1 **Genomic characterization and computational phenotyping of nitrogen-fixing**
2 **bacteria isolated from Colombian sugarcane fields**

3

4 Luz K. Medina-Cordoba[a,b], Aroon T. Chande[a,b,c], Lavanya Rishishwar[a,b,c], Leonard W.
5 Mayer[b,c], Lina C. Valderrama-Aguirre[b,d], Augusto Valderrama-Aguirre[b,e,f], John Christian
6 Gaby[a,g], Joel E. Kostka[a,b]# and I. King Jordan[a,b,c]#

7

8 [a]School of Biological Sciences, Georgia Institute of Technology, Atlanta, Georgia, USA

9 [b]PanAmerican Bioinformatics Institute, Cali, Valle del Cauca, Colombia

10 [c]Applied Bioinformatics Laboratory, Atlanta, Georgia, USA

11 [d]Laboratory of Microorganismal Production (Bioinoculums), Department of Field
12 Research in Sugarcane, INCAUCA S.A.S., Cali, Valle del Cauca, Colombia

13 [e]Biomedical Research Institute (COL0082529), Cali, Valle del Cauca, Colombia

14 [f]Universidad Santiago de Cali, Cali, Colombia.

15 [g] Faculty of Chemistry, Biotechnology and Food Science, Norwegian University of Life
16 Sciences, Ås, Norway

17 Running Head: Computational phenotyping of nitrogen-fixing bacteria

18

19 #Address computational correspondence to king.jordan@biology.gatech.edu

20 *Present address: 950 Atlantic Dr NW, Atlanta, GA 30332

21

22 #Address microbiology correspondence to joel.kostka@biology.gatech.edu

23 *Present address: 310 Ferst Dr NW, Atlanta, GA 30332

24

25 **KEYWORDS:** biofertilizer, nitrogen fixation, plant growth promoter, genome
26 sequencing, computational phenotyping

27 **ABSTRACT** Previous studies have shown that the sugarcane microbiome harbors diverse

28 plant growth promoting (PGP) microorganisms, including nitrogen-fixing bacteria, and the

29 objective of this study was to design a genome-enabled approach to prioritize sugarcane

30 associated nitrogen-fixing bacteria according to their potential as biofertilizers. Using a

31 systematic high throughput approach, 22 pure cultures of nitrogen-fixing bacteria were isolated

32 and tested for diazotrophic potential by PCR amplification of nitrogenase (*nifH*) genes, common

33 molecular markers for nitrogen fixation capacity. Genome sequencing confirmed the presence

34 of intact nitrogenase *nifH* genes and operons in the genomes of 18 of the isolates. Isolate

35 genomes also encoded operons for phosphate solubilization, siderophore production operons,

36 and other PGP phenotypes. *Klebsiella pneumoniae* strains comprised 14 of the 22 nitrogen-

37 fixing isolates, and four others were members of closely related genera to *Klebsiella*. A

38 computational phenotyping approach was developed to rapidly screen for strains that have high

39 potential for nitrogen fixation and other PGP phenotypes while showing low risk for virulence

40 and antibiotic resistance. The majority of sugarcane isolates were below a genotypic and

41 phenotypic threshold, showing uniformly low predicted virulence and antibiotic resistance

42 compared to clinical isolates. Six prioritized strains were experimentally evaluated for PGP

43 phenotypes: nitrogen fixation, phosphate solubilization, and the production of siderophores,

44 gibberellic acid and indole acetic acid. Results from the biochemical assays were consistent

45 with the computational phenotype predictions for these isolates. Our results indicate that

46 computational phenotyping is a promising tool for the assessment of benefits and risks

47 associated with bacteria commonly detected in agricultural ecosystems.

48 **IMPORTANCE**   A genome-enabled approach was developed for the prioritization of

49   native bacterial isolates with the potential to serve as biofertilizers for sugarcane fields

50   in Colombia's Cauca Valley.  The approach is based on computational phenotyping,

51   which entails predictions related to traits of interest based on bioinformatic analysis of

52   whole genome sequences.  Bioinformatic predictions of the presence of plant growth

53   promoting traits were validated with experimental assays and more extensive genome

54   comparisons, thereby demonstrating the utility of computational phenotyping for

55   assessing the benefits and risks posed by bacterial isolates that can be used as

56   biofertilizers.  The quantitative approach to computational phenotyping developed here

57   for the discovery of biofertilizers has the potential for use with a broad range of

58   applications in environmental and industrial microbiology, food safety, water quality, and

59   antibiotic resistance studies.

## INTRODUCTION

60

61 The human population is expected to double in size within the next 50 years,

62 which will in turn lead to a massive increase in the global demand for food (1). Given

63 the scarcity of arable land worldwide, an increase in agricultural production of this

64 magnitude will require vast increases in cropping intensity and yield (2). It has been

65 estimated that as much as 90% of the increase in global crop production will need to

66 come from increased yield alone (3). At the same time, climate change and other

67 environmental challenges will necessitate the development of agricultural practices that

68 are more ecologically friendly and sustainable.

69 Chemical fertilizers that provide critical macronutrients to crops – such as

70 nitrogen (N), phosphorus (P), potassium (K), and sulfur (S) – are widely used to

71 maximize agricultural yield (4). The application of chemical fertilizers represents a

72 major cost for agricultural companies and also contributes to environmental damage, in

73 the form of eutrophication, hypoxia, harmful algal blooms, and air pollution through the

74 formation of microparticles (5). Biological fertilizers (biofertilizers) are comprised of

75 microbial inoculants that promote plant growth, thereby representing an alternative or

76 complementary approach for increasing crop yield, which is more sustainable and

77 environmentally friendly. Biofertilizers augment plant growth through nutrient

78 acquisition, hormone production, and by boosting immunity to pathogens (6).

79 Sugarcane is a tall, perennial grass cultivated in tropical and warm temperate

80 regions around the world, which is capable of producing high concentrations of sugar

81 (sucrose) and diverse byproducts (7). Sugarcane is consistently ranked as one of the

82 top ten planted crops in the world (8). Sugarcane agriculture plays a vital role in the

4

83  economy of Colombia by supporting the production of food products and biofuel

84  (ethanol).  The long-term goals of this work are to develop more effective and

85  sustainable sugarcane cropping practices in Colombia by simultaneously (i) increasing

86  crop yield, and (ii) decreasing the reliance on chemical fertilizers via the discovery,

87  characterization, and application of endemic (native) biofertilizers to Colombian

88  sugarcane fields.

89      Most sugarcane companies in Colombia currently use commercially available

90  biofertilizers, consisting primarily of nitrogen-fixing bacteria, which were discovered and

91  isolated from other countries (primarily Brazil), with limited success.  We hypothesized

92  that indigenous bacteria should be better adapted to the local environment and thereby

93  serve as more effective biofertilizers for Colombian sugarcane.  The use of indigenous

94  bacteria as biofertilizers should also mitigate potential threats to the environment posed

95  by non-native, and potentially invasive, species of bacteria.  Finally, indigenous bacteria

96  represent a renewable resource that agronomists can continually develop through

97  isolation and cultivation of local strains.

98      The advent of next-generation sequencing technologies has catalyzed the

99  development of genome-enabled approaches to harness plant microbiomes in

100  sustainable agriculture (9, 10).  The objective of this study was to use genome analysis

101  to predict the local bacterial isolates that have the greatest potential for plant growth

102  promotion while representing the lowest risk for virulence and antibiotic resistance.

103  Putative biofertilizer strains were isolated and cultivated from Colombian sugarcane

104  fields, and computational phenotyping was employed to predict their potential utility as

105  biofertilizers.  We then performed a laboratory evaluation of the predicted plant growth

106    promoting properties of the prioritized bacterial biofertilizer isolates, with the aim of

107    validating our computational phenotyping approach.

108

109    **RESULTS**

110        **Initial genome characterization of putative nitrogen-fixing bacteria.**  A

111    systematic cultivation approach, incorporating seven carbon substrates in nitrogen-free

112    media (Fig. S1), was employed to isolate putative nitrogen-fixing bacteria from four

113    different sugarcane plant compartments, and isolates were screened for nitrogen

114    fixation potential through PCR amplification of *nifH* genes.  This initial screening

115    procedure yielded several hundred clonal isolates of putative nitrogen-fixing bacteria,

116    and Ribosomal Intergenic Spacer Analysis (RISA) was subsequently used to identify the

117    (presumably) genetically unique strains from the larger set of clonal isolates.  A total of

118    22 potentially unique strains of putative nitrogen-fixing bacteria were isolated in this way

119    and selected for genome sequence analysis.

120        Genome sequencing and assembly summary statistics for the 22 isolates are

121    shown in Table 1.  Isolate genomes were sequenced to an average of 67x coverage

122    (range: 50x – 88x) and genome sizes range from 4.5Mb to 6.1Mb.  GC content varies

123    from 41.82% – 66.69%, with a distinct mode at ~57%.  The genome assemblies are

124    robust with a range of 24 – 294 contigs ≥500bp in length and averages of

125    N50=310,166bp and L50=8.4.  Genome sequence assemblies, along with their

126    functional annotations, can all be found using the NCBI BioProject PRJNA418312.

127  Individual BioSample, Genbank Accession, and Assembly Accession numbers for the

128  22 isolates are shown in Table S1.

129  **Comparative genomic analysis.**  Average nucleotide identity (ANI; Fig. 1) and

130  16S rRNA gene sequence analysis (Fig. S2) were employed in the taxonomic

131  assignment of nitrogen-fixing isolates and the results of both approaches were highly

132  concordant (Table 2), with ANI yielding superior resolution to 16S rRNA gene sequence

133  analysis.  A total of eight different species and seven different genera were identified

134  among the 22 isolates characterized.  Analysis of *nifH* gene sequences also gave

135  similar results; however, four of the isolates were not found to encode *nifH* genes,

136  despite their (apparent) ability to grow on nitrogen-free media and the positive *nifH* PCR

137  results.  This could be due to false-positives in the original PCR analysis for the

138  presence of *nifH* genes, or to changes in the composition of (possibly mixed) bacterial

139  cultures during subsequent growth steps after the initial isolation on nitrogen-free

140  media.

141  The majority of isolates, 14 of 22, were characterized as *Klebsiella pneumoniae*,

142  consistent with previous studies showing that *K. pneumoniae* strains are capable of

143  fixing nitrogen (11); in fact, the canonical *nif* operons were defined in the *K. pneumoniae*

144  type strain 342 genome sequence (12).  *K. pneumoniae* is also known to be an

145  opportunistic pathogen that can cause disease in immunocompromised human hosts

146  (13), which raises obvious safety concerns regarding its application to crops as part of a

147  biofertilizer inoculum.  We performed a comparative sequence analysis between the

148  endophytic nitrogen-fixing *K. pneumoniae* type strain 342, which is capable of infecting

149  the mouse urinary tract and lung (14), and five of the isolates identified as *K.*

7

150   *pneumoniae* here.  All genomes were shown to contain the *nif* cluster, which contains

151   five functionally related *nif* operons involved in nitrogen fixation (Fig. 2).  In contrast, the

152   four most critical pathogenicity islands implicated in the virulence of *K. pneumoniae* 342

153   were all missing in the environmental *K. pneumoniae* isolates characterized here (PAI

154   1-4 in Fig. 2A).  The absence of pathogenicity islands in the genome of the endophytic

155   nitrogen-fixer *K. michiganensis* Kd70 was associated with an inability to infect the

156   urinary tract in mice (15).  Our results indicate that nitrogen-fixing *K. pneumoniae*

157   environmental isolates from Colombian sugarcane fields do not pose a health risk

158   compared to clinical and environmental isolates that have previously been associated

159   with pathogenicity.  We explore this possibility in more detail in the following section on

160   computational phenotyping.

161         The *nifH* genes from the *Klebsiella* isolates characterized here form two distinct

162   phylogenetic clusters (Fig. 3).  This finding is consistent with previous results showing

163   multiple clades of *nifH* among *Klebsiella* genome sequences (16-18) and underscores

164   the potential functional diversity, with respect to nitrogen fixation, for the sugarcane

165   isolates.

166         **Computational phenotyping.**  Computational phenotyping, also referred to as

167   reverse genomics, was used to evaluate the potential of the bacterial isolates

168   characterized here to serve as biofertilizers for Colombian sugarcane fields.  For the

169   purpose of this study, computational phenotyping entails the prediction of specific

170   organismal phenotypes, or biochemical capacities, based on the analysis of functionally

171   annotated genome sequences (19).  The goal of the computational phenotyping

172   performed here was to identify isolates that show the highest predicted capacity for

173    plant growth promotion while presenting the lowest risk to human populations.

174    Accordingly, bacterial isolate genome sequences were screened for gene features that

175    correspond to the desirable (positive) characteristics of (i) nitrogen fixation and (ii) plant

176    growth promotion and the disadvantageous (negative) characteristics of (iii) virulence

177    and (iv) antimicrobial resistance.  Genome sequences were scored and ranked

178    according to the combined presence or absence of these four categories of gene

179    features as described in the Materials and Methods.  To compute genome scores, the

180    presence of nitrogenase and plant growth promoting genes contribute positive values,

181    whereas the presence of virulence factors and predicted antibiotic resistance yield

182    negative values.  Scores for each of the four specific phenotypic categories were

183    normalized and combined to yield a single composite score for each bacterial isolate

184    genome.  The highest scoring isolates are predicted as best candidates to be included

185    as part of a sugarcane biofertilizer inoculum (Figure 4; Table S2).  The predicted

186    biochemical capacities of the highest scoring isolates were subsequently experimentally

187    validated.

188        Isolates are ranked according to their composite genome scores, with a value of

189    10.87 observed as the highest potential for biofertilizer production (Figure 4).  Individual

190    gene and phenotype scores are color coded for each genome, and the four functional-

191    specific categories are shown separately.  The *nif* gene presence/absence profiles were

192    found to be highly similar for all but four of the bacterial isolates characterized here,

193    those which are not members of the *Klebsiella* genus, or closely related species, and do

194    not encode any *nif* genes.  The four non-nitrogen fixing isolates represent bacterial

195    species that are commonly found in soil (20-23), but they are not predicted to be viable

196    biofertilizers.  The *Kosakonia radicincitans* genome encodes the largest number of *nif*

197    genes (*n*=17) observed for any of the Colombian sugarcane isolates.  This is consistent

198    with previous studies showing that isolates of this species are capable of fixing nitrogen

199    (24).  The 14 characterized *K. pneumoniae* genomes all contain 16 out of 21 *nif* genes,

200    including the core *nifD* and *nifK* genes, which encode the heterotetramer core of the

201    nitrogenase enzyme, and the *nifH* gene, which encodes the dinitrogenase reductase

202    subunit (25).  These genomes also all encode the nitrogenase master regulators *nifA*

203    and *nifL*.  The missing *nif* genes for the *K. pneumoniae* isolates correspond to

204    accessory structural and regulatory proteins that are not critical for nitrogen fixation.

205    Accordingly, all of *K. pneumoniae* isolate genomes are predicted to encode the capacity

206    for nitrogen fixation, consistent with previous results (14, 26).  The single *Raoultella*

207    *ornithinolytica* isolate characterized here also contains the same 16 *nif* genes;

208    *Raoultella* species have previously been isolated from sugarcane (27) and have also

209    been demonstrated to fix nitrogen (28).

210        Initially, a total of 29 canonical bacterial plant growth promoting genes were

211    mined from the literature, 25 of which were found to be present in at least one of the

212    bacterial isolate genome sequences characterized here.  These 25 plant growth

213    promoting genes were organized into six distinct functional categories: phosphate

214    solubilization, indolic acetic acid (IAA) production, siderophore production, 1-

215    aminocyclopropane-1-carboxylate (ACC) deaminase, acetoin butanediol synthesis, and

216    peroxidases (Table S3).  For the purposes of visualization (Fig. 4), each functional

217    category is deemed to be present in an isolate genome sequence if all required genes

218    for that function can be found, but the weighted scoring for these categories is based on

219    individual gene counts as described in the Materials and Methods.  The *R.*

220    *ornithinolytica* isolate shows the highest predicted capacity for plant growth promotion,

221    with 5 of the 6 functional categories found to be fully present.  The majority of *K.*

222    *pneumoniae* isolates also show similar, but not identical, plant growth promoting gene

223    presence/absence profiles, with 3 or 4 functional categories present.  The capacity for

224    siderophore production is predicted to vary among *K. pneumoniae* isolates.  The *K.*

225    *radicincitans* genome also encodes 4 functional categories of plant growth promoting

226    genes, but differs from the *K. pneumoniae* isolates with respect to absence of

227    phosphate solubilization genes and the presence of acetoin butanediol synthesis genes.

228    Three of the four species found to lack *nif* genes also do not score present for any of the

229    plant growth promoting gene categories, further underscoring their predicted lack of

230    utility as biofertilizers.

231        Initially, a total of ~2,500 virulence factor genes were mined from the Virulence

232    Factor Database (VFDB) (29), 44 of which were found to be present in at least one of

233    the bacterial isolate genome sequences characterized here.  These 44 virulence factors

234    were organized into six distinct functional categories related to virulence and toxicity:

235    adherence, invasion, capsules, endotoxins, exotoxins, and siderophores.  The weighted

236    scores for these categories were computed based on individual gene presence/absence

237    patterns (Fig. 4).  In contrast to the *K. pneumoniae* clinical isolates which have

238    previously been characterized as opportunistic pathogens, the *K. pneumoniae*

239    environmental isolates showed uniformly low virulence scores.  The virulence factor

240    genes found among the *K. pneumoniae* environmental isolates correspond to

241    adherence proteins, capsules, and siderophores.  As shown in Fig. 2, genomes of

242   environmental isolates lack coding capacity for important invasion and toxin proteins,

243   including the Type IV secretion system, which are found in clinical *K. pneumoniae*

244   isolates.  The *R. ornithinolytica* and *K. radicincitans* isolates, both of which show high

245   scores for nitrogen fixation and plant growth promotion, gave higher virulence scores in

246   comparison to the environmental *K. pneumoniae* isolates.  Whereas *Bacillus pumilus*

247   had the lowest virulence score for any of the isolates, the remaining three non-nitrogen

248   fixing isolates had the highest virulence scores and were shown to encode well-known

249   virulence factors, such as Type IV, hemolysin, and fimbria secretion systems.

250        The predicted antibiotic resistance phenotypes for all characterized isolates were

251   fairly similar across the 20 classes of antimicrobial compounds for which predictions

252   were made.  The majority of the *K. pneumoniae* genomes, along with the relatively high

253   scoring *R. ornithinolytica* and *K. radicincitans* isolate genomes, indicated predicted

254   susceptibility to 10 of the 20 classes of antimicrobial compounds, intermediate

255   susceptibility for 2-4, and predicted resistance to 5-8.  The highest level of predicted

256   antibiotic resistance was seen for *Serratia marcescens*, with resistance predicted for 8

257   compounds and intermediate susceptibility predicted for 4.

258        Computational phenotyping scores for the four categories were normalized and

259   combined into a final score, with respect to their potential as biofertilizers (Fig. 4).  Most

260   of the top positions are occupied by *K. pneumoniae* isolates, with the exception of the

261   second-ranked *R. ornithinolytica* and the third-ranked *K. radicincitans*.  The results of a

262   similar analysis of four additional plant associated *Klebsiella* genomes are shown in Fig.

263   S3.

264      **Virulence comparison**.  The results described in the previous section indicate

265    that the majority of the *K. pneumoniae* strains isolated from Colombian sugarcane fields

266    have the highest overall potential as biofertilizers, including a low predicted potential for

267    virulence.  Nevertheless, the fact that strains of *K. pneumoniae* have previously been

268    characterized as opportunistic pathogens (30) raises concerns when considering the

269    use of *K. pneumoniae* as part of a bioinoculum that will be applied to sugarcane fields.

270    With this in mind, we performed a broader comparison of the predicted virulence profiles

271    for Colombian sugarcane isolates along with a collection of 28 clinical isolates of *K.*

272    *pneumoniae* and several other closely related species (See Table S5 for isolate

273    accession numbers).  For this comparison, the same virulence factor scoring scheme

274    described in the previous section was applied to all 50 genome sequences (Fig. 5).

275    Perhaps most importantly, a very clear distinction was observed in the virulence score

276    distribution, whereby all 28 clinical strains show a substantially higher predicted

277    virulence (from 4.45 to 2.11) in comparison to the environmental isolates (1.55 to 0.00).

278    Furthermore, the three environmental isolates that show the highest predicted virulence

279    correspond to species with low predicted capacity for both nitrogen fixation and plant

280    growth promotion; as such, these isolates would not be considered as potential

281    biofertilizers.  In particular, the *K. pneumoniae* environmental isolates showed uniformly

282    low predicted virulence compared to clinical isolates of the same species. Thus, the

283    results support, in principle, the use of the environmental *K. pneumoniae* isolates as

284    biofertilizers for Colombian sugarcane fields.

285      **Experimental validation of prioritized isolates**.  The top six scoring isolates

286    from the computational phenotyping were subjected to a series of cultivation-based

13

287    phenotypic assays in order to validate their predicted biochemical activities: (i)

288    acetylene reduction (a proxy for nitrogen fixation), (ii) phosphate solubilization, (iii)

289    siderophore production, (iv) gibberellic acid production, and (v) indole acetic acid

290    production.

291        Nitrogen fixation activity, as determined by acetylene reduction to ethylene, was

292    observed in all six isolates, three of which had higher levels in comparison to the

293    positive control (Fig. 6A).  All six of the isolates showed high levels of phosphate

294    solubilization (Fig. 6B & C) and siderophore production (Fig. 6D & E) compared to the

295    respective negative controls.  All six isolates showed the ability to produce gibberellic

296    acid (Fig. 6F), whereas none were able to produce indole acetic acid. The biochemical

297    assay results are consistent with the computational phenotype predictions for these

298    isolates.

299

## DISCUSSION

301        Members of the *Enterobacteriaceae* are often observed in cultivation-

302    independent studies of sugarcane and nitrogen-fixing *Enterobacteriaceae* are often

303    isolated from sugarcane plants worldwide (31-35).  The majority of isolates that were

304    obtained in this study from Colombian sugarcane belonged to the family

305    *Enterobacteriaceae*, with the *Klebsiella* as the most abundant genus along with *Serratia*,

306    *Kluyvera, Stenotrophomonas*, and *Bacillus*.  *Klebsiella* are Gram-negative, facultatively

307    anaerobic bacteria found in soils, plants, or water (36).  *Klebsiella* species have been

308    isolated from a large variety of crops worldwide, such as sugarcane, rice, wheat, and

309    maize (36-38).  *Klebsiella* species associated with plants have been shown to fix

14

310    nitrogen and express other plant growth promoting traits (37, 39).  Specifically,

311    *Klebsiella* species are abundant amongst the cultivable strains of *Enterobacteriaceae*

312    obtained from sugarcane (31).  For example, a survey of sugarcane in Guangxi, China

313    observed that *Klebsiella* was the most abundant plant-associated nitrogen-fixing

314    bacterial group (31), and among the strains isolated, *K. variicola* was shown to colonize

315    sugarcane and promote plant growth (37).  In addition, endophytic *Klebsiella* spp. have

316    been isolated from commercial sugarcane in Brazil, and their potential for plant growth

317    promotion was evaluated *in vitro* (40).  Finally in Pakistan, the phenotypic diversity of

318    plant growth promoting associated with sugarcane was determined, with *Klebsiella* also

319    appearing as one of the most abundant bacteria found (33). At the same time, Klebsiella

320    and other groups of *Enterobacteriaceae* commonly detected in agricultural systems are

321    abundant in the human microbiome and often contain closely related members that are

322    known opportunistic pathogens (41-44). The coexistence of microbial species that

323    contain plant beneficial traits with closely related strains that potentially cause human

324    diseases presents a challenge for the development of sustainable agriculture.  How can

325    we effectively perform a risk-benefit analysis of bacterial strains for potential use in the

326    agricultural biotechnology industry?  Thus, the overall goal of this study was to develop

327    high throughput methods for the isolation and screening of nitrogen-fixing bacteria for

328    their potential as biofertilizers.

329

330        **Computational phenotyping for the prioritization of potential biofertilizers.**

331    A computational phenotyping approach was developed for the screening of plant growth

332    promoting bacteria for their potential to serve as biofertilizers.  Computational

15

333    phenotyping entails the implementation of a variety of bioinformatic and statistical

334    methods to predict phenotypes of interest based on whole genome sequence analysis

335    (45, 46).  This approach has been used for a variety of applications in the biomedical

336    sciences: prediction of clinically relevant phenotypes, study of infectious diseases,

337    identification of opportunistic pathogenic bacteria in the human microbiome, and cancer

338    treatment decisions (47, 48).  To our knowledge, this study represents the first time

339    computational phenotyping has been used for agricultural applications.  To implement

340    computational phenotyping for the prioritization of potential biofertilizers, we developed

341    a scoring scheme based on the genome content of four functional gene categories of

342    interest: nitrogen-fixing genes, other plant growth promoting genes, virulence factor

343    genes, and antimicrobial resistance genes.

344         The results of the computational phenotyping predictions, confirmed by

345    laboratory experiments, support the potential use of selected bacterial strains isolated

346    from Colombian sugarcane fields as biofertilizers with minimum health risk to the human

347    population.  In particular, all isolates with higher scores (5.53 to 10.87, Fig. 4) in our

348    scheme were found to demonstrate the potential to fix nitrogen and to promote plant

349    growth in other ways, while lacking many of the important known virulence factors and

350    antibiotic resistance genes that can be found in clinical isolates of the same species.  In

351    general, isolates SCK7, SCK14, and SCK19 appeared to possess more potent plant

352    growth promoting properties compared to isolates SCK9, SCK16, and SCK21 (Fig. 4).

353    Our computational phenotyping scheme also has valuable negative predictive value.

354    Isolates that contained few or none of the beneficial traits that characterize biofertilizers,

355    *Bacillus pumilus* SCK3 and *Stenotrophomonas maltophilia* SCK1, had the lowest scores

356    (-10 and -11 respectively).  Finally, it is also worth reiterating that the computationally

357    predicted biochemical activities related to plant growth promotion were all validated by

358    experimental results (Fig. 6).

359            **Virulence profiling for the prioritization of potential biofertilizers.**

360    Opportunistic pathogens are microorganisms that usually do not cause disease in a

361    healthy host, but rather colonize and infect an immunocompromised host (49, 50).  For

362    example, *Klebsiella spp.* including *Klebsiella pneumoniae*, *Klebsiella oxytoca*, and

363    *Klebsiella granulomatis* were associated with nosocomial diseases (51) and other

364    hospital-acquired infections, primarily in immunocompromised persons (52).  The

365    potential for virulence, along with the presence of antimicrobial resistance genes, is an

366    obvious concern when proposing to use *Klebsiella* spp. as biofertilizers.  Importantly, we

367    found that the environmental *Klebsiella* isolates did not contain pathogenicity islands

368    associated with many virulence factor genes usually found in clinical isolates of

369    *Klebsiella* spp. (Fig. 2).  Our results are corroborated by a previous study of *Klebsiella*

370    *michiganensis* Kd70 isolated from the intestine of larvae of *Diatraea saccharalis*, for

371    which the genome was shown to contain multiple genes associated with plant growth

372    promotion and root colonization, but lacked pathogenicity islands in its genome (15).  In

373    order to shed further light on this problem, we extended our study of environmental

374    isolates from Colombian sugarcane to comparisons with genomes of *Klebsiella* clinical

375    isolates associated with opportunistic infections in humans along with a number other

376    environmental isolates with available genome sequences (Fig. 5).  The virulence factor

377    profiles for all of the environmental isolates were clearly distinct from the clinical strains,

17

378  which show uniformly higher virulence profile scores, underscoring the relative safety of

379  *Klebsiella* environmental isolates for use as biofertilizers.

380  **Potential for the use of computational phenotyping in other microbiology**

381  **applications.**  The results obtained from the computational phenotyping approach

382  developed in this study serve as a proof of principle in support of genomic guided

383  approaches to sustainable agriculture.  In particular, computational phenotyping can

384  serve to substantially narrow the search space for potential plant growth promoting

385  bacterial isolates, which can be further interrogated via experimental methods.

386  Computational phenotyping can be used to simultaneously identify beneficial properties

387  of plant associated bacterial isolates while avoiding potentially negative characteristics.

388  In principle, this approach can be applied to a broad range of potential plant growth

389  promoting isolates, or even assembled metagenomes, from managed agricultural

390  ecosystems.

391  We can also envision a number of other potential applications for computational

392  phenotyping of microbial genomes.  The computational phenotyping methodology

393  developed here has broad potential including diverse applications in agriculture, plant

394  and animal breeding, food safety, water quality microbiology along with other industrial

395  microbiology applications such as bioenergy, quality control/quality assurance, and

396  fermentation microbiology as well as human health applications such as pathogen

397  antibiotic resistance, virulence predictions, and microbiome characterization.  For

398  instance, computational phenotyping could be useful in food safety related to vegetable

399  crop production.  Vegetables harbor a diverse bacterial community dominated by the

400  family *Enterobacteriaceae,* Gram-negative bacteria that include a huge diversity of plant

18

401 growth promoting bacteria and enteric pathogens (53). Vegetables such as lettuce,

402 spinach, and carrots are usually consumed raw, which increases the concern of

403 bacterial infections or human disease outbreaks associated with consumption of

404 vegetables (49).

405      Increasing antibiotic resistance, generated by the abuse of antibiotics in

406 agriculture as well as medicine, is another major threat to human health (54), and the

407 food supply chain creates a direct connection between the environmental habitat of

408 bacteria and human consumers (55). Our computational phenotyping approach could

409 provide for an additional food safety solution, which could be used to prevent the spread

410 of antibiotic resistance pathogens genes present in the food chain.

411

**MATERIALS AND METHODS**

413 **Sampling and cultivation of putative nitrogen-fixing bacteria from**

414 **sugarcane.** INCAUCA is a Colombian sugarcane company located in the Cauca River

415 Valley in the southwest region of the country between the western and central Andes

416 mountain ranges (http://www.incauca.com/). Samples of leaves, rhizosphere soil, stem,

417 and roots were collected from the sugarcane fields 32T and 37T of the INCAUCA San

418 Fernando farm located in the Cauca Valley (3°16'30.0"N 76°21'00.0"W). A high-

419 throughput enrichment approach was developed to enable the cultivation of multiple

420 strains of putative nitrogen-fixing bacteria from sugarcane field samples; details of this

421 approach can be found in the Supplementary Material (Supplementary Methods and

422 Fig. S1).

19

423       A total of 22 distinct *nifH* PCR+ isolates that passed the initial cultivation and

424    screening steps were grown in LB medium (Difco) at 37°C for subsequent genomic DNA

425    extraction.  The E.Z.N.A. bacterial DNA kit (Omega Bio-Tek) was used for genomic

426    DNA extraction, and paired-end fragment libraries (~1,000bp) were constructed using

427    the Nextera XT DNA library preparation kit (Illumina).

428       **Genome sequencing, assembly, and annotation.**  Isolate genomic DNA

429    libraries were sequenced on the Illumina MiSeq platform using V3 chemistry, yielding

430    approximately 400,000 paired-end 300bp sequence reads per sample.  A list of all

431    genome sequence analysis programs that were used for this study is provided Table

432    S4.  Sequence read quality control and trimming were performed using the programs

433    FastQC version0.11.5 (56) and Trimmomatic (v.0.35) (57).  *De novo* sequence

434    assembly was performed using the program SPAdes (v.3.6) (58).  Assembled genome

435    sequences were annotated using the Rapid Annotations using Subsystems Technology

436    (RAST) Web server (59, 60) and NCBI Prokaryotic Genome Annotation Pipeline

437    (PGAP) (61).  The 15 *Klebsiella* isolates characterized in this way were briefly described

438    in a Genome Announcement (62), and the analysis here includes 7 additional non-

439    *Klebsiella* isolates.

440       **Comparative genomic analysis.**  Average Nucleotide Identity (ANI) was

441    employed to assign the taxonomy of the bacterial isolates characterized here (63, 64).

442    Taxonomic assignment was also conducted by targeting small subunit ribosomal RNA

443    (SSU rRNA) gene sequences.  Nitrogenase enzyme encoding *nifH* gene sequences

444    were extracted from isolate genome sequences, clustered, and taxonomically assigned

445    using the TaxaDiva (v.0.11.3) method developed by our group (12).  Whole genome

446    sequence comparisons between bacterial isolates characterized here and the *K.*

447    *pneumoniae* type strain 342 were performed using BLAST+ (v.2.2.28) (65) and

448    visualized with the program CGView (v.1.0) (66).  Details of the methods used

449    comparative genomic analysis can be found in the Supplementary Methods section.

450        **Computational phenotyping**.  Computational phenotyping was performed by

451    searching the bacterial isolate genome sequences characterized here for the

452    presence/absence of genes or features related to four functional classes of interest, with

453    respect to their potential as biofertilizers: (i) nitrogen fixation (NF), (ii) plant growth

454    promotion (PGP), (iii) virulence factors (iv), and (4) antimicrobial resistance (AMR).

455    Gene panels were manually curated by searching the literature (NCBI PubMed) for

456    genes implicated in nitrogen fixation and plant growth promotion.  The Virulence Factors

457    Database (VFDB) was used to curate the virulence factor gene panel (29).  AMR levels

458    were quantified using the PATRIC3/mic prediction tool (67).  A composite score was

459    developed to characterize each bacterial isolate genome sequence with respect to the

460    presence/absence of genes from the NF, PGP, and VF gene panels along with the

461    predicted AMR levels.  Details on the gene panels, AMR level, and the composite

462    scoring system can be found in the Supplementary Methods.

463        **Experimental validation.**  Predictions made by computational phenotyping were

464    validated using five distinct experimental assays: (1) Acetylene reduction assay for

465    nitrogen fixation activity, (2) Phosphate solubilization assay, (3) Siderophore production

466    assay, (4) Gibberellic acid production assay, and (5) Indole acetic acid production

467    assay.  Details of each experimental assay can be found in the Supplementary

468    Methods.

469    **Figure Legends**

470    FIG 1 **Phylogeny of the bacterial isolates characterized here (SCK numbers)**

471    **together with their most closely related bacterial type strains.**  The phylogeny was

472    reconstructed using pairwise average nucleotide identities between whole genome

473    sequence assemblies, converted to p-distances, with the neighbor-joining method.

474    Horizontal branch lengths are scaled according the p-distances as shown.

475

476    FIG 2 **Comparison of the *K. pneumoniae* type strain 342 to *K. pneumoniae***

477    **sugarcane isolates characterized here.**  (A) BLAST ring plot showing synteny and

478    sequence similarity between *K. pneumoniae* 342 and five *K. pneumoniae* sugarcane

479    isolates.  The *K. pneumoniae* 342 genome sequence is shown as the inner ring, and

480    syntenic regions of the five *K. pneumoniae* sugarcane isolates are shown as rings with

481    strain-specific color-coding according to the percent identity between regions of *K.*

482    *pneumoniae* 342 and the sugarcane isolates.   The genomic locations of *nif* operon

483    cluster along with four important pathogenicity islands (PAIs) are indicated.  PAI1 – type

484    IV secretion and aminoglycoside resistance, PAI2 hemolysin and fimbria secretion,

485    heme scavenging, PAI3 – radical S-adenosyl-L-methionine (SAM) and antibiotic

486    resistance pathways, PAI4 – fosfomycin resistance and hemolysin production.  (B) A

487    scheme of the *nif* operon cluster present in both *K. pneumoniae* 342 and the five *K.*

488    *pneumoniae* sugarcane isolates.

22

489 FIG 3 **Phylogeny of the *nifH* genes for the *Klebsiella* bacterial isolates**

490 **characterized here (SCK numbers).**  The phylogeny was reconstructed using pairwise

491 nucleotide p-distances between *nifH* genes recovered from the isolate genome

492 sequences using the neighbor-joining method.  Horizontal branch lengths are scaled

493 according the p-distances as shown.

494

495 FIG 4 **Computational phenotyping of the sugarcane bacterial isolates**

496 **characterized here.**  The presence (red) and absence (blue) profiles for nitrogen

497 fixation genes, plant growth promoting genes, and virulence factor genes are shown for

498 the 22 bacterial isolates.  Results are shown for all *n*=21 nitrogen-fixing genes.  Results

499 for plant growth promoting genes (*n*=25) and virulence factor genes (*n*=44) are merged

500 into six gene categories each.  Predicted antibiotic resistance profiles are shown for

501 *n*=20 antibiotic classes.  Detailed results for gene presence/absence and predicted

502 antibiotic resistance profiles are shown in Table S2.  The results for all four phenotypic

503 classes of interest were merged into a single priority score for each isolates (right side

504 of plot), as described in the Materials and Methods, and used to rank the isolates with

505 respect to their potential as biofertilizers.

506

507 FIG 5 **Comparison of predicted virulence profiles for clinical *K. pneumoniae***

508 **isolates compared to the environmental (sugarcane) bacterial isolates**

509 **characterized here**.  As in Fig. 4, predicted virulence profiles for six classes of

510 virulence factor genes are shown for each isolate.  Isolate-specific virulence factor

23

511    scores are shown for each isolate are based on the presence/absence profiles for the

512    *n*=44 virulence factor genes as described in the Materials and Methods.  The virulence

513    factor genes are used to rank the genomes from most (left) to least (right) virulent.

514    Clinical versus environmental samples are shown to the left and right, respectively, of

515    the red line, based on their virulence scores.

516

517    FIG 6 **Experimental validation of prioritized biofertilizer isolates**.  The

518    computationally predicted plant growth promoting phenotypes for the top six isolates

519    were experimentally validated.  All six strains were capable of acetylene reduction, i.e.

520    ethylene production (A), phosphate solubilization (B&C), siderophore production (D&E),

521    and gibberellic Acid production (F).

## ACKNOWLEDGMENTS

527

## REFERENCES

529    1.    Fess TL, Kotcon JB, Benedito VA. 2011. Crop breeding for low input agriculture:

530          A sustainable response to feed a growing world population. Sustainability

531          3:1742-1772.

532    2.    Bargaz A, Lyamlouli K, Chtouki M, Zeroual Y, Dhiba D. 2018. Soil microbial

533          resources for improving fertilizers efficiency in an integrated plant nutrient

534          management system. Front Microbiol 9:1606.

535    3.    Tilman D, Balzer C, Hill J, Befort BL. 2011. Global food demand and the

536          sustainable intensification of agriculture. Proc Natl Acad Sci U S A 108:20260-4.

537    4.    Stewart WM, Dibb DW, Johnston AE, Smyth TJ. 2005. The contribution of

538          commercial fertilizer nutrients to food production. Agronomy Journal 97:1-6.

539    5.    Savci S. 2012. Investigation of effect of chemical fertilizers on environment.

540          International Conference on Environmental Science and Development 1:287-

541          292.

542    6.    Bhardwaj D, Ansari MW, Sahoo RK, Tuteja N. 2014. Biofertilizers function as key

543          player in sustainable agriculture by improving soil fertility, plant tolerance and

544          crop productivity. Microbial Cell Factories 13.

545    7.    Cherubin MR, Karlen DL, Cerri CE, Franco AL, Tormena CA, Davies CA, Cerri

546          CC. 2016. Soil quality indexing strategies for evaluating sugarcane expansion in

547          Brazil. PLoS One 11:e0150860.

548    8.    Selman-Housein G, Lopez MA, Ramos O, Carmona ER, Arencibia AD,

549          Menendez E, Miranda F. 2000. Towards the improvement of sugarcane bagasse

550          as raw material for the production of paper pulp and animal feed. Plant Genetic

551          Engineering: Towards the Third Millennium 5:189-193.

552    9.    Dong M, Yang Z, Cheng G, Peng L, Xu Q, Xu J. 2018. Diversity of the bacterial

553          microbiome in the roots of four *Saccharum* species: *S. spontaneum, S.

554          robustum, S. barberi,* and *S. officinarum.* Front Microbiol 9:267.

555    10.   Li HB, Singh RK, Singh P, Song QQ, Xing YX, Yang LT, Li YR. 2017. Genetic

556          diversity of nitrogen-fixing and plant growth promoting *Pseudomonas* species

557          isolated from sugarcane rhizosphere. Front Microbiol 8:1268.

558    11.   Postgate JR. 1982. Biological nitrogen fixation: fundamentals. Philos Trans R

559          Soc Lond B Biol Sci 296:375-385.

560    12.   Gaby JC, Rishishwar L, Valderrama-Aguirre LC, Green SJ, Valderrama-Aguirre

561          A, Jordan IK, Kostka JE. 2018. Diazotroph community characterization via a

562          high-throughput *nifH* amplicon sequencing and analysis pipeline. Appl Environ

563          Microbiol 84.

564  13.  Li B, Zhao Y, Liu C, Chen Z, Zhou D. 2014. Molecular pathogenesis of *Klebsiella*

565      *pneumoniae*. Future Microbiol 9:1071-81.

566  14.  Fouts DE, Tyler HL, DeBoy RT, Daugherty S, Ren Q, Badger JH, Durkin AS,

567      Huot H, Shrivastava S, Kothari S, Dodson RJ, Mohamoud Y, Khouri H, Roesch

568      LF, Krogfelt KA, Struve C, Triplett EW, Methe BA. 2008. Complete genome

569      sequence of the N$_2$-fixing broad host range endophyte *Klebsiella pneumoniae*

570      342 and virulence predictions verified in mice. PLoS Genet 4:e1000141.

571  15.  Dantur KI, Chalfoun NR, Claps MP, Tortora ML, Silva C, Jure A, Porcel N,

572      Bianco MI, Vojnov A, Castagnaro AP, Welin B. 2018. The endophytic strain

573      *Klebsiella michiganensis* Kd70 lacks pathogenic island-like regions in its genome

574      and is incapable of infecting the urinary tract in mice. Front Microbiol 9:1548.

575  16.  Rosenblueth M, Martinez L, Silva J, Martinez-Romero E. 2004. *Klebsiella*

576      *variicola*, a novel species with clinical and plant-associated isolates. Systematic

577      and Applied Microbiology 27:27-35.

578  17.  Raymond J, Siefert JL, Staples CR, Blankenship RE. 2004. The natural history of

579      nitrogen fixation. Mol Biol Evol 21:541-54.

580  18.  Zehr JP, Jenkins BD, Short SM, Steward GF. 2003. Nitrogenase gene diversity

581      and microbial community structure: a cross-system comparison. Environ

582      Microbiol 5:539-54.

583  19.  Weimann A, Mooren K, Frank J, Pope PB, Bremges A, McHardy AC. 2016. From

584      genomes to phenotypes: Traitar, the microbial trait analyzer. mSystems

585      1:e00101-16.

586  20.  Deredjian A, Alliot N, Blanchard L, Brothier E, Anane M, Cambier P, Jolivet C,

587       Khelil MN, Nazaret S, Saby N, Thioulouse J, Favre-Bonte S. 2016. Occurrence of

588       *Stenotrophomonas maltophilia* in agricultural soils and antibiotic resistance

589       properties. Research in Microbiology 167:313-324.

590  21.  Caulier S, Gillis A, Colau G, Licciardi F, Liepin M, Desoignies N, Modrie P,

591       Legreve A, Mahillon J, Bragard C. 2018. Versatile antagonistic activities of soil-

592       borne *Bacillus* spp. and *Pseudomonas* spp. against *Phytophthora infestans* and

593       other potato pathogens. Front Microbiol 9:143.

594  22.  Badran S, Morales N, Schick P, Jacoby B, Villella W, Lorenz T. 2018. Complete

595       genome sequence of the *Bacillus pumilus* phage Leo2. Genome Announc 6.

596  23.  Pavan ME, Franco RJ, Rodriguez JM, Gadaleta P, Abbott SL, Janda JM,

597       Zorzopulos J. 2005. Phylogenetic relationships of the genus *Kluyvera*: transfer of

598       *Enterobacter intermedius* Izard et al. 1980 to the genus *Kluyvera* as *Kluyvera*

599       *intermedia* comb. nov. and reclassification of *Kluyvera cochleae* as a later

600       synonym of *K. intermedia*. Int J Syst Evol Microbiol 55:437-42.

601  24.  Berger B, Wiesner M, Brock AK, Schreiner M, Ruppel S. 2015. *K. radicincitans*, a

602       beneficial bacteria that promotes radish growth under field conditions. Agronomy

603       for Sustainable Development 35:1521-1528.

604  25.  Stacey G, Burris RH, Evans HJ. 1992. Biological Nitrogen Fixation. Chapman

605       and Hall, New York.

606  26.  Scott KF, Rolfe BG, Shine J. 1981. Biological nitrogen fixation: primary structure

607       of the *Klebsiella pneumoniae nifH* and *nifD* genes. J Mol Appl Genet 1:71-81.

608    27.    Luo T, Ou-Yang XQ, Yang LT, Li YR, Song XP, Zhang GM, Gao YJ, Duan WX,

609            An Q. 2016. *Raoultella* sp. strain L03 fixes $N_2$ in association with

610            micropropagated sugarcane plants. J Basic Microbiol 56:934-40.

611    28.    Schicklberger M, Shapiro N, Loque D, Woyke T, Chakraborty R. 2015. Draft

612            genome sequence of *Raoultella terrigena* R1Gly, a diazotrophic endophyte.

613            Genome Announc 3.

614    29.    Chen L, Zheng D, Liu B, Yang J, Jin Q. 2016. Hierarchical and refined dataset for

615            big data analysis--10 years on. Nucleic Acids Res 44:D694-697.

616    30.    Holt KE, Wertheim H, Zadoks RN, Baker S, Whitehouse CA, Dance D, Jenney A,

617            Connor TR, Hsu LY, Severin J, Brisse S, Cao HW, Wilksch J, Gorrie C, Schultz

618            MB, Edwards DJ, Nguyen KV, Nguyen TV, Dao TT, Mensinke M, Minh VL, Nhu

619            NTK, Schultsz C, Kuntaman K, Newton PN, Moore CE, Strugnell RA, Thomson

620            NR. 2015. Genomic analysis of diversity, population structure, virulence, and

621            antimicrobial resistance in *Klebsiella pneumoniae*, an urgent threat to public

622            health. Proceedings of the National Academy of Sciences of the United States of

623            America 112:E3574-E3581.

624    31.    Lin L, Li Z, Hu C, Zhang X, Chang S, Yang L, Li Y, An Q. 2012. Plant growth-

625            promoting nitrogen-fixing enterobacteria are in association with sugarcane plants

626            growing in Guangxi, China. Microbes Environ 27:391-8.

627    32.    Magnani GS, Didonet CM, Cruz LM, Picheth CF, Pedrosa FO, Souza EM. 2010.

628            Diversity of endophytic bacteria in Brazilian sugarcane. Genet Mol Res 9:250-8.

629    33.    Mehnaz S, Baig DN, Lazarovits G. 2010. Genetic and phenotypic diversity of

630           plant growth promoting rhizobacteria isolated from sugarcane plants growing in

631           pakistan. J Microbiol Biotechnol 20:1614-23.

632    34.    Ferrara FID, Oliveira ZM, Gonzales HHS, Floh EIS, Barbosa HR. 2012.

633           Endophytic and rhizospheric enterobacteria isolated from sugar cane have

634           different potentials for producing plant growth-promoting substances. Plant and

635           Soil 353:409-417.

636    35.    Taule C, Mareque C, Barlocco C, Hackembruch F, Reis VM, Sicardi M, Battistoni

637           F. 2012. The contribution of nitrogen fixation to sugarcane (Saccharum

638           officinarum L.), and the identification and characterization of part of the

639           associated diazotrophic bacterial community. Plant and Soil 356:35-49.

640    36.    Bagley ST. 1985. Habitat association of *Klebsiella* species. Infect Control 6:52-8.

641    37.    Wei CY, Lin L, Luo LJ, Xing YX, Hu CJ, Yang LT, Li YR, An QL. 2014.

642           Endophytic nitrogen-fixing *Klebsiella variicola* strain DX120E promotes

643           sugarcane growth. Biology and Fertility of Soils 50:657-666.

644    38.    Ji SH, Gururani MA, Chun SC. 2014. Isolation and characterization of plant

645           growth promoting endophytic diazotrophic bacteria from Korean rice cultivars.

646           Microbiol Res 169:83-98.

647    39.    Lin L, Wei C, Chen M, Wang H, Li Y, Li Y, Yang L, An Q. 2015. Complete

648           genome sequence of endophytic nitrogen-fixing *Klebsiella variicola* strain

649           DX120E. Stand Genomic Sci 10:22.

650    40.    Beneduzi A, Moreira F, Costa PB, Vargas LK, Lisboa BB, Favreto R, Baldani JI,

651           Passaglia LMP. 2013. Diversity and plant growth promoting evaluation abilities of

652          bacteria isolated from sugarcane cultivated in the South of Brazil. Applied Soil

653          Ecology 63:94-104.

654  41.    Denton M, Kerr KG. 1998. Microbiological and clinical aspects of infection

655          associated with Stenotrophomonas maltophilia. Clin Microbiol Rev 11:57-80.

656  42.    Downing KJ, Leslie G, Thomson JA. 2000. Biocontrol of the sugarcane borer

657          Eldana saccharina by expression of the Bacillus thuringiensis cry1Ac7 and

658          Serratia marcescens chiA genes in sugarcane-associated bacteria. Appl Environ

659          Microbiol 66:2804-10.

660  43.    Ribeiro VB, Zavascki AP, Rozales FP, Pagano M, Magagnin CM, Nodari CS, da

661          Silva RC, Dalarosa MG, Falci DR, Barth AL. 2014. Detection of bla(GES-5) in

662          carbapenem-resistant Kluyvera intermedia isolates recovered from the hospital

663          environment. Antimicrob Agents Chemother 58:622-3.

664  44.    Juhnke ME, des Jardin E. 1989. Selective medium for isolation of Xanthomonas

665          maltophilia from soil and rhizosphere environments. Appl Environ Microbiol

666          55:747-50.

667  45.    Richesson RL, Sun JM, Pathak J, Kho AN, Denny JC. 2016. Clinical phenotyping

668          in selected national networks: demonstrating the need for high-throughput,

669          portable, and computational methods. Artificial Intelligence in Medicine 71:57-61.

670  46.    Drouin A, Giguere S, Deraspe M, Marchand M, Tyers M, Loo VG, Bourgault AM,

671          Laviolette F, Corbeil J. 2016. Predictive computational phenotyping and

672          biomarker discovery using reference-free genome comparisons. Bmc Genomics

673          17.

674    47.    Berger AH, Brooks AN, Wu X, Shrestha Y, Chouinard C, Piccioni F, Bagul M,

675            Kamburov A, Innielinski M, Hogstrom L, Zhu C, Yang X, Pantel S, Sakai R,

676            Kaplan N, Root D, Narayan R, Natoli T, Lahr D, Tirosh I, Tamayo P, Getz G,

677            Wong B, Doench J, Subramanian A, Golub TR, Meyerson M, Boehm JS. 2016.

678            High-throughput phenotyping of lung cancer somatic mutations. Cancer

679            Research 76.

680    48.    Bone WP, Washington NL, Buske OJ, Adams DR, Davis J, Draper D, Flynn ED,

681            Girdea M, Godfrey R, Golas G, Groden C, Jacobsen J, Kohler S, Lee EMJ, Links

682            AE, Markello TC, Mungall CJ, Nehrebecky M, Robinson PN, Sincan M, Soldatos

683            AG, Tifft CJ, Toro C, Trang H, Valkanas E, Vasilevsky N, Wahl C, Wolfe LA,

684            Boerkoel CF, Brudno M, Haendel MA, Gahl WA, Smedley D. 2016.

685            Computational evaluation of exome sequence data using human and model

686            organism phenotypes improves diagnostic efficiency. Genetics in Medicine

687            18:608-617.

688    49.    Berg G, Erlacher A, Smalla K, Krause R. 2014. Vegetable microbiomes: is there

689            a connection among opportunistic infections, human health and our 'gut feeling'?

690            Microb Biotechnol 7:487-95.

691    50.    Fishman JA. 2013. Opportunistic infections--coming to the limits of

692            immunosuppression? Cold Spring Harb Perspect Med 3:a015669.

693    51.    Rosenblueth M, Martinez L, Silva J, Martinez-Romero E. 2004. Klebsiella

694            variicola, a novel species with clinical and plant-associated isolates. Syst Appl

695            Microbiol 27:27-35.

696  52.  Podschun R, Ullmann U. 1998. Klebsiella spp. as nosocomial pathogens:

697       epidemiology, taxonomy, typing methods, and pathogenicity factors. Clin

698       Microbiol Rev 11:589-603.

699  53.  Osterblad M, Pensala O, Peterzens M, Heleniusc H, Huovinen P. 1999.

700       Antimicrobial susceptibility of Enterobacteriaceae isolated from vegetables. J

701       Antimicrob Chemother 43:503-9.

702  54.  Canica M, Manageiro V, Abriouel H, Moran-Gilad J, Franz CMAP. 2019.

703       Antibiotic resistance in foodborne bacteria. Trends in Food Science &

704       Technology 84:41-44.

705  55.  Bengtsson-Palme J. 2017. Antibiotic resistance in the food supply chain: where

706       can sequencing and metagenomics aid risk assessment? Current Opinion in

707       Food Science 14:66-71.

708  56.  Andrews S.  FastQC a quality control tool for high throughput sequence data.

709       http://www.bioinformatics.babraham.ac.uk/projects/fastqc/. Accessed

710       07/31/2017.

711  57.  Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for

712       Illumina sequence data. Bioinformatics 30:2114-2120.

713  58.  Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,

714       Nikolenko SI, Pham S, Prjibelski AD. 2012. SPAdes: a new genome assembly

715       algorithm and its applications to single-cell sequencing. Journal of Computational

716       Biology 19:455-477.

717  59.  Aziz RK, Bartels D, Best AA, DeJongh M, Disz T, Edwards RA, Formsma K,

718       Gerdes S, Glass EM, Kubal M. 2008. The RAST Server: rapid annotations using

719       subsystems technology. BMC genomics 9:75.

720  60.  Wattam AR, Abraham D, Dalay O, Disz TL, Driscoll T, Gabbard JL, Gillespie JJ,

721       Gough R, Hix D, Kenyon R. 2013. PATRIC, the bacterial bioinformatics database

722       and analysis resource. Nucleic Acids Research 42:D581-D591.

723  61.  Tatusova T, DiCuccio M, Badretdin A, Chetvernin V, Nawrocki EP, Zaslavsky L,

724       Lomsadze A, Pruitt KD, Borodovsky M, Ostell J. 2016. NCBI prokaryotic genome

725       annotation pipeline. Nucleic Acids Res 44:6614-24.

726  62.  Medina-Cordoba LK, Chande AT, Rishishwar L, Mayer LW, Marino-Ramirez L,

727       Valderrama-Aguirre LC, Valderrama-Aguirre A, Kostka JE, Jordan IK. 2018.

728       Genome sequences of 15 *Klebsiella* sp. isolates from sugarcane fields in

729       Colombia's Cauca Valley. Genome Announc 6.

730  63.  Konstantinidis KT, Tiedje JM. 2005. Genomic insights that advance the species

731       definition for prokaryotes. Proceedings of the National Academy of Sciences of

732       the United States of America 102:2567-2572.

733  64.  Goris J, Konstantinidis KT, Klappenbach JA, Coenye T, Vandamme P, Tiedje

734       JM. 2007. DNA–DNA hybridization values and their relationship to whole-

735       genome sequence similarities. International Journal of Systematic and

736       Evolutionary Microbiology 57:81-91.

737  65.  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden

738       TL. 2009. BLAST+: architecture and applications. BMC Bioinformatics 10:421.

739    66.    Grant JR, Arantes AS, Stothard P. 2012. Comparing thousands of circular

740           genomes using the CGView comparison tool. BMC Genomics 13:202.

741    67.    Nguyen M, Brettin T, Long SW, Musser JM, Olsen RJ, Olson R, Shukla M,

742           Stevens RL, Xia F, Yoo H, Davis JJ. 2018. Developing an in silico minimum

743           inhibitory concentration panel test for *Klebsiella pneumoniae*. Sci Rep 8:421.

744 **Tables**

745 Table 1. **Genome assembly statistics for the isolates characterized here.**

746

| Sample ID | Genome Length (bp) | N50[a] | L50[b] | GC(%) | # of Contigs[c] |
|---|---|---|---|---|---|
| SCK1 | 4,522,541 | 402,304 | 4 | 66.79 | 24 |
| SCK2 | 5,231,439 | 417,927 | 5 | 59.33 | 53 |
| SCK3 | 3,824,428 | 670,745 | 3 | 41.82 | 150 |
| SCK4 | 4,511,030 | 223,239 | 8 | 66.79 | 55 |
| SCK5 | 5,774,634 | 162,673 | 13 | 53.1 | 98 |
| SCK6 | 6,094,823 | 117,689 | 15 | 56.73 | 294 |
| SCK7 | 5,693,007 | 282,996 | 7 | 57.03 | 50 |
| SCK8 | 5,695,902 | 281,292 | 9 | 57.03 | 50 |
| SCK9 | 5,579,618 | 311,650 | 6 | 57.03 | 42 |
| SCK10 | 5,591,472 | 614,324 | 3 | 57.03 | 34 |
| SCK11 | 5,696,136 | 382,597 | 5 | 57.15 | 268 |
| SCK12 | 5,817,089 | 176,655 | 10 | 57.02 | 79 |
| SCK13 | 5,476,221 | 358,490 | 5 | 57.34 | 33 |
| SCK14 | 5,465,811 | 300,899 | 5 | 57.34 | 41 |
| SCK15 | 5,564,330 | 330,579 | 5 | 57.15 | 43 |
| SCK16 | 5,795,921 | 478,592 | 3 | 54.06 | 84 |
| SCK17 | 5,475,984 | 358,490 | 4 | 57.34 | 35 |
| SCK18 | 5,476,135 | 422,400 | 3 | 57.34 | 32 |
| SCK19 | 5,688,396 | 270,585 | 7 | 57.09 | 56 |
| SCK20 | 5,500,801 | 82,111 | 20 | 57.45 | 165 |
| SCK21 | 5,324,920 | 112,078 | 15 | 55.26 | 100 |
| SCK22 | 5,847,607 | 65,329 | 29 | 57.02 | 181 |

747

748 [a] When the contigs of an assembly are arranged from largest to smallest, N50 is the

749 length of the contig that makes up at least 50% of the genome

750 [b] L50 is the number of contigs equal to or longer than N50 In other words, L50, for

751 example, is the minimal number of contigs that cover half the assembly

752 [c] Number of contigs ≥500bp in length

753 Table 2. **Identity of the most closely related species (genus) for the isolates**

754 **characterized here.** Species (genus) identification was performed using average

755 nucleotide identity (ANI), 16S rRNA and *nifH* sequence comparisons.

756

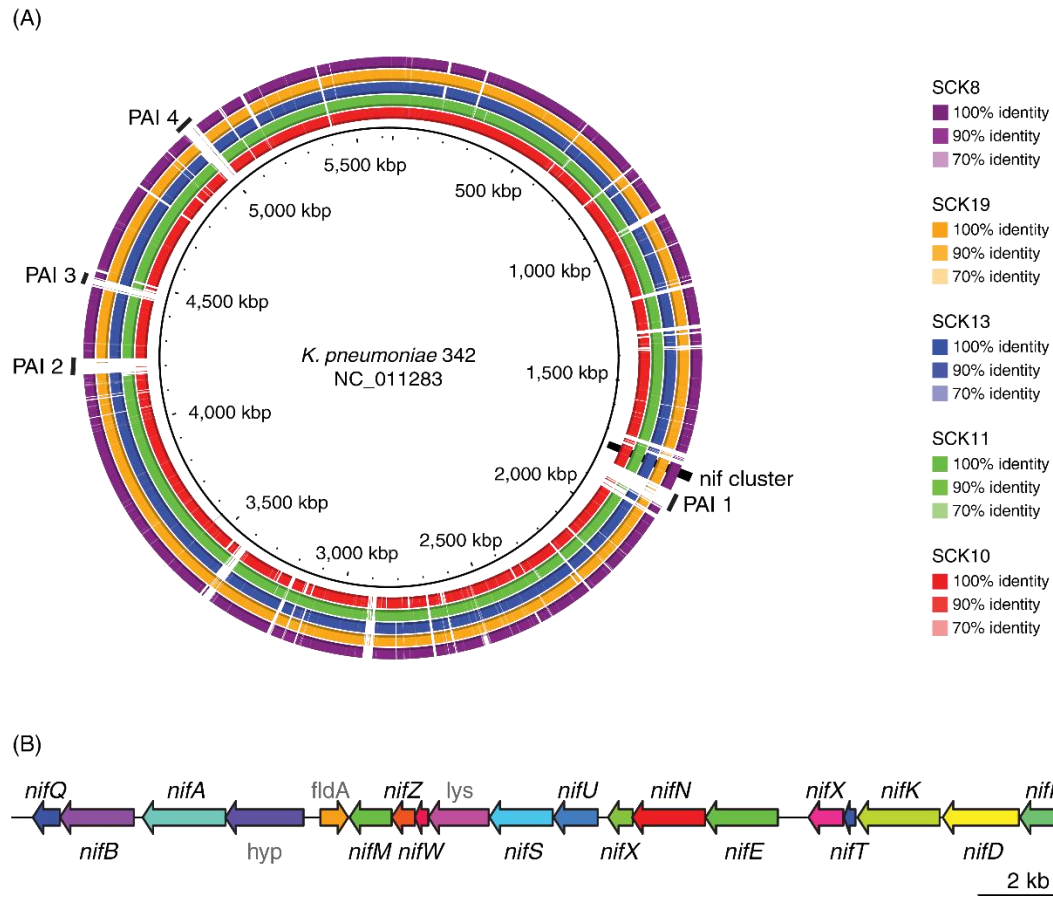| Strain | ANI | 16S | *nifH* |
|---|---|---|---|
| SCK1 | *Stenotrophomonas maltophilia* | *Stenotrophomonas* | NA |
| SCK2 | *Serratia marcescens* | *Serratia* | NA |
| SCK3 | *Bacillus pumilus* | *Bacillus* | NA |
| SCK4 | *Stenotrophomonas maltophilia* | *Stenotrophomonas* | NA |
| SCK5 | *Kluyvera intermedia* | *Kluyvera* | *Kluyvera* |
| SCK6 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK7 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK8 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK9 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK10 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK11 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK12 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK13 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK14 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK15 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK16 | *Kosakonia radicincitans* | *Kosakonia* | *Kosakonia* |
| SCK17 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK18 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK19 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK20 | *Klebsiella pneumoniae* | *Klebsiella* | *Klebsiella* |
| SCK21 | *Raoultella ornithinolytica* | *Raoultella* | *Raoultella* |
| SCK22 | *Klebsiella variicola* | *Klebsiella* | *Klebsiella* |

757
758
759
760
761
762
763
764
765
766
767
768

769

FIG 1 **Phylogeny of the bacterial isolates characterized here (SCK numbers)**

**together with their most closely related bacterial type strains**. The phylogeny was

reconstructed using pairwise average nucleotide identities between whole genome

sequence assemblies, converted to p-distances, with the neighbor-joining method.

Horizontal branch lengths are scaled according the p-distances as shown.

775

FIG 2 **Comparison of the *K. pneumoniae* type strain 342 to *K. pneumoniae* sugarcane isolates characterized here**. (A) BLAST ring plot showing synteny and sequence similarity between *K. pneumoniae* 342 and five *K. pneumoniae* sugarcane isolates. The *K. pneumoniae* 342 genome sequence is shown as the inner ring, and syntenic regions of the five *K. pneumoniae* sugarcane isolates are shown as rings with strain-specific color-coding according to the percent identity between regions of *K. pneumoniae* 342 and the sugarcane isolates. The genomic locations of nif operon cluster along with four important pathogenicity islands (PAIs) are indicated. PAI1 – type IV secretion and aminoglycoside resistance, PAI2 hemolysin and fimbria secretion, heme scavenging, PAI3 – radical S-adenosyl-L-methionine (SAM) and antibiotic resistance pathways, PAI4 – fosfomycin resistance and hemolysin production. (B) A scheme of the nif operon cluster present in both *K. pneumoniae* 342 and the five *K. pneumoniae* sugarcane isolates.

39
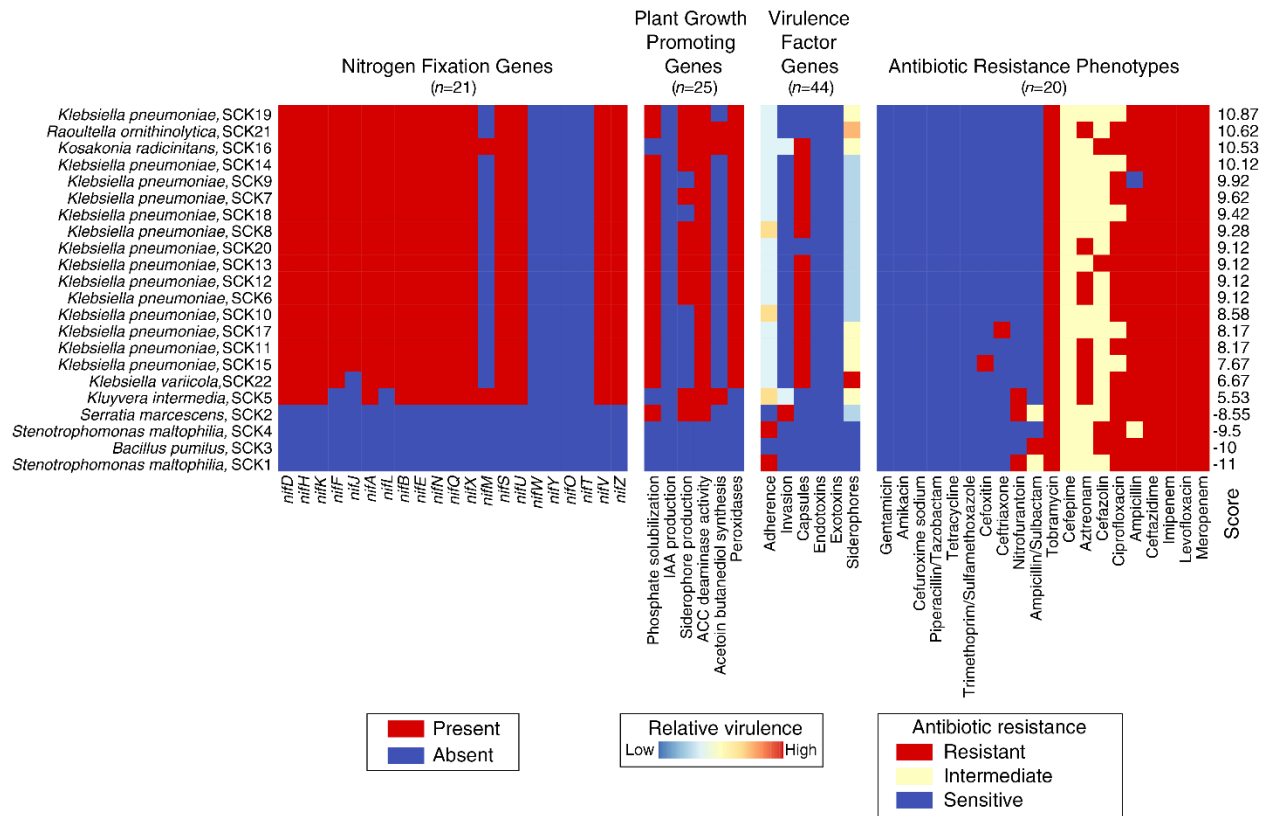
FIG 3 **Phylogeny of the *nifH* genes for the *Klebsiella* bacterial isolates**

**characterized here (SCK numbers).** The phylogeny was reconstructed using pairwise

nucleotide p-distances between *nifH* genes recovered from the isolate genome

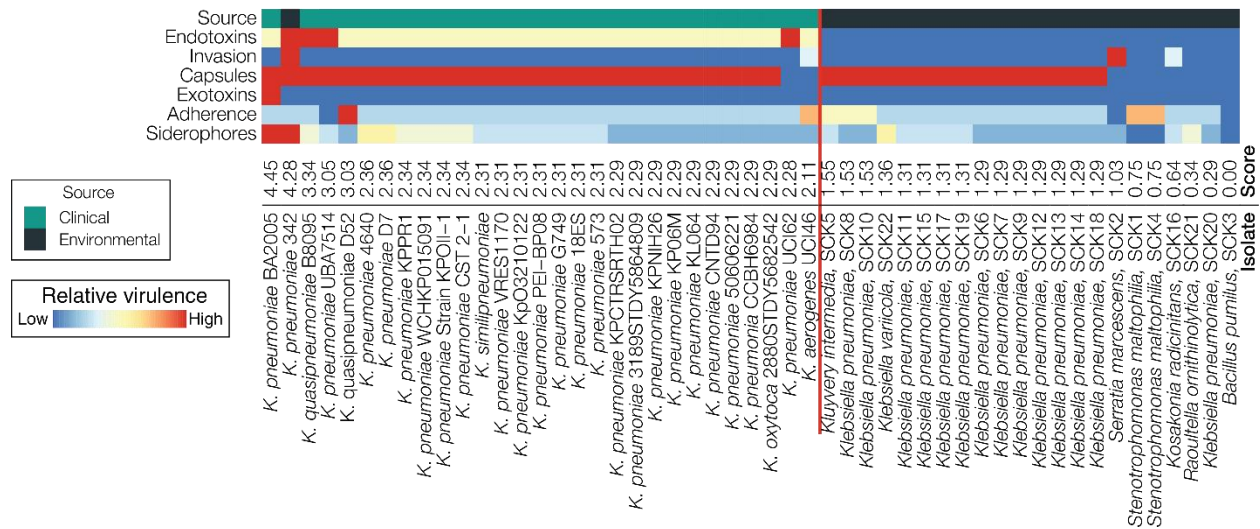sequences using the neighbor-joining method. Horizontal branch lengths are scaled

according the p-distances as shown.

794

FIG 4 **Computational phenotyping of the sugarcane bacterial isolates**

**characterized here.** The presence (red) and absence (blue) profiles for nitrogen

fixation genes, plant growth promoting genes, and virulence factor genes are shown for

the 22 bacterial isolates. Results are shown for all $n$=21 nitrogen-fixing genes. Results

for plant growth promoting genes ($n$=25) and virulence factor genes ($n$=44) are merged

into six gene categories each. Predicted antibiotic resistance profiles are shown for

$n$=20 antibiotic classes. Detailed results for gene presence/absence and predicted

antibiotic resistance profiles are shown in Table S2. The results for all four phenotypic

classes of interest were merged into a single priority score for each isolates (right side

of plot), as described in the Materials and Methods, and used to rank the isolates with

respect to their potential as biofertilizers.

806

FIG 5 **Comparison of predicted virulence profiles for clinical *K. pneumoniae***

**isolates compared to the environmental (sugarcane) bacterial isolates**

**characterized here**. As in Fig. 4, predicted virulence profiles for six classes of

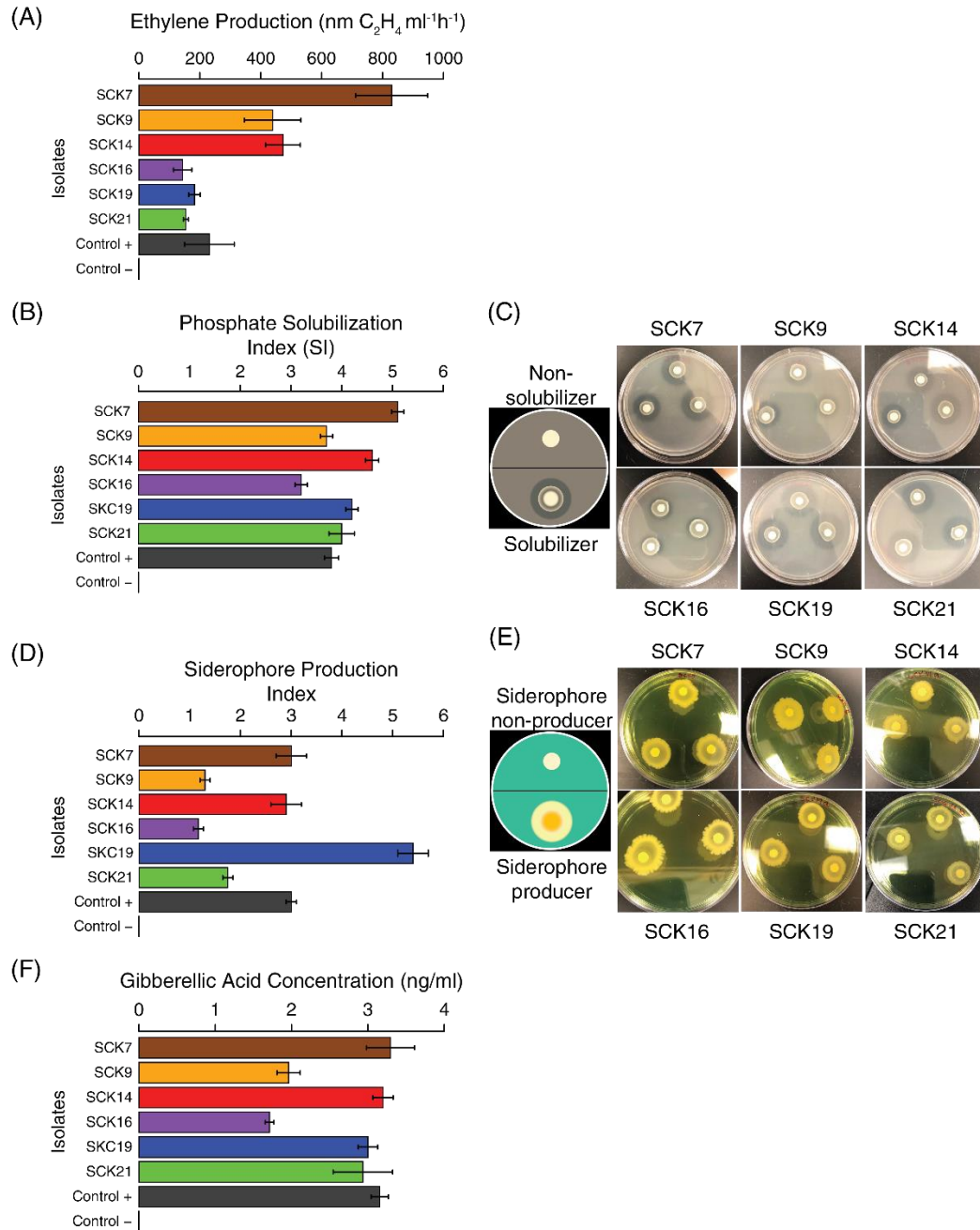virulence factor genes are shown for each isolate. Isolate-specific virulence factor

scores are shown for each isolate are based on the presence/absence profiles for the

$n$=44 virulence factor genes as described in the Materials and Methods. The virulence

factor genes are used to rank the genomes from most (left) to least (right) virulent.

Clinical versus environmental samples are shown to the left and right, respectively, of

the red line, based on their virulence scores.

FIG 6 **Experimental validation of prioritized biofertilizer isolates**. The computationally predicted plant growth promoting phenotypes for the top six isolates were experimentally validated. All six strains were capable of acetylene reduction, i.e. ethylene production (A), phosphate solubilization (B&C), siderophore production (D&E), and gibberellic acid production (F).

43