTree 2 – Approximately Maximum-
627    Likelihood Trees for Large Alignments. PLOS ONE. 5:e9490.

628    Qhobela, M., Claflin, L. E., and Nowell, D. C. 1990. Evidence that Xanthomonas campestris
629    pv. zeae can be distinguished from other pathovars capable of infecting maize by restriction
630    fragment length polymorphism of genomic DNA. Canadian Journal of Plant Pathology.
631    12:183–186.

632    R Core Team. 2013. *R: A language and environment for statistical computing*.

633    Rodriguez-R, L. M., and Konstantinidis, K. T. 2016. *The enveomics collection: a toolbox for*
634    *specialized analyses of microbial genomes and metagenomes*. PeerJ Preprints. Available at:
635    https://peerj.com/preprints/1900/ [Accessed February 7, 2017].

636    Salipante, S. J., and Hall, B. G. 2011. Inadequacies of Minimum Spanning Trees in Molecular
637    Epidemiology. Journal of Clinical Microbiology. 49:3568–3575.

638    Sanko, T. J., Kraemer, A. S., Niemann, N., Gupta, A. K., Flett, B. C., Mienie, C., et al. 2018.
639    Draft Genome Assemblages of 10 Xanthomonas vasicola pv. zeae Strains, Pathogens
640    Causing Leaf Streak Disease of Maize in South Africa. Genome Announc. 6:e00532-18.

641    Schliep, K. P. 2011. phangorn: phylogenetic analysis in R. Bioinformatics. 27:592–593.

642    Seemann, T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics. 30:2068–

643    2069.

644    Stamatakis, A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of

645    large phylogenies. Bioinformatics. 30:1312–1313.

646    Szöllősi Gergely J., Davín Adrián Arellano, Tannier Eric, Daubin Vincent, and Boussau

647    Bastien. 2015. Genome-scale phylogenetic analysis finds extensive gene transfer among

648    fungi. Philosophical Transactions of the Royal Society B: Biological Sciences. 370:20140335.

649    Tushemereirwe, W., Kangire, A., Ssekiwoko, F., Offord, L. C., Crozier, J., Boa, E., et al. 2004.

650    First report of Xanthomonas campestris pv. musacearum on banana in Uganda. Plant

651    Pathology. 53:802–802.

652    USDA-NASS. 2017. Crop production Summary 2016, United States Department of

653    Agriculture, National Agricultural Statistics Service. Washington, D.C.□: United States

654    Department of Agriculture, Statistical Reporting Service, Crop Reporting Board□: [Supt. of

655    Docs., U.S. G.P.O., distributor]. Available at: http://purl.access.gpo.gov/GPO/LPS1137.

656    Varani, A. de M., Souza, R. C., Nakaya, H. I., Lima, W. C. de, Almeida, L. G. P. de, Kitajima,

657    E. W., et al. 2008. Origins of the Xylella fastidiosa Prophage-Like Regions and Their Impact in

658    Genome Differentiation. PLOS ONE. 3:e4059.

659    Varani, A. M., Monteiro-Vitorello, C. B., Nakaya, H. I., and Van Sluys, M.-A. 2013. The Role of

660    Prophage in Plant-Pathogenic Bacteria. Annual Review of Phytopathology. 51:429–451.

661    White, F. F., Potnis, N., Jones, J. B., and Koebnik, R. 2009. The type III effectors of
662    *Xanthomonas*. Molecular Plant Pathology. 10:749–766.

663    Wick, R. R., Judd, L. M., Gorrie, C. L., and Holt, K. E. 2017. Unicycler: Resolving bacterial
664    genome assemblies from short and long sequencing reads. PLOS Computational Biology.
665    13:e1005595.

666    Xie, Z., and Tang, H. 2017. ISEScan: automated identification of insertion sequence elements
667    in prokaryotic genomes. Bioinformatics. 33:3340–3347.

668    Yin, T., Cook, D., and Lawrence, M. 2012. ggbio: an R package for extending the grammar of
669    graphics for genomic data. Genome biology. 13:R77.

670    Young, J. M., Bradbury, J. F., Davis, R. E., Dickey, R. S., Ercolani, G. L., Hayward, A. C., et al.
671    1991. Nomenclatural revisions of plant pathogenic bacteria and list of names 1980-1988.
672    Review of Plant Pathology. 70:211–221.

673    Yu, G., Smith, D. K., Zhu, H., Guan, Y., and Lam, T. T.-Y. 2017. ggtree: an R package for
674    visualization and annotation of phylogenetic trees with their covariates and other associated
675    data. Methods in Ecology and Evolution. 8:28–36.

676    Zaworski, F. 2016. Winter Breeding Programs in South America. SeedWorld. Available at:
677    https://seedworld.com/winter-breeding-programs-south-america/ [Accessed February 21,
678    2019].

679 **Figure Legends**

680 **Figure 1. Phylogeny of *X. vasicola* strains.** A) Parsimony tree based on pan-genome SNPs

681 from 91 draft and fully sequenced *X. vasicola* genomes built using kSNP3 (Gardner, Slezak,

682 and Hall 2015). Four main pathovars/groups are indicated by solid colored lines. Colors in

683 tree tips indicate country of isolation and tip letters indicate the plant host of the isolate. Bar

684 shows the tree scale. Dotted line to the out-group (*X. oryzae* pv. *oryzae*) indicates this

685 distance was scaled down tenfold to improve visualization. Colors in the tip letters (black or

686 white) are for readability and do not indicate any feature. 74% of the nodes had support over

687 70% (as calculated by kSNP3 (Gardner, Slezak, and Hall 2015). B) Consensus minimum

688 spanning tree based on core genome SNPs with PXO99A as a reference. Circle colors

689 indicate seven groups of interest of *X. vasicola.* The consensus tree was bootstrapped 500

690 times and edge colors indicate bootstrap percentages.

691 **Figure 2. Phenotypic characterization of *Xvv* strains.** Disease caused by *Xvv* and *Xvh* on

692 corn (hybrid P1151). Three week-old plants were infiltrated with $10^8$ CFU mL$^{-1}$ of each isolate,

693 and disease was assessed at seven days post inoculation (dpi). Lesion lengths indicate

694 expansion beyond the infiltration site. The experiment was replicated three times and

695 combined data from all replications is shown here. Letters designate significance groups at *p*

696 value <0.0001 using one-way ANOVA+Tukey's HSD using square root transformation of data

697 and sample size of at least seven replicates per isolate (90% statistical power).

698 **Figure 3. Over and under-represented genes in U.S. *Xvv*.** Heatmap shows presence (white

699 or colored) or absence (grey) of all ortholog groups found in at least 10 strains in the genome

700 set, each vertical line represents an ortholog group (4912 total). A hypergeometric test was

701 applied to each group to determine over or under representation in the group of U.S. *Xvv*

34

702  strains versus all other groups. Colors in the heatmap (blue to red) indicate 1 - the adjusted *p*

703  value for these tests. Negative values indicate the result of an under-representation test

704  (adjusted *p* value * -1). The vertical lines in the heatmap are ordered according to the

705  presence of each ortholog group in the genome of the U.S. *Xvv* strain CO-5, numbers below

706  the heatmap (1M to 5M) indicate the position (in million base pairs) of the genes in the

707  reference genome. Vertical lines after the 5M mark are orthologs not present in CO-5 ordered

708  according to their frequency in other genomes. Bars at the left of the heatmap indicate the

709  group of the strain as in Figure 1B. The dendrogram to the left corresponds to a MLSA tree

710  based on all orthologs (Supplementary Figure 1). Five gene clusters of over-represented

711  genes (A-E) were identified and indicated with letters below the heatmap, these clusters

712  defined as groups of at least 10 genes found less than 5 Kb from each other in the CO-5

713  genome and with an adjusted *p* value for over-representation <0.05 (0.95 in the figure).

714  **Figure 4. Genomic islands, predicted type 3 effectors and insertion sequences in *X.***

715  ***vasicola* genomes.** Circular representation of eight complete *X. vasicola* genomes with

716  annotated regions of interest. Legend in black square at the bottom right indicates what each

717  circle represents. Genomic scale in million base pairs is shown. The outermost circles show

718  presence of annotated genes in each genome. Clusters of genes identified as over-

719  represented in US-*Xvv*, and their orthologues, are shown in colors. Predicted Genomic

720  Islands using three methods integrated in Island Viewer (Bertelli et al. 2017) are shown in red

721  colors, multiple predictions for a same region are shown stacked, SIGI-HMM is based on

722  sequence composition, IslandPath-DIMOB is based on dinucleotide bias and presence of

723  mobility genes, and IslandPick is based on phylogenetic comparisons. Predictions of type 3

724  secreted proteins are shown in green colors, four methods integrated in the EffectiveDB

725   (Eichinger et al. 2016) suite are shown: EffectiveT3 predicts Type 3 secretion signals,

726   Effective CCBD, conserved binding domains of Type 3 chaperones, EffectiveELD, eukaryotic-

727   like domains, and Predotar predicts plant subcellular localization. Proteins having a significant

728   score with at least EffectiveT3 are shown. Insertion sequence (IS) elements are shown in

729   black as identified using ISEScan (Xie and Tang 2017).

730   **Figure 5. Prophages in *X. vasicola* genomes. A).** Circular representation of identified

731   prophages in eight complete *X. vasicola* genomes, genomic scale in million base pairs is

732   shown. Intact prophage regions as identified by PHASTER (Arndt et al. 2016) are shown in

733   blue and named P+number according to their position. Incomplete or questionable regions

734   are shown in green (regions that are close together are shown stacked to improve

735   readability). **B)** Diagram showing the genes in the predicted prophage corresponding to

736   Cluster E in genomes of *Xvh* and *Xvv.* The genes are grouped according to their strand and

737   colored according to their annotation. The diagram for *Xvv* strain CO-5 P1 shows for each

738   gene the top hit against known genes in the Virus and Prophage database of PHASTER

739   (Arndt et al. 2016), genes with no annotation had no significant hits.

740   **Figure 6. Horizontal Gene Transfer events predicted between *X. vasicola* genomes.** A

741   strain tree based on all core genome orthologs and ortholog gene trees obtained with

742   orthofinder (Emms and Kelly 2015) were reconciled using Ranger-DTL (Bansal et al. 2018) to

743   identify possible horizontal gene transfer events. Tips of the trees indicate country and host of

744   isolation for each genome and branch colors indicate the four main *X. vasicola* pathovars.

745   Tree to the left shows the results for all ortholog trees, and to the right for genes assigned to

746   U.S. *Xvv* cluster E. Arrow thickness and color indicate predicted cumulative frequency of each

747   event, a frequency of one would mean an event was identified in all 100 evaluated reconciled

748    trees for all genes analyzed. Top 10 events with highest probability for each set are shown.

749    Arrow heads indicate the direction of the predicted event.

750   **Supplementary Materials Legends.**

751   **Supplementary Figure 1. Phylogeny of *X. vasicola* strains obtained using different**

752   **methods.** Trees are shown using different methods: Orthofinder and Pan-X build trees based

753   on protein sequences of core genome trees using the STAG method and RaxML+FasTree

754   respectively. CSI phylogeny builds trees based on core genome SNPs, and the MLSA tree

755   was generated based on concatenated sequences of housekeeping genes identified with

756   AMPHORA, aligned with Muscle. The tree was built by using the R package phangorn. When

757   present, the trees are rooted using *X. oryzae* pv. *oryzae* (*Xoo*) PXO99A as an outgroup,

758   otherwise the tree was rooted using *Xvh* NCPPB 1060. *Xoo* PXO99A was excluded from pan-

759   X analyses. Bars to the left indicate branch length as generated by each program.

760   **Supplementary Figure 2**. **Average Nucleotide Identity (ANI) between pairs of *X. vasicola***

761   **genomes.** The heatmap shows pairwise ANI values between *X. vasicola* genomes, with *Xoo*

762   PXO99A included for comparison. Dendrograms to the top and left show hierarchical

763   clustering of genomes based on ANI. Bars to the left indicate Host, Pathovar and Country of

764   isolation.

765   **Supplementary Figure 3. Shared orthologs between different *X. vasicola* groups. A)**

766   UpSet visualization of intersections between orthologs present in each relevant *X. vasicola*

767   group. Orthologs were identified using orthofinder in each genome, an ortholog group was

768   said to be present in a group if it was present in at least 30% of the strains evaluated. Vertical

769   bars show the intersection between groups with bold circles below. First bar corresponds to

770   the intersection of all groups (core genome). Horizontal bars indicate the number of genes

771   found in each group. Highlighted in blue is the intersection between corn *Xvv* from the U.S.

772   and Argentina with *Xvh*, and highlighted in purple are genes exclusive to corn *Xvv*. **B)** Core

773 genome statistics obtained from Pan-X. Percentage of core and accessory genes is shown

774 (left), then number of strains containing groups of orthologs (middle) and the distribution of

775 gene length in all genomes (right).

776 **Supplementary Figure 4. Identification of over-represented genes in U.S. Xvv**. A) The

777 frequency of each gene (each point) in the U.S. *Xv* population (x axis) is compared to their

778 frequency in sets of randomly chosen *X. vasicola* genomes (y axis). The average frequency of

779 each gene in 100 groups is shown and error bars indicate standard deviation. Dot colors

780 indicate whether a given gene was identified as over or under-represented. B) Density plot

781 showing uniform size distribution for random sets of genomes (100 per gene) chosen as

782 comparison groups for hypergeometric tests to determine over or under representation when

783 compared to U.S. *Xvv* genomes C) Density plot shows the pathovar composition of the

784 random sets.

785 **Supplementary Figure 5. Gene Ontology (GO) term enrichment in over-represented**

786 **U.S. *Xvv* genes**. Go terms identified as statistically enriched in the group of over-represented

787 genes and their genomic clusters are shown. No terms were found enriched for Clusters A, B

788 or D. Dot color and size indicate *p* value of enrichment as determined using GoFuncR. GO

789 annotations were obtained using Blast2GO.

790 **Supplementary Figure 6. Known Type 3 effectors in *X. vasicola* genomes.** Heatmap

791 shows copy number of known T3 effectors as determined by blast of each genome against

792 consensus *Xanthomonas* T3 effector sequences. Copy number is shown to a maximum of 3.

793 The only effector with a higher copy number is AvrBs3 (TAL effectors) in *Xoo* PXO99A.

794 Dendrogram at the top corresponds to the KSNP3 tree in Figure 1. Dendrogram to the left

795 shows hierarchical clustering of effectors according to their presence/absence pattern in the

796　genomes. Color bars at the top indicate Pathovar, Host and Country of isolation for each

797　genome.

798　**Supplementary Figure 7. Taxonomic distribution of over-represented genes in possible**

799　**horizontally transferred clusters in U.S. *Xvv*.** Krona plots obtained using Kaiju showing the

800　taxonomic assignation of genes in each over-represented U.S. *Xvv* cluster in the strain CO-5

801　as well as in the whole genome. Each gene was matched to its closest sequence in the

802　progenomes database and assigned a taxonomic group accordingly. Colors in each plot are

803　ordered according to percentages and do not correspond to the same taxa across clusters.

804　**Supplementary Figure 8. Other prophages identified in *X. vasicola* genomes.** Diagrams

805　show genes found in the predicted prophages in *X. vasicola* genomes different from the

806　prophage corresponding to Cluster E. Gare grouped according to their strand and colored

807　according to their annotation in phaster. For each gene the top hit against known genes in the

808　Virus and Prophage database of PHASTER is shown; genes with no annotation had no

809　significant hits.

810　**Supplementary Figure 9. Genetic distances of genes in ortholog groups assigned to**

811　**over-represented clusters in U.S. *Xvv*.** Phylogenetic trees obtained with orthofinder were

812　analyzed to find the distances between all tips in the tree (strains containing each gene) using

813　the cophenetic function from the ape package. Boxplots show the distribution of distances for

814　each tree, boxplots showing means around zero indicate that all the tips were found at the

815　same distance, meaning the gene sequence was identical across strains.

816　**Supplementary Figure 10. Horizontal Gene Transfer events predicted between *X.***

817　***vasicola* genomes using ALE.** A strain tree based on all core genome orthologs and

818　ortholog gene trees obtained with orthofinder were reconciled using ALE to identify possible

40

819  horizontal gene transfer events. Tips of the trees indicate country and host of isolation for

820  each genome and branch colors indicate the four main *X. vasicola* pathovars. Tree to the left

821  shows the results for all ortholog trees, and to the right for genes assigned to U.S. *Xvv* cluster

822  E. Arrow thickness and color indicate predicted cumulative frequency of each event, a

823  frequency of one would mean an event was identified in all 100 evaluated reconciled trees for

824  all genes analyzed. Arrow head indicate the direction of the predicted event.

825  **Supplementary Figure 11**. **Reassembly of Xvv strains from South Africa.** Mauve multiple

826  alignment shows the south African *Xvv* strain SAM119 (top), the publicly available assembly

827  of strain Xvz45 (GCF_003111905.1) (middle), and a reassembly of strain Xvz45 used in this

828  paper using the corresponding raw reads (SAMN10286417) (bottom). Long vertical red lines

829  indicate contig limits. Similar results were obtained with other genomes form this set,

830  indicated with accession numbers SAMN- in Supplementary Table 1.
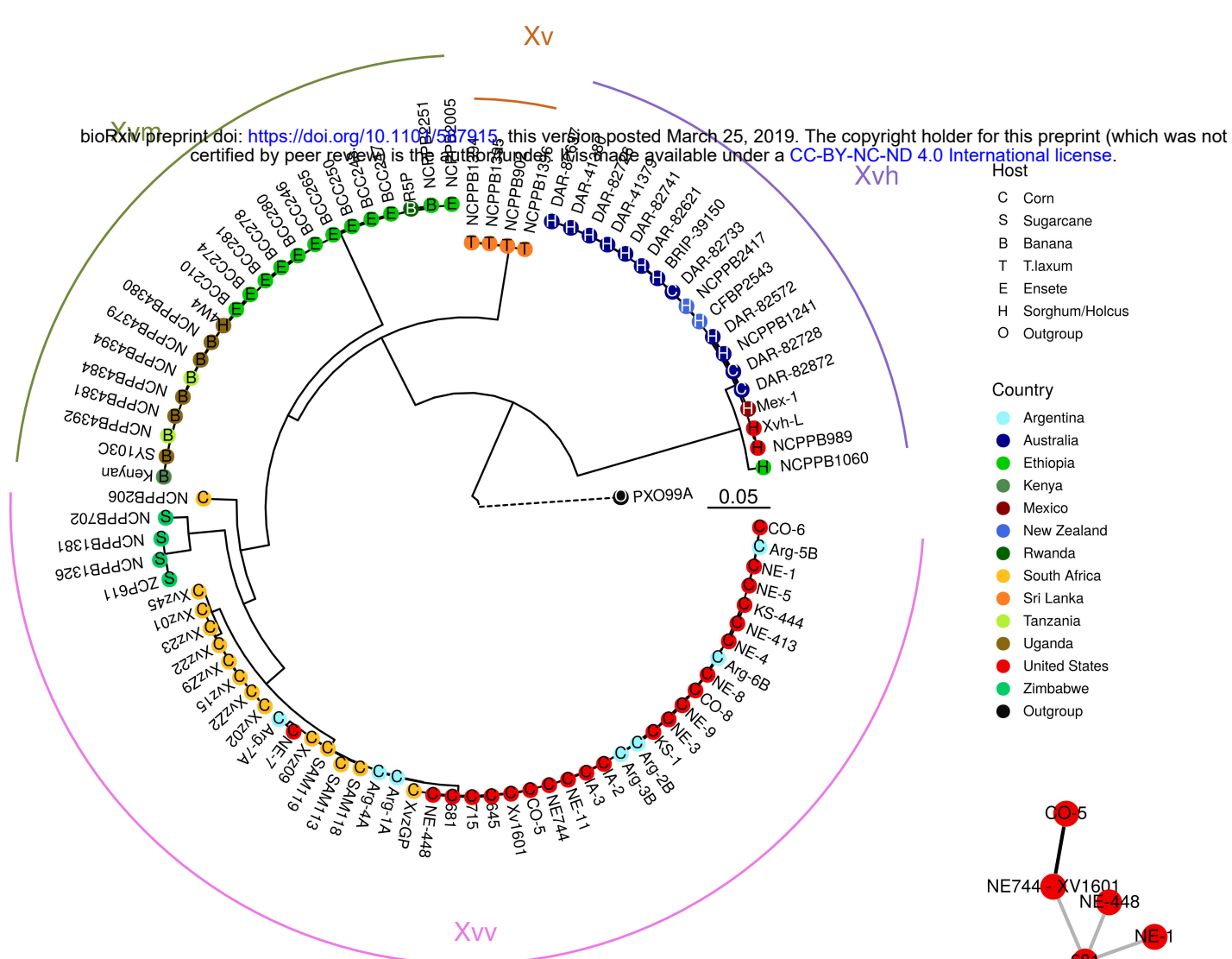
831  **Supplementary Table 1. Inventory of genomic sequences used in this work.**

832  **Supplementary Table 2. Ortholog groups determined as over or under-represented in**
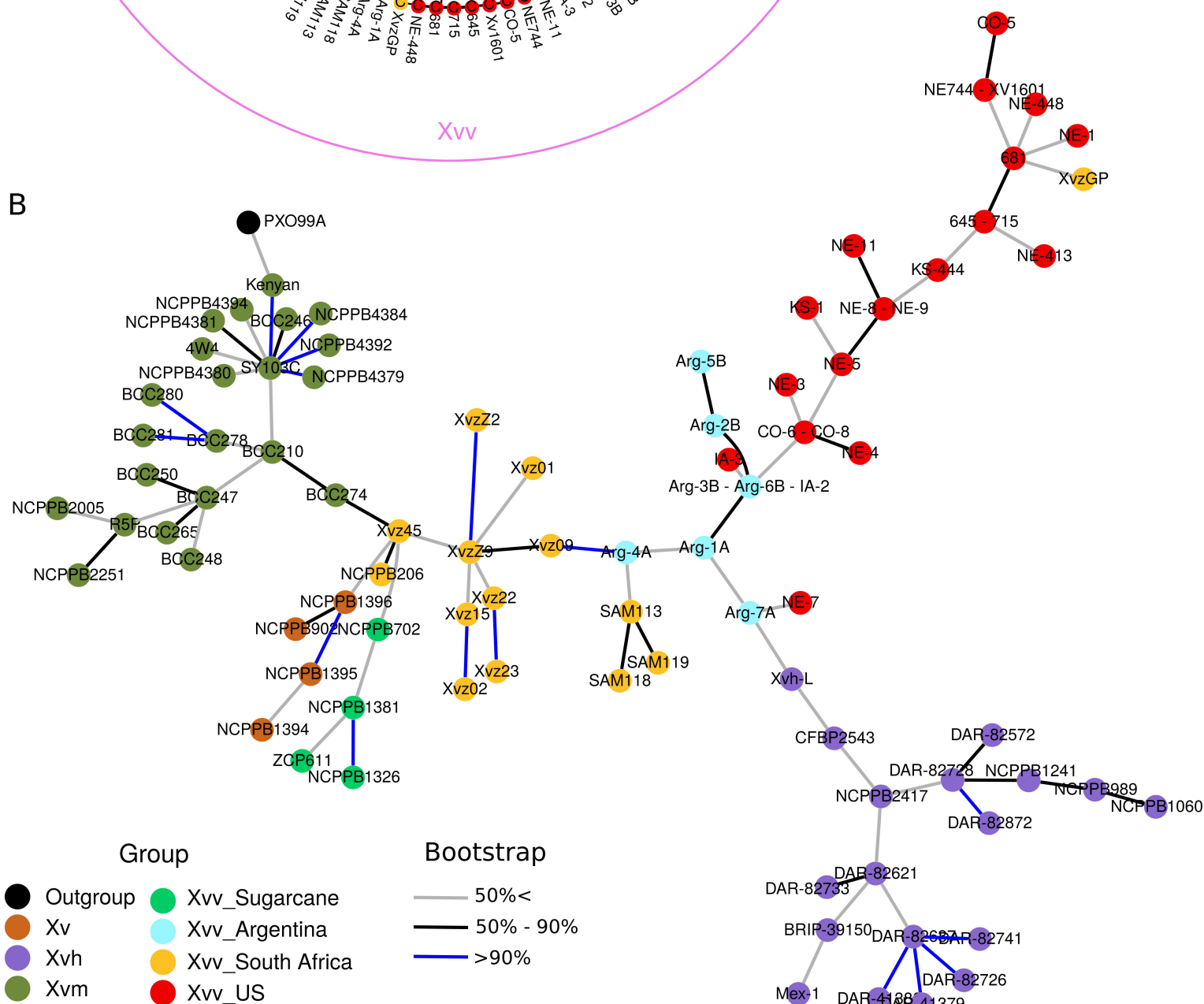
833  **U.S. *Xvv* strains**

834  **Supplementary Table 3. Annotation of genes assigned to over-represented clusters in**
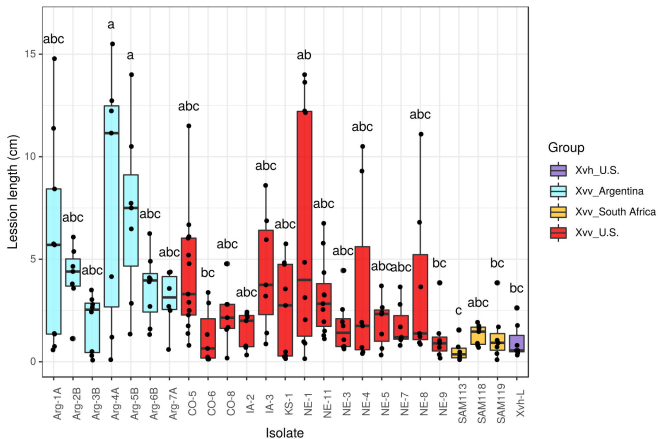
835  ***Xvv* strain CO-5.**

836  **Supplementary Table 4. Primers used for molecular detection of *Xvv.***

41

A

B

Annotation

- Cluster_A
- Cluster_B
- Cluster_C
- Cluster_D
- Cluster_E
- Method=IslandPath-DIMOB
- Method=IslandPick
- Method=SIGI-HMM
- Method=EffectiveT3
- Method=EffectiveCCBD
- Method=EffectiveELDs
- Method=Predotar
- Method=ISEScan

+ Strand
- Strand
Genomic Islands
T3SS
IS

**A**

US_Xvv C0-5

US_Xvv NE744

US_Xvv XV1601

SA_Xvv SAM119

Sugarcane_Xvv ZCP611

Xv NCPPB902

Xvh NCPPB1060

Xvh Mex-1

PHASTER prediction
- Incomplete/Questionable
- intact

**B**

**C0-5(US_Xvv) P1**

Entero_nrEp400

Xantho_Xp15   Pseudo_F10   Xantho_CP2
Xantho_CP1   Xantho_CP1   Shewan_1/44   Stenot_S1   Entero_c   Entero_c
Xantho_CP1
Stenot_S1   Xantho_CP1   Xantho_CP1   Xantho_CP2   Xantho_CP2   Stenot_S1   Stenot_S1   Entero_cdtI   Stenot_S1

4510000   4520000   4530000   4540000

**XV1601(US_Xvv) P1**

4510000   4520000   4530000   4540000

**NE744(US_Xvv) P1**

4510000   4520000   4530000   4540000

Annotation
- Attachment_site
- Head_protein
- Hypothetical_protein
- Integrase
- Phage-like_protein
- Portal_protein
- Tail_protein
- Terminase

**Mex-1(Xvh) P5**

4550000   4560000   4570000   4580000

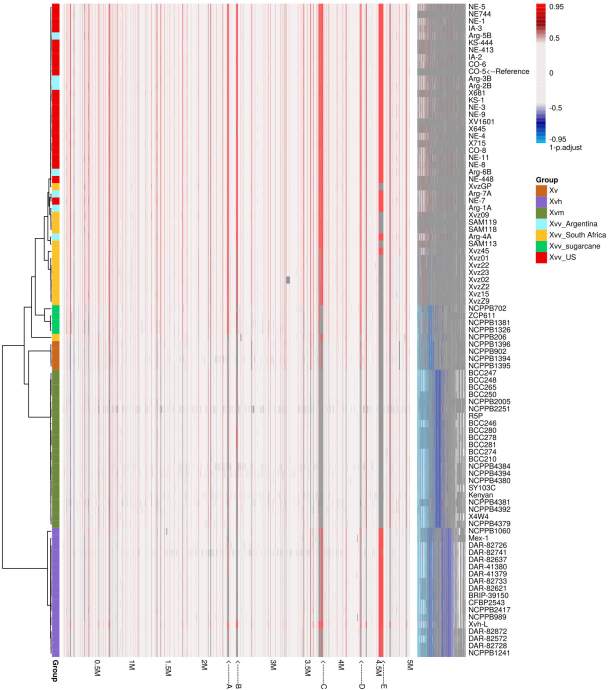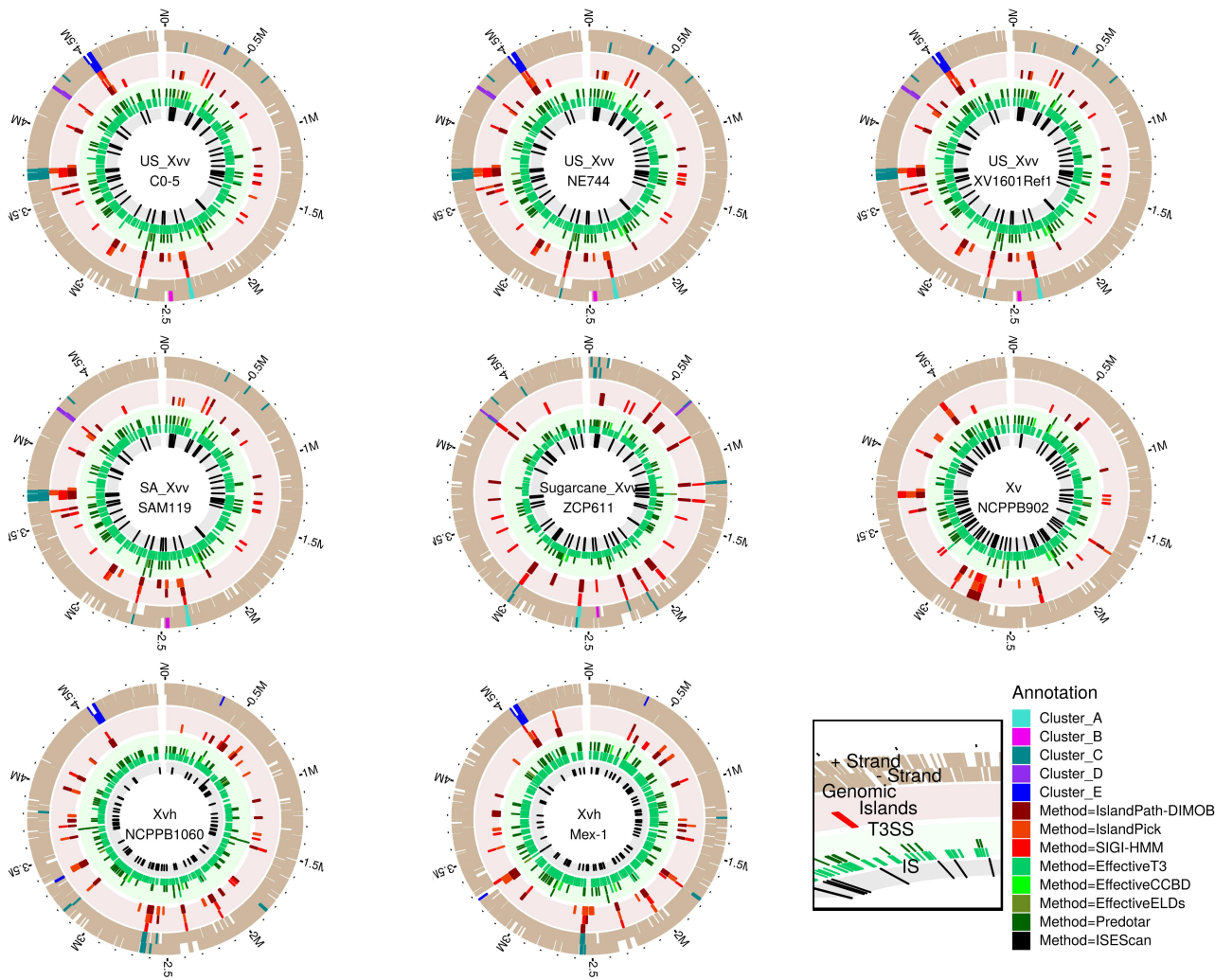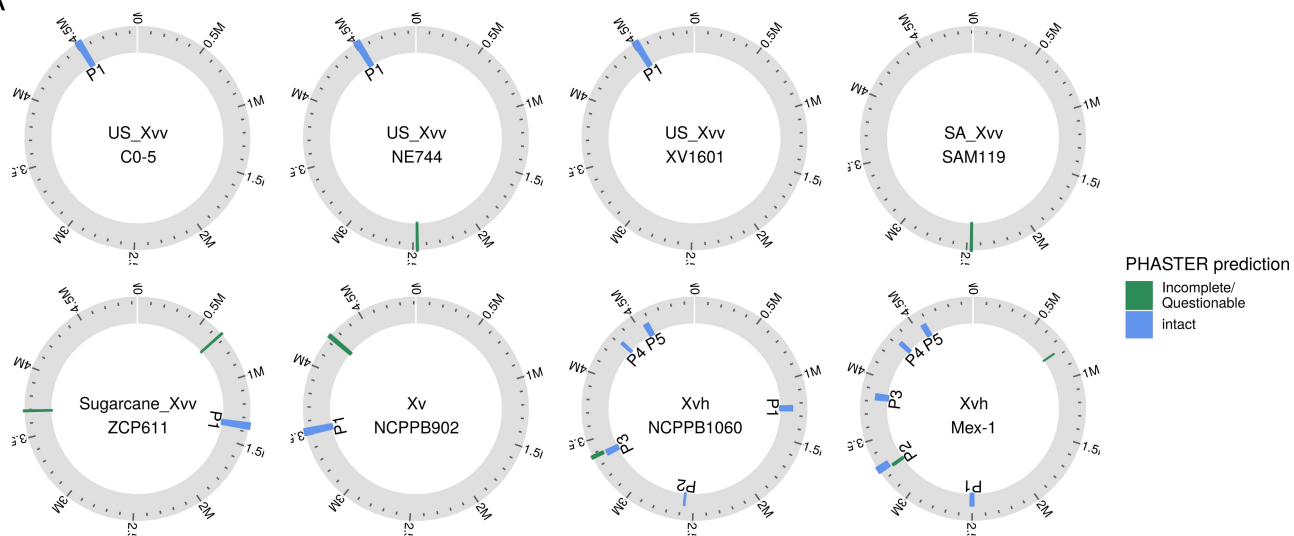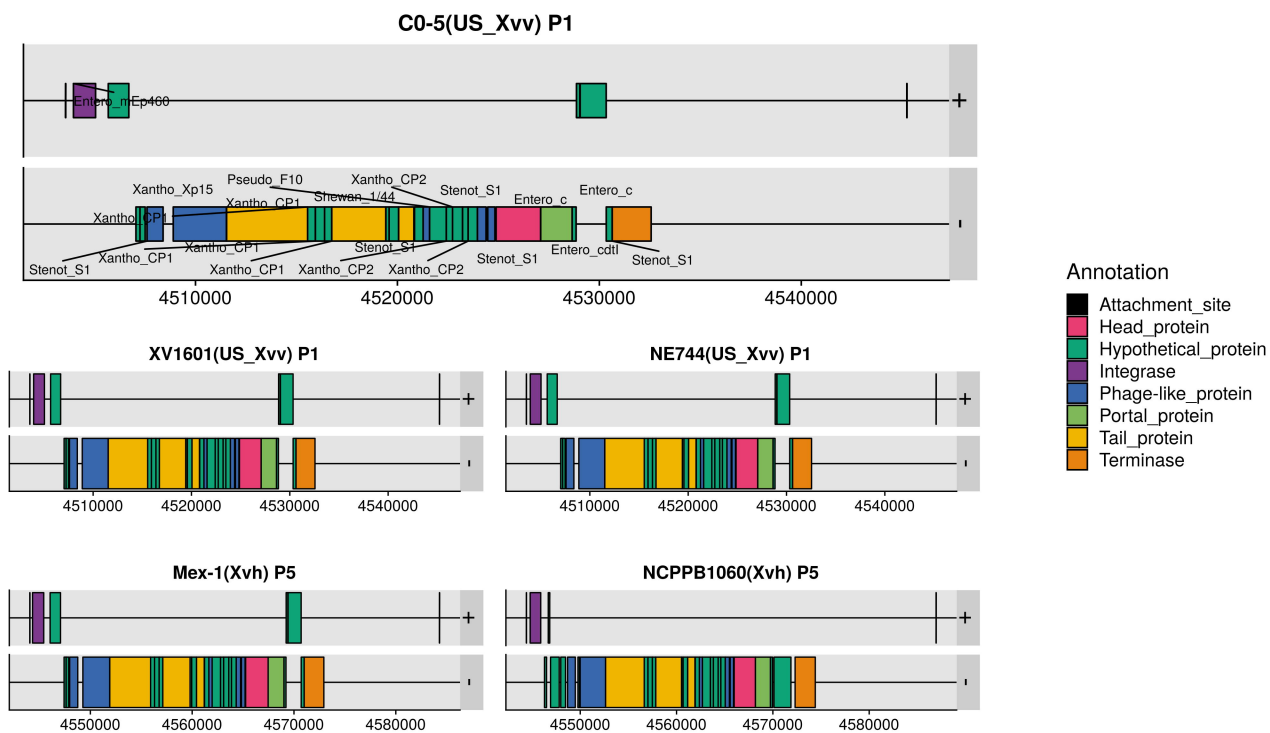**NCPPB1060(Xvh) P5**

4550000   4560000   4570000   4580000

All genes (4739)    Cluster_E genes (57)

Country
- C Argentina
- Australia
- Ethiopia
- Kenya
- Mexico
- New Zealand
- Rwanda
- South Africa
- Sri Lanka
- Tanzania
- Uganda
- United States
- Zimbabwe

Host
- C  Corn
- S  Sugarcane
- B  Banana
- T  T.laxum
- E  Ensete
- H  Sorghum/Holcus

HGT Frequency
0.3        0.08

Xvv

Xv

Xvm

Xvh