1    **Title:**Genomic characterization of a diazotrophic microbiota associated with maize aerial root

2    mucilage

3

4    **Author names:**

5    Shawn M. Higdon[1], Tania Pozzo[1], Nguyet Kong[2], Bihua Huang[2,3], Mai Lee Yang[2], Richard

6    Jeannotte[2], C. Titus Brown[2], Alan B. Bennett[1*], Bart C. Weimer[2,3*]

7

8    **Affiliation:**

9    [1]Department of Plant Sciences, University of California, Davis, California 95616

10   [2]Department of Population Health and Reproduction, University of California, Davis, California

11   95616

12   [3]100K Pathogen Genome Project, University of California, Davis, California 95616

13   *To whom correspondence should be sent.

14

15   Corresponding authors:

16   abbennett@ucdavis.edu

17   bcweimer@ucdavis.edu

18

19   Keywords: Biological Nitrogen Fixation, Diazotrophic Bacteria, Mucilage Polysaccharide,

20   Nitrogen Fixation Gene; *Zea mays*

## Abstract

A geographically isolated maize landrace cultivated on nitrogen-depleted fields without synthetic

fertilizer in the Sierra Mixe region of Oaxaca, Mexico utilizes nitrogen derived from the

atmosphere and develops an extensive network of mucilage-secreting aerial roots that harbors a

diazotrophic microbiota. Targeting these diazotrophs, we selected nearly 600 microbes from a

collection isolated from these plants and confirmed their ability to incorporate heavy nitrogen

($^{15}N_2$) metabolites *in vitro*. Sequencing their genomes and conducting comparative bioinformatic

analyses showed that these genomes had substantial phylogenetic diversity. We examined each

diazotroph genome for the presence of *nif* genes essential to nitrogen fixation (*nif*HDKENB) and

carbohydrate utilization genes relevant to the mucilage polysaccharide digestion. These analyses

identified diazotrophs that possessed canonical *nif* gene operons, as well as many other operon

configurations with concomitant fixation and release of >700 different $^{15}N$ labeled metabolites.

We further demonstrated that many diazotrophs possessed alternative *nif* gene operons and

confirmed their genomic potential to derive chemical energy from mucilage polysaccharide to

fuel nitrogen fixation. These results confirm that some diazotrophic bacteria associated with

Sierra Mixe maize were capable of incorporating atmospheric nitrogen into their small molecule

extracellular metabolites through multiple *nif* gene configurations while others were able to fix

nitrogen without the canonical (*nif*HDKENB) genes.

2

40 **Data Summary**

41 Genetic resources, including biological materials and nucleic acid sequences, were accessed

42 under an Access and Benefit Sharing (ABS) Agreement between the Sierra Mixe community and

43 the Mars Corporation, and with authorization from the Mexican government. An internationally

44 recognized certificate of compliance has been issued by the Mexican government under the

45 Nagoya Protocol for such activities (ABSCH-IRCC-MX-207343-3). Any party seeking access to

46 the nucleic acid sequences underlying the analysis reported here is subject to the full terms and

47 obligations of the ABS agreement and the authorization from the government of Mexico.

48 Individuals wishing to access nucleic acid sequence data for scientific research activities should

49 contact Mars Incorporated Chief Science Officer at CSO@effem.com.

50

51

## Introduction

53 Nitrogen is an essential macroelement for plant productivity that is often limiting to plant growth

54 when the natural abundance of its bio-available forms is depleted in the environment. Exogenous

55 nitrogen is currently provided for maize cultivation either through synthetic Haber-Bosch

56 fertilizer produced at high environmental and economic cost (1), or from crop rotation with

57 legumes that replenish field nitrogen levels by symbiotic association with diazotrophs, bacteria

58 capable of biological nitrogen fixation (BNF) (2, 3). Because maize is a crop of immense

59 agricultural importance, the establishment of conventional varieties capable of meeting their

60 nitrogen demands through mutualistic associations with free-living diazotrophic bacteria would

61 be of significant value to the goal of achieving global food security through sustainable

62 intensification without relying on fertilization (4). One strategy for the discovery of useful maize

63 diazotrophic plant-microbe associations involves exploring the microbiome of cultivated maize

64 landraces near the center of the maize origin of domestication (5).

65 A recent report demonstrated that an indigenous landrace of maize found in Totontepec

66 Villa de Morelos in the Sierra Mixe region of Mexico acquires 28-82% of its nitrogen from the

67 air and exhibits an extensive system of aerial roots with heavy secretion of a mucilage composed

68 of unique complex polysaccharides (6). Analysis of a public, low coverage shotgun metagenome

69 sequences from the roots, stems, and aerial root mucilage revealed the aerial root mucilage

70 microbiota to be enriched in taxa with many known species that are diazotrophic (6). In addition,

71 the mucilage was the only plant tissue type to be enriched for homologs of the canonical nitrogen

72 fixation genes (*nif*HDKENB), as previously proposed by Dos Santos et al., to be essential for a

73 bacterium to be diazotrophic (6, 7). The demonstration that the Sierra Mixe mucilage harbors a

74 diazotrophic microbial community, that it exhibits reduced taxonomic complexity, and the

4

75    absence of soil from aerial root mucilage suggests that it could be a useful model system for

76    elucidating associative mechanisms between free-living bacteria and cereal crops with mucilage-

77    secreting aerial roots, such as maize.

78         Following investigations reported by Van Deynze et al. (6), we hypothesized that free-

79    living diazotrophs from the aerial root mucilage microbiota utilize mucilage derived

80    carbohydrates as an energy source for BNF. To address this, we cultured many bacteria by

81    targeting diazotrophic bacteria specifically associated with Sierra Mixe maize. Subsequenctly,

82    we characterized 588 microbial diazotrophic isolates to verify fixation and other traits using

83    whole genome sequencing (WGS). Measuring the ability to incorporate heavy dinitrogen gas

84    ($^{15}N_2$) into secreted metabolites with tandem mass spectrometry confirmed that the isolates were

85    diazotrophic and produced a variety of compounds containing the label. Subsequent WGS

86    analysis using comparative genomics with each diazotrophic isolate genome included assessing

87    differences in nucleotide composition, assigning taxonomic classifications, and estimating

88    percent recovery from the mucilage microbiome. To elucidate the genomic determinants for

89    BNF by mucilage-derived diazotrophs, we examined their genomes for the presence of features

90    related to mucilage polysaccharide utilization, the canonical *nif* genes based on the Dos Santos

91    model (7) with the *Klebsiella pneumoniae* NIF regulon as the model framework, and known

92    alternative *nif* genes. Our results indicate that the mucilage microbial isolates contained the

93    capacity to utilize the mucilage complex polysaccharide and, surprisingly, that many of the

94    diazotrophic isolates did not possess recognizable homology for known *nif* genes – yet were

95    diazotrophic. These findings suggest the presence of novel mechanisms of nitrogen fixation by

96    many phylogenomic groups of bacteria, several of which were not previously associated with

97    this trait.

## Methods

98

### Bacterial isolation

99

100    Roots, stems and mucilage (200–500 μL) collected from different fields of the Sierra Mixe

101    region in Mexico were spread on 1.5% BHI (BD, catalogue number 211059; Franklin Lakes, NJ,

102    USA) or modified nitrogen-free M9 agar (BD) with and without a 1% (w/v) D-arabinose,

103    galactose or xylose at pH 5, 5.8 or 7. Plant tissues were blended in 1×PBS prior to culturing on

104    medium and the blender decontaminated with 10% bleach followed by 70% ethanol between

105    samples (v/v). Cultures were incubated at 25°C or 37°C, aerobically and anaerobically, for up to

106    4 weeks. Once colonies appeared, they were sub-cultured on the same medium to ensure purity.

107    Each organism was grown in BHI broth at the respective condition and resuspended in 5% non-

108    fat dry milk and glycerol and stored cryogenically for further use.

### Biological Nitrogen Fixation Assay

109

110    To assay for Microbial $^{15}N_2$ assimilation, isolates were first grown twice overnight in the

111    respective growth condition prior to collection and washed twice with 0.9% (w/v) saline solution

112    before re-suspension in Fahraeus medium containing 1% D-glucose at pH 5.8 to determine the

113    nitrogen fixation capacity. Prior to the fixation assay, dissolved oxygen was removed from the

114    medium by sparging with argon gas for 1.5 hours while stirring and a vacuum pump was used to

115    remove any oxygen in the headspace. Each isolate ($OD_{600} = 2$; 2 mL) was added to an airtight 4

116    mL glass vial. Addition of the heavy atom was achieved by removing 20 mL of headspace gas

117    and replacing it with 5 mL of either $^{15}N_2$ or $^{14}N_2$ nitrogen gas directly into the culture. The

118    cultures were incubated at 37 °C anaerobically for 6 - 48 hours, depending on the growth rate

6

119    and collected at the beginning of stationary phase for each culture. All experiments were done in

120    triplicate.

**Microbial metabolite extraction and quantitation**

122    Subsequent to growth the metabolites were extracted from cell pellets as described by Villas-

123    Bôas (8). Bacterial cultures were transferred to 2 mL tubes and centrifuged at 14,800 rpm for 10

124    min at –9 °C. After collection of the cell pellet 500 μL of cold methanol (-20°C) was added

125    before lysing the cells with bead beating (9, 10). After adding 0.4 g of 0.1 mm glass beads cells

126    were lysed by two cycles of bead beating with 30 s per cycle, 1 min rest on ice between each

127    cycle [9,10]. The lysed samples were centrifuged at 14,800 rpm for 10 min at –9 °C after which 50

128    ml of each supernatant was transferred to LC vials for metabolite analysis. Samples were stored

129    in -80 °C until analysis using LC/TOF-MS. In order to confirm the enrichment by $^{15}$N, a subset

130    of residual pellets (50 mg of dried pellets), after metabolite extraction, were submitted to the UC

131    Davis Stable Isotope Facility for Isotope Ratio Mass Spectrometry (IRMS) analysis ($^{15}$N/$^{14}$N

132    ratio). $^{15}$N-labeled metabolite analysis was performed using LC-TOF G6230A (Agilent

133    Technologies) instrument equipped with 1290 Infinity HPLC system. Chromatographic

134    separation was performed on a Zorbax Eclipse XDB-C18 (2.1×15 mm, 1.8 μm) with a flow of

135    500 μL·min$^{-1}$ and the following elution gradient: 0 min, 10 % B; 2.5 min, 80 % B; 4.0 min, 100

136    % B; 4.5 min, 100 % B; 5.0 min, 10 % B; 6.0 min, 10 %. Solvent A was water and solvent B was

137    acetonitrile, both containing 0.1 % formic acid with a column temperature of 40 °C and an

138    injection volume of 1-5 μL. This HPLC system was connected to an Agilent 6230 time-of-flight

139    analyzer with an Agilent Jet Stream electrospray (ESI) interface operating in positive ion mode

140    under the following conditions: capillary 3500 V, nebulizer 35 psi g, drying gas 8 L·min$^{-1}$, gas

141    temperature 350 °C, skimmer voltage 80 V, fragmentor voltage 135 V, octapole RF 750 V. The

7

142  mass axis was calibrated using the mixture provided by the manufacturer in the m/z 50–1700

143  range. Acquisition rate was set to 1 spectrum per second (13,593 transients/spectrum). A

144  reference solution provided continuous calibration using the following reference masses:

145  121.0509 and 922.0098 m/z. Accurate mass spectra from 70 to 1700 m/z were recorded and

146  processed with MassHunter Workstation software (B.04.00). Statistical analysis was performed

147  using GeneSpring-MassProfiler Pro (version 12.1) software from Agilent Technologies, and

148  MetaboAnalyst (http://www.metaboanalyst.ca/) (11).

149  **Biomarkers of nitrogen-fixation**

150  The basis of this approach is that as a microbe incorporates $^{15}N$ by fixation, $^{15}N$ will be used in

151  the biosynthesis of small molecules and macromolecules, such as nucleic acids and proteins,

152  shifting their masses of 1 unit per atom of nitrogen replaced. A given bacteria fixing nitrogen and

153  exposed to $^{15}N_2$ gas will have a very different spectrum compared to the same bacteria exposed

154  to $^{14}N_2$ only.

155  The mass spectrometry analysis of each extract generated an average spectrum per

156  sample that contains thousands of masses. All the spectra were aligned and assembled in one

157  data matrix using SpecAlign software. Using the data from all the isolates, we performed a

158  statistical analysis (t-test, in MetaboAnalyst) (11) to determine the features (masses) that were

159  significantly changing across isolates when controls and treated samples were compared. This

160  approach allows us to identify biomarkers of nitrogen fixation that could be common to all the

161  isolates, totally or partially (some isolates could have all the biomarkers identified, some others

162  only a subset). More than 700 masses were significantly different using a q value (a p-value

163  adjusted by False Discovery Rate (FDR); this statistical approach allows to correct for possible

164  false positives) of 0.05 as threshold (q value $\leq$ 0.05 was determined to be significant). Masses

8

165     with q≤0.05 and fold-change (intensity of given mass in 15N-treated samples vs intensity of the

166     same mass in $^{14}$N-treated samples) of >1 were considered in the following calculations. Then for

167     each isolate, the relative intensities (percentage of each peak raw intensity over total raw signal)

168     for all the biomarkers were summed. Sums of the relative intensities for the biomarkers in

169     control and treated samples, for a given isolate, were computed and ratio $^{15}$N/$^{14}$N was calculated.

170     Isolates with BNF ratios greater than or equal to 1 were considered as sufficient $N_2$-fixers, where

171     the sum of peak intensities under $^{15}N_2$-enriched atmosphere was found to be equal to that of the

172     unenriched control. Following this logic, isolates with BNF ratios greater than 1 were considered

173     to be more efficient $N_2$-fixers (i.e. higher $^{15}$N ratios indicated a higher detected abundance of $^{15}$N

174     atom incorporation into N-containing biomarkers) while those with ratios lower than 1 were

175     considered low-fixing.

176     **Bacterial whole genome sequencing**

177     Each Sierra Mixe microbial isolate was recovered from cryogenic storage by streaking cells onto

178     Luria-Bertani (LB) agar medium plates and incubating for one to two days at 28 °C. Single

179     colonies were sub-cultured in liquid LB medium at 28 °C to an $OD_{600}$ value of 0.7. Genomic

180     DNA (gDNA) was extracted from the cell culture pellet of each isolate using the *Mo Bio*

181     Ultraclean Microbial DNA extraction kit (QIAGEN, Inc). Sequencing libraries were

182     subsequently constructed using the KAPA HyperPlus DNA library preparation kit (Roche, Inc)

183     by following the instructions of the technical datasheet provided. A gDNA input of 100 ng was

184     fragmented enzymatically for 9 minutes to achieve an average insert size of 450bp. The inserts

185     were ligated to customized dual-indexed barcode adapters (Integrated DNA Technologies), and

186     the library was size-selected by using KAPA Pure beads to carry out the kit's dual-SPRI protocol

187     to generate an average adapter-ligated gDNA insert molecule size of 600 bp. The size-selected

188    libraries were then PCR amplified over a total of five cycles. Average library molecular size was

189    determined using the DNA High Sensitivity Assay kit with the Agilent 2100 Bioanalyzer

190    (Agilent Technologies). The Library was then used to generate paired end reads over 150 cycles

191    at the UC Davis DNA Sequencing Technologies Core facility on the Illumina HiSeq 4000

192    system.

193    **Isolate Genome Sequence Analysis**

194    The paired-end FASTQ files of each isolate library were quality trimmed using Trimmomatic

195    0.36 using the following settings: ILLUMINACLIP:TruSeq3-PE.fa:2:40:15; LEADING: 2;

196    TRAILING:2; SLIDINGWINDOW:4:15; MINLEN:50 (12). The trimmed reads were

197    subsequently assembled using MEGAhit 1.2 with default settings (13). Assembly metrics were

198    obtained with the default settings of QUAST 4.1, the quality assessment tool for genome

199    assemblies (14), and the output for each assembly is summarized in S2 Table. Genome binning

200    analysis to assess the purity of each isolate genome was carried out using the program Metabat

201    with the default settings (15). The number of bins generated by Metabat for each isolate genome

202    are displayed in S2 Table. Values for genomic coverage were generated by aligning trimmed

203    reads to the resulting assemblies with BWA followed by the use of the depth function from

204    Samtools (16, 17). Code for the Snakemake workflow used to conduct the computational

205    analysis is available at: (https://github.com/shigdon/snakemake_mucilage-isolates).

206    **Genome distance analysis and taxonomic classification**

207    Whole genome assemblies were classified and compared using Sourmash 3.1.0 (18), which

208    provides implementation of both the MinHash and Lowest Common Ancestor (LCA) algorithms

209    to carry out whole genome comparisons and taxonomic classification of microbial isolates in a

10

210    fast, efficient and lightweight computational fashion (18-20). The complete assembly files output

211    from MEGAhit 1.2 for each isolate genome were used to generate MinHash signatures, also

212    referred to as sketches, using the program Sourmash 3.1.0 (https://github.com/dib-lab/sourmash).

213    The chosen k-mer size for each isolate genome's MinHash signature was set to 31 (k-31). These

214    sketches served as genomic fingerprint signatures that were used to carry out an all-by-all

215    comparison at the whole-genome level by using the 'compare' function of Sourmash to calculate

216    Jaccard Similarity Index (JSI) values for each pairwise comparison, which was output as a

217    matrix in csv format. This csv file was then used to generate the all by all comparative matrix

218    and associated dendrogram in Fig. 1 using the ComplexHeatmap package in R (21). For

219    taxonomic assignment of total genome assemblies, the k-31 signatures were queried against a

220    database of k-31 MinHash signatures that correspond to the curated microbial genomes within

221    the Genome Taxonomy Database (GTDB) v89 using the 'lca search' command of Sourmash

222    (available at: https://osf.io/wxf9z/). K-31 MinHash signatures were also generated using

223    Sourmash for the genome bins of each isolate genome that were created using Metabat. The

224    MinHash signature of each genome bin was classified using the 'lca search' function of

225    Sourmash using the aforementioned prepared database. Results from bin classification using

226    Sourmash are presented in S4 Table. Quantification of full taxonomies generated using

227    Sourmash LCA classification data from isolate genome bins derived was visualized as a Heat

228    Tree using MetacodeR 0.3.1 in R (22).Code used to generate, compare and classify MinHash

229    genome sketches is included in the Snakemake workflow hosted at:

230    (https://github.com/shigdon/snakemake_mucilage-isolates). Code used for analysis of Sourmash

231    output and figure generation in R is available at: (https://github.com/shigdon/R-Mucilage-

232    isolates-sourmash).

233    **Mucilage metagenome taxonomic classification**

234    Paired end Illumina sequence data from Sierra Mixe aerial root mucilage metagenome sample

235    OLMM00 was downloaded from Figshare

236    (https://figshare.com/s/04997ae7f7d18b53174a#/articles/6615497) and analyzed to characterize

237    the breadth of microbial diversity present within the mucilage environment. The shotgun

238    metagenomic reads were quality filtered using Trimmomatic 0.36 and the surviving reads were

239    separated into microbial and non-microbial fractions using the classify function of Kraken2

240    2.0.8_beta with the Refseq complete databases for Bacteria, Archaea, and Viruses (23, 24). The

241    microbial component of OLMM00 classified with Kraken2 was subsequently visualized using

242    the R package MetacodeR at the Phylum, Class, Order and Family levels, which is presented in

243    Fig S1 (25). The relative abundance of each microbial taxon classified at the genus level was

244    computed after performing Bayesian re-estimation of hits using Bracken2 (26) and normalization

245    of read classifications for each taxon with the counts per million method using the R package

246    Phyloseq (S6 Table) (27). Prior to analyzing the microbial community, the table of classified

247    microbial taxa output by Bracken2 was filtered to remove taxa for which the number of

248    classified reads was below 500, which resulted in a total of 609 unique genera identified within

249    the OLMM00 metagenome (S7 Table). Source code for analysis and figure generation is

250    available at: (https://github.com/shigdon/R-Mucilage-Metagenome).

251    *Nif* **and alternative** *nif* **gene mining**

252    Protein coding sequences were predicted for each microbial isolate genome by using the

253    corresponding MEGAhit-assembled contigs as input files for the prokaryotic genome annotation

254    program Prokka 1.12 (28). The multi-FASTA amino acid files output for each isolate genome

12

255   were scanned against profile hidden markov models (pHMMs) corresponding to *nif* genes of the

256   *K. pneumoniae* NIF regulon using the 'hmmscan' function of HMMER 3.1b (29). These were

257   acquired from the Pfam and TIGRFAM libraries of pHMMs (30, 31). HMM hits for each *nif*

258   gene were stringently filtered in R using the dplyr package to retain query-subject hits that

259   maintained model coverage greater than or equal to 75 % and a maximum e-value of $1e^{-9}$ (32).

260   Visualization of *nif* gene profiles for all pure isolates depicted in Fig 3 was achieved using the

261   Complex Heatmap package in R by clustering pure isolates based their relative MinHash

262   distances and displaying counts of unique coding sequences that were found to match each *nif*

263   HMM (21). TIGRFAMs used to scan for canonical *nif* genes of the *K. pneumoniae* NIF regulon

264   included: TIGR01817, TIGR02938, TIGR02176, TIGR01287, TIGR01282, TIGR01286,

265   TIGR01283, TIGR01285, TIGR01290, TIGR02000 TIGR03402, TIGR02660, TIGR02933 and

266   TIGR01752. Pfams used to scan for *nif* gene mining included: PF04891.11 and PF03206.13.

267   TIGRFAMs used to scan for alternative *nif* gene mining included: TIGR01860, TIGR02930,

268   TIGR02932, TIGR01861, TIGR02929 and TIGR02931. The corresponding hmmscan results for

269   alternative *nif* genes were filtered to retain query-model matches with maximum e-values of $1e^{-06}$

270   and 85 % minimum model coverage. Source code for bacterial genome mining analyses and

271   figure generation is available at: (https://github.com/shigdon/R-Mucilage-isolates-nif) and

272   (https://github.com/shigdon/R-alt-nif-analysis).

273   **CAZyme gene mining**

274   The multi-FASTA amino acid files for each microbial isolate genome that were generated by

275   Prokka were each used as input for the dbCAN2 analytical pipeline (33). This was achieved

276   using a local installation of the source code for the dbCAN2 pipeline hosted on Github

277   (https://github.com/linnabrown/run_dbcan). Output files in CSV format were read into R and

13

278    filtered using the R packages within tidyverse 1.2.1 (34). Circular heatmap plots were made

279    using the ggtree package (35). Source code for analysis and figure generation is available at:

280    (https://github.com/shigdon/R-Mucilage-isolates-dbCAN2).

**Pan-genome Analysis**

282    Genomic features predicted by Prokka (36) for each microbial genome included in the isolate

283    sub-population study were aggregated in GenBank feature format and collectively used as input

284    for pan-genome analysis using the program Roary 3.12.0 (37). Configuration for running the

285    Roary microbial pan-genomic pipeline included use of the "-e" flag to generate a multi-FASTA

286    alignment of core genes using PRANK and a minimum blastp identity value of 95 percent. To

287    visualize the pan-genome of the isolate set presented in Fig 5C, the gene presence and absence

288    output file, the associated dendrogram and an isolate-genus mapping file were uploaded to the

289    Phandango web server (38). Source code for analysis and figure generation is available at:

290    (https://github.com/shigdon/R-alt-nif-analysis).

# Results

**Diazotrophic isolates were confirmed by functional assay of $^{15}N_2$ incorporation**

293    We isolated putative diazotrophic bacteria in samples collected from Sierra Mixe maize plants

294    grown using a nitrogen-deficient basal medium supplemented with sugars corresponding to the

295    monosaccharide composition of aerial root mucilage (S1 Table). Culturing each isolate in N-

296    deficient liquid media under an atmosphere containing $^{15}N_2$ gas and measuring their ability to

297    incorporate $^{15}N$ atoms into small molecule metabolites ( i.e. <1000 Da) by Time of Flight mass

298    spectrometry confirmed that the isolates were diazotrophic and produced a large number of

299    compounds with different masses and chemical structures. Summation of peak intensities for N-

300    containing compounds common to enriched and control (compressed air) cultures enabled each

301    isolate's BNF capacity to be measured as a ratio of $^{15}N/^{14}N$ (BNF ratio). Overall, BNF ratios

302    obtained for all pure isolates assayed ranged from 0.6 to 4.6 (Table S2). While most isolates

303    exhibited moderate BNF ratios between 1 and 2, ~5% of the isolates demonstrated N-fixation

304    with BNF ratios >2 (Table 1). The observed BNF ratio variation among these confirmed

305    diazotrophs prompted investigation of the underlying genomic determinants for BNF of each

306    isolate.

307    **Whole genome analysis revealed significant phylogenetic diversity**

308    The selected bacterial isolates were subjected to WGS and resulted in a collection of draft

309    genome assemblies with fold coverages that ranged from 14 – 330X (S3 Table). Analysis of

310    mucilage isolates revealed an unexpected range of diversity in nucleotide composition and

311    taxonomy. All-by-all comparison of MinHash sketches for each isolate genome depicted the

312    relative genomic distances of all pairings that verified the diversity of genomes (Fig 1).

313    Complete taxonomic classification for each bacterial genome (S4 Table) at the maximum sketch

314    size found 33 known bacterial taxa among the 472 isolate genomes, and 116 genomes that were

315    unidentified (Fig 1). Possible explanations for unidentified isolates included lack of a database

316    accession match or the presence of multiple bacterial genomes within a WGS MinHash sketch

317    that triggered disagreement within the genomic classification structure of the lowest common

318    ancestor algorithm (LCA).

319        To assess whether isolate genomes were pure or derived from a mixed culture that

320    appeared pure during isolation, we used Metabat (15) to bin each WGS assembly and identify

321    isolates comprised of multiple organisms. This resulted in 492 isolate genomes with single bins

15

322    of single organism DNA sequences (Fig 1, S6 Table) – indicating pure cultures. WGS assemblies

323    with 2, 3, 4 and 5 bins had frequencies of 72, 19, 3 and 2, respectively (Fig 1 and S3 Table),

324    indicating that what appeared to be a single colony contained multiple organisms and that further

325    WGS analysis was needed to deconvolute respective sequences. Reexamination of the

326    deconvoluted genomes for taxonomic classification of each genome bin increased the resolution

327    of microbial diversity and augmented the diversity of the taxa present and capable of fixing

328    nitrogen (S5 Table).

329          Visualization of the classified genome bins indicated that the selected isolates were

330    primarily comprised of Proteobacteria, a substantial number of Firmicutes, and relatively few

331    Actinobacteria (S1 Fig). While deconvoluted genomes largely classified as

332    Gammaproteobacteria, relatively few deconvoluted genomes were classified to the

333    Alphaproteobacteria or Betaproteobacteria classes. Congruent with the findings of Carvalho et

334    al., several deconvoluted genomes from our study were classified as *Burkholderia*, along with

335    other Betaproteobacteria that included *Achromobacter*, *Acidovorax* and *Herbaspirillum* (39).

336    However, deconvoluted genomes classified as *Enterobacter*, *Klebsiella*, *Metakosakonia*,

337    *Rahnella*, *Raoultella,* and *Pseudomonas* were among the most abundant in the mixed cutures.

338    Membership of deconvoluted genomes classified to Firmicutes included a substantial number of

339    *Lactococcus* and several were identified as *Enterococcus* and *Bacillus*. Included in the few

340    Actinobacteria genomes sequenced, deconvoluted genome analysis found *Curtobacterium*,

341    *Leifsonia*, *Microbacterium*, *Micrococcus* and *Rhodococcus* as well.

342          Comparison of the deconvoluted genomes and pure genomes for taxonomic content with

343    the OLMM00 mucilage metagenome reported by Van Deynze et al. (6) indicated that the

344    culturing strategy enriched the isolates that fixed nitrogen and obtained a small fraction of the

16

345    possible mucilage microbiome reported from the low sequence coverage metagenome. Using

346    609 genera identified in OLMM00 as a benchmark for bacterial diversity (S7 Table), the unique

347    genera classified among isolate WGS assemblies comprised ~5% of genera in the mucilage

348    microbiome. In addition, analysis of OLMM00 metagenome provided further insight to the

349    phylogenetic diversity of mucilage microbiota associated with this landrace (S8 Table).

350    Proteobacteria, Bacteroidetes, Actinobacteria and Firmicutes were the most abundant phyla in

351    the mucilage microbiome (S2 Fig, S8 Table). However, confirmation of multiple organisms

352    contributing to mixed cultures (i.e. composite genomes) limited our ability to attribute observed

353    BNF phenotypes to a distinct organism within co-cultured isolates. This observation prompted

354    genomic profiling of each pure isolate genome for carbohydrate utilization and *nif* features to

355    address the hypothesis that mucilage diazotrophs derive energy from mucilage polysaccharide to

356    fuel BNF.

357    **Diazotrophic isolates possessed CAZymes and sugar transporters relevant for mucilage**

358    **digestion**

359    Examining isolate genomes for glycosyl hydrolase (GH) genes relevant to the composition of

360    aerial root mucilage polysaccharide (6, 40) was done using Hidden Markov Models (HMMs) of

361    GH families in the Carbohydrate Active Enzymes (CAZy) database (S9 Table) (41). This

362    analysis revealed that the pure culture diazotrophs contained genes supporting the genomic

363    potential to degrade and derive energy from mucilage polysaccharides. Targeting GHs with

364    arabinofuranosidase, fucosidase, galactosidase, glucuronidase, mannosidase and xylosidase

365    activities revealed that diazotrophic genomes with small differences in genome diversity

366    contained similar GH profiles spanning 12 functional GH groups (Fig 2A). Comparison of GH

367    groups conferring arabinofuranosidase and/or xylosidase activities demonstrated that the more

17

368     promiscuous 'Ara/Xyl' GH group had the highest abundance with increased genome copy

369     number for the majority of classified genomes. GH groups with exclusive galactose or mannose

370     substrate specificities were also abundant in the isolates examined, where the sum of the isolates

371     with genes in these GH groups was determined to be 366 of the 492 genomes (S10 Table). In

372     contrast to the plethora of genomes found to possess pentose and/or hexose cleaving GHs, those

373     with strict glucuronate and fucose specificities were far less abundant in the pure cultures.

374     Interestingly, most genomes possessed genes in GH groups with promiscuous substrate

375     specificities that encompassed the complete range of mucilage polysaccharide compositional

376     diversity across five different GH families (GH1, GH2, GH31, GH4, GH30).

377         In addition to generating GH profiles, querying genomes for the presence of sugar

378     transport genes relevant to monosaccharides that contribute to mucilage polysaccharide structure

379     revealed that isolated diazotrophs possess the machinery necessary for transport of mucilage-

380     derived monosaccharides obtained from the digestion of mucilage, indicating that the initiating

381     step of catabolism was present in the genome (Fig 2B). Utilizing a list of mucilage relevant

382     accessions (S11 Table) from the Transporter Classification Database (TCDB) (42), we generated

383     sugar transport profiles for each genome. Summarizing genome counts by genus level

384     classification demonstrated that those classified to the most common Gammaproteobacteria

385     exhibited sugar transporters for all six monosaccharide moieties derived from the mucilage

386     polysaccharide (S12 Table). Additionally, isolates of the most commonly classified genera

387     possessed multiple genes and/or mechanisms for transport of each monosaccharide type in

388     mucilage. Genomes assigned to less abundant genera tended to exhibit higher variation in sugar

389     transporter profiles, where the absence of known carbohydrate transport systems corresponding

390     to some, but not all components of mucilage polysaccharide was observed. This observation may

18

391  explain how the culturing strategy resulted in reflecting abundant members of the mucilage

392  microbiome.

**Diazotrophic isolates displayed genomic variation in canonical *nif* gene features**

394  The genetic basis for BNF was established following more than 100 years of research, where

395  numerous *nif* genes have been implicated as contributing factors to the phenotype with various

396  operon configurations (43). We investigated the genomic mechanism for the diazotrophic

397  phenotype (i.e. BNF) by examining the predicted coding sequences using HMMs for the six *nif*

398  genes of the Dos Santos model (7) within the context of seven genetic operons comprising the

399  *K. pneumoniae* NIF regulon, which included: 1) the operon of *nif* genes involved in regulation of

400  the *nif* pathway, *nif*RLA; 2) the catalytic operon, *nif*HDK; 3) operons involved in formation of

401  the functional Fe-Mo protein, *nif*EN and *nif*BQ; 4) an operon of genes involved in assembly of

402  the functional enzyme complex, *nif*USVM; and 4) operons conferring genes associated with

403  mediating electron transfer, *nif*J and *nif*WF (44, 45). Results from this extensive analysis

404  generated *nif* gene profiles and revealed three distinct groups of diazotrophic isolates (NIF

405  groups) based on *nif* gene content and variation in structure (Fig 3). NIF groups included a subset

406  of 193 genomes positive for the presence of homologous protein-coding sequences to HMMs for

407  all *nif* genes in the Dos Santos model (DS-positive, DSP), a smaller subset of 66 isolates with *nif*

408  gene profiles reflecting a semi-complete set of Dos Santos model *nif* genes (Semi-DS, SDS) and

409  a subset of 233 isolates that completely lacked genes with HMM homology for all Dos Santos

410  model *nif* genes (DS-negative, DSN), yet phenotypically displayed diazotrophy.

411       Although each NIF group included genomes classified from a range of bacterial genera,

412  each group also included diazotrophs with an "unassigned" taxonomic classification (S6 Table).

413  However, DSP genomes positively classified to known genera were comprised entirely of

19

414      Gammaproteobacteria assigned to *Enterobacter*, *Klebsiella*, *Kosakonia*, *Metakosakonia*,

415      *Pseudomonas*, *Rahnella* or *Raoultella*. SDS isolates had much higher taxonomic diversity, where

416      SDS group membership was attributed by pure isolates from Actinobacteria, Firmicutes and

417      Proteobacteria. Within these three phyla, SDS isolate genera included *Acidovorax*,

418      *Acinetobacter*, *Bacillus*, *Curtobacterium*, *Herbaspirillum*, *Leifonia*, *Micrococcus*, *Pseudomonas*

419      and *Stenotrophomonas*. In a fashion similar to SDS isolates, DSN genomes were also composed

420      of Actinobacteria, Firmicutes and Proteobacteria. Interestingly, DSN genomes displayed the

421      highest taxonomic diversity among the three NIF groups by including *Acinetobacter*,

422      *Agrobacterium*, *Atlantibacter*, *Citrobacter*, *Curtobacterium*, *Enterobacter*, *Erwinia*, *Escherichia*,

423      *Hafnia*, *Lactococcus*, *Lelliottia*, *Metakosakonia*, *Microbacterium*, *Morganella*, *Pantoea*,

424      *Pseudomonas*, *Rahnella*, *Rhodococcus*, *Serratia* and *Staphylococcus*. While genomes classified

425      as *Enterobacter*, *Metakosakonia* and *Rahnella* were found in both the DSP and DSN groups,

426      *Pseudomonas* genomes were present in all three NIF groups. In addition to *Pseudomonas*,

427      commonalities between genera identified within the SDS and DSN groups included membership

428      to *Acinetobacter* and *Curtobacterium*.

429         Every genome from a diazotroph in the DSP group possessed homologous protein coding

430      regions to *nif* genes in the *K. pneumoniae* NIF regulon (Fig 3, S3 Fig). Importantly, diazotrophs

431      in this group possessed homologs to the six *nif* genes of the Dos Santos Model and exhibited

432      BNF ratios that confirmed their ability to fix atmospheric nitrogen. The majority of diazotrophs

433      in the DSP group had moderate BNF ratio values within the inclusive range of 1 to 2, and four

434      isolates exhibited capacity ratios > 2 (Fig 4A). While the 21 *Rahnella* genomes were the only

435      subset found to possess homologs for all 16 *nif* genes investigated, the remaining 172 genomes

436      lacked homologs to either the *nif*J, *nif*L, *nif*Q or *nif*W genes in variable degrees and/or

20

437  combinations. However, these diazotrophs exhibited nearly identical *nif* gene profile

438  compositions with the exception of slight variations in gene copy number. In the case of DSP

439  isolates classified as *Enterobacteriaceae*, distinguished clades of *Enterobacter* and *Klebsiella*

440  genomes each lacked homologous genes to *nif*L and *nif*W while clades of *Pseudomonas* and

441  most *Rahnella* genomes were the only diazotrophs with homologs for the *nif*W gene. With

442  respect to the *nif*H gene encoding the dinitrogenase reductase protein, 150 genomes in the DSP

443  catagory had single copy homologs and 43 exhibited the presence of two copies. Overall, the *nif*

444  gene content and BNF ratios of diazotrophs in the DSP group demonstrated that many mucilage

445  diazotrophs adhered to the *K. pneumoniae* NIF regulon and Dos Santos models to conduct BNF.

446  All 66 members of the SDS group contained homologs to at least one, but not all, of the

447  essential *nif* genes in the Dos Santos model (Fig 3, S4 Fig) and fixed nitrogen with BNF ratios

448  similar to diazotrophs in the DSP group (Fig 4B). In a similar fashion to the DSP isolates, all

449  SDS isolates were found to possess homologs for at least one copy of the *nif*H gene but

450  interestingly two copies were detected in 15 diazotrophs of the SDS group. Genes homologous to

451  dinitrogenase component I, *nif*D and *nif*K, were only found in a single isolate of the SDS group.

452  Regarding the three *nif* genes involved with biosynthesis of FeMoCo, only a single SDS isolate

453  (BCW-200147) possessed homologs to *nif*E and *nif*N, and genes matching the HMM for *nif*B

454  were not detected in any SDS diazotroph genomes. Beyond Dos Santos' model of essential *nif*

455  genes, many SDS isolates possessed homologs for several genes in the *nif*RLA and *nif*USVM

456  operons of *K. pneumoniae*, but genes involved with electron transfer (*nif*F and *nif*J) were not

457  detected among the majority of isolates in this group. Despite lacking the complete set of *nif*

458  genes in the Dos Santos Model, BNF ratios for isolates in this group ranged from 0.8 to 3.0.

459  Taken together, *nif* gene analysis combined with the diazotrophic phenotype (i.e. BNF ratios) in

460    the SDS group revealed that many mucilage isolates exhibited BNF activity without the presence

461    of any essential *nif* genes defined by the Dos Santos model, suggesting that a novel mechanism

462    of diazotrophy may be expressed in the microbiome of this landrace.

463           Contrary to the DSP and SDS NIF groups, the 233 diazotrophs in the DSN group

464    completely lacked the presence of homologous protein coding sequences for all *nif* genes in the

465    Dos Santos model (Fig 3, S5 Fig) and exhibited BNF ratios that rivaled those of diazotrophs in

466    the other NIF groups containing gene matches to HMMs for all or part of the *nif* genes in the Dos

467    Santos model (Fig 4C). Members of the DSN group lacked homologs for many *nif* genes

468    constituting the NIF regulon of *K. pneumoniae*, and nearly all of them possessed coding

469    sequences resembling genes of the *nif*USVM operon. While many DSN genomes encoded

470    homologous genes to the BNF regulatory protein *nif*A, members of this group contained gene

471    sequences that matched the *nif*L HMM to a much lesser extent. Contrary to the observed *nif* gene

472    profiles of diazotrophs in the SDS group, observed trends for DSN genomes included presence

473    of homologous sequenecs to the *nif*F and *nif*J genes involved with electron transfer. Similar to

474    observations made with the other two NIF Groups, 188 DSN diazotophs exhibited BNF ratios

475    between 1 and 2. Surprisingly, among all three NIF Groups, the DSN group presented the largest

476    number of diazotrophs with BNF ratio values > 2. Collectively, *nif* gene profiles of DSN

477    genomes and their corresponding BNF assay results demonstrated that these diazotrophs were

478    capable of BNF without employment of any *nif* genes in the Dos Santos model and only a subset

479    of the *K. pneumoniae* NIF regulon.

480    **Alternative *nif* genes were detected in isolates with substantial genome variation**

481    Following queries for canonical *nif* genes of the Dos Santos model, we investigated whether the

482    bacterial genomes encoded *nif* genes for known alternative nitrogenase systems that either

483    strictly utilize iron (*anf*) or incorporate vanadium in place of molybdenum (*vnf*) as metal co-

484    factors of dinitrogenase (43). Utilizing TIGRFAMs for the *anf*D, *anf*K, *anf*G, *vnf*D, *vnf*K *and*

485    *vnf*G nitrogenase genes along with those of the Mo-Fe type nitrogenase (*nif*HDK), HMM

486    analysis of predicted protein sequences from each genome revealed a small subset of diazotrophs

487    with alternative *nif* genes. This resulted in the identification of 42 genomes with coding

488    sequences that matched all nine *nif* HMMs (Fig 5A). Investigation of these *nif* genes also

489    confirmed 146 diazotrophs in possession of the *nif*HDK operon without genes matching

490    alternative nitrogenase HMMs, and 63 that had genes matching only the HMM model for *nif*H.

491    Investigating the genomes with alternative *nif* genes revealed that each was previously assigned

492    to the DSP group. This observation warranted further investigation of genomic similarities and

493    differences between the 42 genomes with alternative *nif* genes.

494        WGS comparison of diazotrophs that contained alternative *nif* genes uncovered

495    substantial phylogenomic diversity within the group. Computation of genomic distances between

496    the 42 previously identified genomes revealed 12 distinct groupings of highly similar diazotrophs

497    with JSI of nearly 1 (Fig 5B). Cross-referencing previously generated taxonomy for these

498    alternative *nif*-possessing diazotrophs revealed two genera classifications. Among these

499    taxonomic assignments, 38 isolates were classified to be *Raoultella*, 2 isolates were classified as

500    *Metakosakonia*, and 2 were classified as *Enterobacteriaceae*. This indicated that the majority of

501    diazotrophs with homologs to alternative *nif* genes had genomes with significant nucleotide

502    similarity to reference genomes in the Genome Taxonomy Database (GTDB) classified as

503    *Raoultella* (46). Interestingly, diazotrophs classified as *Raoultella* exhibited broad genomic

504    diversity and formed multiple taxonomic clusters, with the two "unassigned" genomes

505    interspersed among them, suggesting that they are near relatives of *Raoultella*. Comparison of

506    the JSI values between genomes classified as *Raoultella* presented values ranging from 0.1 to 1.

507    Additionally, the two *Metakosakonia* genomes presented strong dissimilarity to the other 40

508    isolates with JSI values close to zero for each pairing. These observations indicated large

509    variation in genome composition for this subset of isolated diazotrophs and prompted subsequent

510    exploration of the pan-genome among isolates that lack classical *nif* genome construction yet fix

511    nitrogen.

512        Observed differences in nucleotide composition among genomes with alternative *nif*

513    genes were expanded by elucidating the pan-genome for this group of diazotrophs. Annotated

514    protein coding features of each genome served as inputs for pan-genome analysis to determine

515    the core genome among diazotrophic isolates with alternative *nif* genes (37). Pan-genome

516    analysis revealed a narrow core genome comprised by 285 of the 15,353 genes provided as input

517    (S13 Table) with 3,374 soft core genes, 2,532 shell genes and 9,162 cloud genes occurring

518    within 95 – 99, 15 – 95, and 0 – 15 % of diazotrophic genomes, respectively. Genome clustering

519    based on the presence and absence of annotated genomic features (Fig 5C) was highly similar to

520    that observed using MinHash, where the isolate groupings of the phylogenetic tree generated

521    using the pan-genome corresponded with clades determined using genome distance differences

522    (Fig 5B). Although taxonomic annotation of diazotrophs comprising the pan-genome suggested

523    many distinct groups of *Raoultella* genomes (annotated in green), interspersion of the two

524    "unassigned" genomes with small blocks of unique coding features (annotated in purple) among

525    the defined clades of *Raoultella* corroborated findings from the MinHash analysis with blocks of

526    core genes. Visualization of the pan-genome revealed the *Metakosakonia* clade (annotated in

527    orange) of two diazotrophs (BCW201058 and BCW201155) as a near relative to the duo of

528    distinguished *Raoultella* genomes (BCW200600 and BCW201900), which confirmed findings

24

529    from the genome distance analysis. Furthermore, these four genomes possessed large blocks of

530    features absent from the other 38 genomes in the group.

## Discussion

**Diazotrophic isolates represented a small fraction of the mucilage microbiome**

533    The strategy to isolate diazotrophs focused on simulating the native environment of aerial root

534    mucilage (anaerobic/microaerophilic, pH and temperature) in combination with nitrogen

535    deprivation. This enabled providing various carbon sources associated with the mucilage

536    polysaccharide to force expression of the metabolic traits that are likely associated with growth

537    and survival on maize during *in vitro* isolation and selection (S1 Table). This was based on the

538    two-component hypothesis that diazotrophs of the resident microbiota incorporate atmospheric

539    nitrogen into various compounds via BNF, which is biologically powered by ATP when utilizing

540    sugars derived from mucilage polysaccharides to fuel the energy needs of the energetically

541    expensive transformation. Successful generation of a large isolate collection from mucilage with

542    this strategy set the stage for further investigations to confirm the putative diazotrophic isolates.

543    In response, this study established an *in vitro* functional metabolomic assay to quantify each

544    isolate's ability to incorporate heavy nitrogen into various extracellular metabolites, which both

545    confirmed the diazotrophic nature of isolates in this collection and verified the efficacy of the

546    strategy to recover diazotrophs (Table 1, Fig 4, S2 Table).

547        WGS of nearly 600 diazotrophic isolates provided a means to assess the taxonomic

548    diversity of the isolate collection relative to that of the mucilage microbiome. Concerns of isolate

549    misclassification were avoided by using whole genome analysis and composition to assign

550    taxonomy for diazotrophic genomes rather than a conserved marker gene with higher sequence

25

551    conservation (47, 48). Utilizing Kraken to classify genera derived from normalized read counts

552    (49) of the previously reported OLMM00 mucilage metagenome (6) (S7 Table, S8 Table)

553    identified 609 genera, of which the diazotrophic genome collection had 29 in common (S5

554    Table). This revealed ~5% of the bacterial diversity from the aerial root mucilage microbiota is

555    contained within the isolate collection and demonstrated that the cultured subpopulation had 25%

556    of the top 20 most abundant known genera in the OLMM00 metagenome. Although many

557    diazotroph genomes were "unassigned" taxonomically, which highlights the potential novelty of

558    many bacteria in this isolate collection, metagenome sequencing of mucilage samples at a higher

559    depth and re-classification of isolate genomes following expansion of microbial WGS databases

560    should be achieved in the future to verify these results.

561            Comparing taxa classified in the mucilage metagenome to taxonomically classified

562    diazotroph genomes validated our strategy to recover taxa with both high relative abundance in

563    the aerial root mucilage microbiome and functionally important traits. Notably, the majority of

564    genomes in our collection were classified to the Actinobacteria, Firmicutes, and Proteobacteria

565    phyla, which strongly aligns with previous efforts to characterize plant-associated microbiomes

566    (S1 Fig, S4 Table) (50-52). Reads classified to *Pseudmonas* in OLMM00 had the highest relative

567    abundance among genera in the metagenome, and this isolate collection contained several

568    distinct clades of *Pseudomonas* based on the substantial genome dissimilarity observed from all-

569    by-all whole genome sequence comparisons (Fig 1). Whole genome taxonomic classification of

570    diazotroph genomes also revealed presence of the second most abundant genus of OLMM00,

571    *Acidovorax*, in the collection, as well as others assigned to genera with high relative abundance

572    in the mucilage metagenome that include *Agrobacterium*, *Herbaspirillum* and *Burkholderia*.

573    However, the majority of classified diazotrophs were Gammaproteobacteria that exhibited low

26

574    relative abundance in OLMM00 (S1 Fig, S2 Fig, S7 Table). This suggested that diazotrophic

575    contributions to Sierra Mixe maize by the mucilage microbiome may originate from community

576    members of lower abundance, as evidenced by the diverse set of diazotrophic isolates described

577    here. Furthermore, comparison of taxonomic analysis between whole genome sequences of

578    selected diazotrophs and the OLMM00 metagenome suggested that microbial diversity of the

579    mucilage microbiome is much broader than that of the collection. This suggests that diazotrophy

580    may not be a widespread feature among genera detected in the OLMM00 mucilage metagenome.

581    **Diazotrophs exhibited the genomic potential for mucilage polysaccharide utilization**

582    Utilizing the canonical pathway for BNF, one of the most energy-intensive biochemical

583    processes in biology that consumes 16 ATP per reaction cycle to convert a single dinitrogen

584    molecule into ammonia (53), an actively fixing diazotroph associated with Sierra Mixe maize

585    would require a reliable feedstock to produce chemical energy. Based on the diverse

586    monosaccharide composition (arabinose, fucose, galactose, glucuronate, mannose, xylose) of

587    aerial root mucilage polysaccharide (6, 54) and evidence of endogenous GH activity present in

588    fresh mucilage samples (55), we surmised that harnessing it for energy to drive BNF requires

589    bacterial genes encoding both GHs to facilitate polysaccharide catabolism, and those conferring

590    the ability to transport smaller sugars into the cell. We mined isolate genomes for carbohydrate

591    utilization genes and parsed relevant data using manually curated lists of relevant database

592    accessions (S9 Table and S11 Table) (33).

593    　　　GHs are the most abundant class of carbohydrate active enzymes (CAZymes) and consist

594    of over 150 distinct families with documented substrate specificities (56). Importantly, GHs

595    often attribute multiple substrate specificities while maintaining similar protein domain

596    architectures and sequence similarity. This ascribes the potential for substrate promiscuity among

597    GH enzymes classified to a given GH family based on differences in protein structure. The GH

598    profiles of isolate genomes indicated that mucilage diazotrophs possess the genomic potential to

599    liberate monosaccharide components of the mucilage polysaccharide (Fig 2A). A summary of

600    diazotrophic isolate counts for the number of isolates with genes in each GH group by genus

601    classification further suggested that the majority of isolated diazotroph genomes encode highly

602    specific as well as promiscuous GHs (S10 Table). These results indicated that mucilage

603    diazotrophs are capable of liberating multiple polysaccharide derivatives irrespective of

604    taxonomic assignment.

605          While the ability to liberate small carbohydrates from mucilage polysaccharide is

606    necessary for its utilization as an energy source, diazotrophs from this niche must also possess

607    the corresponding sugar transport systems. Bacteria possess multiple mechanisms for

608    monosaccharide transport that primarily consist of membrane bound permeases, symporters,

609    ABC-type porters and phosphotransferase (PTS) systems (57). We found the presence of sugar

610    transporters from these classes with specificities for all six monosaccharide derivatives of

611    mucilage polysaccharide in all of the genomes (Fig 2B, S12). Considering these findings along

612    with observations that mucilage diazotrophs possessed highly promiscuous GHs corresponding

613    to the mucilage composition, we surmised that mucilage bacteria are theoretically capable of

614    utilizing their endogenous carbohydrate utilization genes to derive energy from mucilage

615    carbohydrates. Broadly, this analysis confirmed that the majority of our diazotrophic isolates

616    possess genes that may confer the ability to derive energy from mucilage polysaccharide and

617    provides additional support for the hypothesis that diazotrophs of the mucilage microbiota utilize

618    the polysaccharide to drive BNF.

28

**Diazotrophs formed three distinct nitrogen fixation groups based on genome analysis.**

619

620 Based on the isolation strategy to enrich for diazotrophic bacteria from the mucilage microbiome

621 and the confirmed BNF phenotypes of diazotrophic isolates, we hypothesized that the

622 diazotrophic genomes contain the minimum set of *nif* genes proposed by Dos Santos (7).

623 Remarkably, the collection contained a mixture of diazotrophs that were categorized into three

624 groups: the DSP group of diazotrophs fully adherent to the Dos Santos model for essential *nif*

625 gene content, a smaller group of SDS diazotrophs with incomplete versions of the Dos Santos

626 model, and the DSN group that completely lacked all six essential *nif* genes (Fig 3, S3 Fig, S4

627 Fig and S5 Fig). While the DSP group consisted of diazotrophs that possessed homologous

628 sequences to HMMs for all six essential *nif* genes (*nif*HDKENB) of the Dos Santos model along

629 with matches to the majority of other NIF regulon genes (7), discovery of the DSN and SDS

630 isolates lacking homologous sequences to this set of canonical *nif* genes either entirely, or in-

631 part, was unexpected. Interestingly, absence of matches to the HMM for the *nif*L gene that

632 confers repression of the *nif*-specific transcriptional activator NifA in a large number of DSP

633 diazotroph genomes suggests that these isolates may be acclimatized to high frequencies of

634 nitrogen-fixing conditions in their native environment (58). Furthermore, the *nif*W gene was

635 found to be non-essential for a large number of DSP diazotrophs that lacked presence of a

636 homologous gene in their genome, which is corroborated by a previous report in *nif*W⁻ strains of

637 *K. pneumoniae* (59). However, observations that all confirmed diazotrophs in the DSP group

638 were adherent to the the well established genetic structure of the *K. pneumoniae* NIF regulon

639 (44), and that genomes classified as *Klebsiella* were only assigned to the DSP group validated

640 use of the *Klebsiella* model to examine the diazotrophic isolate genomes for canonical *nif* genes.

29

641     Taxonomic classification of diazotrophic genomes revealed a spectrum of phylogenetic

642     diversity that was not found to be indicative of *nif* gene presence. For example, while

643     gammaproteobacterial genera classified among DSP genomes included *Enterobacter*, *Klebsiella*,

644     *Kosakonia*, *Metakosakonia*, *Pseudomonas*, *Rahnella* and *Raoultella*, the SDS and DSN groups

645     contained genomes that were classified as *Enterobacter*, *Metakosakonia*, *Pseudomonas* and/or

646     *Rahnella* as well. Our discovery of diazotrophs in the DSP group classified as

647     Gammaproteobacteria suggested that bacteria of this taxonomic class from the mucilage

648     environment are likely to contribute to the BNF phenotype of Sierra Mixe maize. This is

649     supported by previous studies describing species from enterobacterial genera classified among

650     genomes in the DSP group (*Enterobacter*, *Klebsiella*, *Kosakonia*, *Rahnella*, and *Raoultella*) as

651     diazotrophic endophytes associated with cereal crops such as sugarcane, rice, and maize (60-64).

652     Recent reports demonstrated the successful engineering of a *Pseudomonas* strain capable of

653     associating with wheat and maize as a diazotrophic endophyte (65), as well as successful growth

654     promotion of maize using a diazotrophic strain of *Pseudomonas* isolated from the rhizosphere of

655     rice (66). However, to the best of our knowledge, a naturally occurring diazotrophic

656     pseudomonad associated with maize endophytically is yet to be reported. Additionally, genomes

657     in the SDS and DSN NIF groups were classified to many other genera outside of

658     Gammaproteobacteria, which indicates that diazotrophs of Sierra Mixe maize exhibit much

659     broader phylogenetic diversity relative to these previous reports of diazotrophs that associate

660     with cereal crops.

661     **Many diazotrophs exhibited high BNF ratios independent of possessing canonical *nif* genes**

662     In contrast to our hypothesis, results from the BNF assay and *nif* gene mining confirmed a

663     substantial portion of the isoated diazotrophs lacked homologous protein coding sequences to

30

664    many, or all, canonical *nif* genes of the Dos Santos and *Klebsiella* models yet exhibited high

665    BNF ratios independent of canonical *nif* genes. Our quantitative assay to detect the incorporation

666    of $^{15}$N-dinitrogen from an enriched atmosphere into secreted metabolites served as a robust

667    alternative to conventional methods of diazotrophic detection, such as colorimetric assays for

668    ammonium secretion and the acetylene reduction assay, which limit detection of evidence for

669    BNF to ammonium accumulation or secondary nitrogenase activity (i.e. production of ethylene

670    through the reduction of acetylene gas), respectively (67, 68). As there has never been a

671    documented case of diazotrophs utilizing atmospheric nitrogen without key components of the

672    nitrogenase enzyme complex, our observations that SDS and DSN diazotrophs lacked protein

673    coding sequences homologous to essential *nif* genes in their genomes (S4 Fig, S5 Fig) lead us to

674    question the metabolic mechanisms that allowed them to be successfully cultured and isolated on

675    nitrogen-free medium in the laboratory.

676        While comparison of *nif* gene profiles (Fig 3) with results from the BNF assay confirmed

677    that DSP isolates utilize atmospheric nitrogen for growth, comparison with BNF assay results for

678    the SDS and DSN NIF groups indicated that these isolates were also capable of incorporating

679    atmospheric nitrogen into secreted metabolites at efficiencies that both rivaled and exceeded

680    those of DSP isolates in some cases (Fig 4). For example, while lactococci are commonly

681    associated with plants (69), our investigation serves as the first report of diazotrophic lactococci

682    based on observations that *Lactococcus* isolates exhibited some of the highest BNF ratios (Fig

683    4C, S2 Table). These results were unexpected due to the total absence of homologous sequences

684    to HMMS for essential *nif* genes within lactococcal isolate genomes (Fig 3, S5 Fig), and

685    suggested that bacteria of the mucilage microbiota lacking essential *nif* genes are capable of

686    incorporating atmospheric nitrogen into their metabolism under N-limiting environmental

687    conditions through metabolic mechanisms outside of the Dos Santos and *Klebsiella* models.

688    Taken together, the genome analysis and BNF assay results revealed that possession of canonical

689    *nif* genes comprising the Dos Santos and *Klebsiella* models were not required for all diazotrophs

690    from Sierra Mixe maize to exhibit BNF activity, suggesting that novel diazotrophic mechanisms

691    exist in this community.

692          Uncovering the genetic underpinnings of the observed BNF phenotype for mucilage

693    diazotrophs lacking canonical *nif* genes will rely on advances in genomic analysis and future

694    experimentation. While HMMs derived from consensus sequences of full-length coding

695    sequences serve as a reliable tool to detect known genomic features in bacteria, they do not invite

696    the possibility of detecting novel protein coding sequences conferring known biological

697    functions through alternative protein domain architecture. Therefore, advances in genome

698    annotation that integrate machine learning algorithms with HMM libraries derived from

699    consensus sequences of protein domains rather than full-length coding sequences, such as

700    *Nanotext*, may enable the discovery of new proteins conferring familiar activities (70).

701    Additionally, implementation of microbial pan-genome association studies using appropriate

702    control groups for DSN isolates with confirmed BNF phenotypes may also shed light on

703    additional significant genes associated with diazotrophy (71).

704    **Alternative nitrogenase genes were not present in SDS and DSN isolate genomes**

705    We queried WGS from diazotrophic isolates for protein coding sequences homologous to known

706    alternative nitrogenase genes in search of an explanation for the discovery that confirmed

707    diazotrophic isolates lacked essential *nif* genes of the Dos Santos and *Klebsiella* models.

708    Environments with limited abundance of molybdenum often harbor diazotrophic bacteria that

709    exhibit genetic operons encoding alternative nitrogenase systems. These include Vanadium-Iron

710    (Vn-Fe) type and Iron-only type nitrogenases (Fe-Fe) that assume quaternary structure without

711    utilization of molybdenum and the assistance of an additional *nif* gene encoding the *gamma*

712    subunit for the catalytic component (43). Additionally, these operons arose over evolutionary

713    time through genetic duplication events and neofunctionalization of the Fe-Mo *nif*HDK operon

714    in response to abiotic stress (43, 53). Referencing previous reports on the nutrient deficient

715    quality of indigenous fields for Sierra Mixe maize cultivation (6), the BNF assay, and *nif* gene

716    mining results, we hypothesized that SDS and DSN diazotrophs possessed alternative *nif* genes

717    and tested it by scanning the protein coding sequences of diazotroph genomes with HMMs for

718    the Vn-Fe *nif* genes (*vnf)* and Fe-Fe *nif* genes (*anf*).

719         While results from this investigation forced the rejection of our hypothesis by confirming

720    that SDS and DSN isolates do not possess alternative *nif* genes, they did reveal discovery of a

721    subset of diazotrophs from the DSP group that possessed genes resembling the *anf* and *vnf*

722    genetic operons. We found 42 diazotrophs with genes matching TIGRFAMs from all three

723    classes of known nitrogenase systems (Fig 5A). Although unexpected, this result corroborates

724    the previous report that alternative *nif* genes were only found to occur in diazotrophs that also

725    possessed the Mo-Fe nitrogenase system (53), and the observation of alternative *nif* gene

726    sequences in Sierra Mixe mucilage (6).

727         Comparison of whole genome nucleotide composition for diazotrophs with homologs to

728    alternative *nif* genes provided evidence that this subset of the DSP NIF group exhibited

729    considerable genomic diversity and contained distinct members with resemblance to previously

730    reported *Metakosakonia* and *Raoultella* reference genomes (Fig 5B). However, this subset of

731    diazotrophic isolates exhibited high genome dissimilarity and the group was found to contain

732    genomes for which assignment to a known genus was unattainable through LCA classification

33

733     using the GTDB. These observations suggested that the mucilage microbiota harbors

734     *Metakosakonia* and *Raoultella* with alternative *nif* genes and variation in metabolic capabilities,

735     as well as potentially novel genera with considerable genomic differences. Further investigation

736     by pan-genome analysis revealed large blocks of genomic features corresponding to the variation

737     in genome composition observed in four isolate genomes that formed a distinct clade (Fig 5C).

738     To our knowledge, this is the first report of maize-associated *Raoultella* exhibiting alternative *nif*

739     genes, and the genomic evidence surrounding this discovery invites the possibility for

740     classification of a new species within the genus.

741         This work reaffirmed the proposal of Sierra Mixe maize as a model system to investigate

742     nitrogen fixation in cereal crops by validating its association with diazotrophic bacteria that

743     possess canonical genetic operons for nitrogen fixation (72). Our investigation emphasized the

744     importance of aerial root mucilage to the nitrogen-fixing phenotype of the system by confirming

745     the presence of classical nitrogen fixing bacteria in the aerial root mucilage microbiota that

746     contained the genomic potential to derive energy for BNF from mucilage polysaccharide. We

747     also demonstrated that mucilage-derived diazotrophs incorporated atmospheric nitrogen into

748     their metabolism through unknown metabolic pathways extending beyond current knowledge

749     that defines BNF as bacterial conversion of dinitrogen to ammonia through the expression of

750     canonical *nif* gene products within the Dos Santos and *Klebsiella* models. We succeeded in

751     recovering and characterizing diazotrophs from the mucilage microbiota and found diazotrophs

752     that did not contain any canonical *nif* genes, suggesting their use of novel genes for the

753     conversion of dinitrogen into organic nitrogen forms that were assimilated into many small

754     molecules exported by the organisms. Collectively, this study demonstrated that specific

755    microbiome members of Sierra Mixe maize display diazotrophy with multiple molecular

756    mechanisms.

## Funding information

765

## Author Contributions

767    BCW and MY and NK carried out the strategy to culture, isolate and store the microbial

768    collection from Sierra Mixe maize. SMH, TP and BH constructed DNA sequencing libraries for

769    WGS of bacterial isolates and BH conducted the genome sequencing. SMH carried out all

770    bioinformatic analyses related to WGS analysis with guidance from BCW and CTB. BCW and

771    RJ designed the BNF assay, established the method and analyzed the associated data, and NK

772    conducted the experiments. SMH wrote the first draft of the manuscript. ABB, BCW, CTB,

773    SMH and TP edited and revised the manuscript.

## Conflict of interest

The authors declare no conflicts of interests. None of the authors are employed by the major funding agency of this work, MARS, Inc.

## References

1. Rosenblueth M, Ormeño-Orrillo E, López-López A, Rogel MA, Reyes-Hernández BJ, Martínez-Romero JC, et al. Nitrogen Fixation in Cereals. Frontiers in Microbiology. 2018;9(1794).

2. Giller KE. Nitrogen fixation in tropical cropping systems: CABI; 2001.

3. Yusuf AA, Iwuafor ENO, Abaidoo RC, Olufajo OO, Sanginga N. Grain legume rotation benefits to maize in the northern Guinea savanna of Nigeria: fixed-nitrogen versus other rotation effects. Nutrient Cycling in Agroecosystems. 2009;84(2):129-39.

4. Triplett EW. Diazotrophic endophytes: progress and prospects for nitrogen fixation in monocots. Plant and Soil. 1996;186(1):29-38.

5. Philippot L, Raaijmakers JM, Lemanceau P, Van Der Putten WH. Going back to the roots: the microbial ecology of the rhizosphere. Nature Reviews Microbiology. 2013;11(11):789.

6. Van Deynze A, Zamora P, Delaux P-M, Heitmann C, Jayaraman D, Rajasekar S, et al. Nitrogen fixation in a landrace of maize is supported by a mucilage-associated diazotrophic microbiota. PLoS biology. 2018;16(8):e2006352.

7. Dos Santos PC, Fang Z, Mason SW, Setubal JC, Dixon R. Distribution of nitrogen fixation and nitrogenase-like sequences amongst microbial genomes. BMC Genomics. 2012;13(1):162-.

795    8.   Villas‑Bôas SG, Højer‑Pedersen J, Åkesson M, Smedsgaard J, Nielsen J. Global
796         metabolite analysis of yeast: evaluation of sample preparation methods. Yeast.
797         2005;22(14):1155-69.

798    9.   Xie Y, Chou LS, Cutler A, Weimer B. DNA macroarray profiling of Lactococcus lactis
799         subsp lactis IL1403 gene expression during environmental stresses. Applied and
800         Environmental Microbiology. 2004;70(11):6738-47.

801    10.  Draper JL, Hansen LM, Bernick DL, Abedrabbo S, Underwood JG, Kong N, et al. Fallacy
802         of the Unique Genome: Sequence Diversity within Single Helicobacter pylori Strains.
803         MBio. 2017;8(1).

804    11.  Xia JG, Sinelnikov IV, Han B, Wishart DS. MetaboAnalyst 3.0-making metabolomics more
805         meaningful. Nucleic Acids Research. 2015;43(W1):W251-W7.

806    12.  Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
807         data. Bioinformatics. 2014:btu170.

808    13.  Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node
809         solution for large and complex metagenomics assembly via succinct de Bruijn graph.
810         Bioinformatics. 2015:btv033.

811    14.  Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome
812         assemblies. Bioinformatics. 2013;29(8):1072-5.

813    15.  Kang DD, Froula J, Egan R, Wang Z. MetaBAT, an efficient tool for accurately
814         reconstructing single genomes from complex microbial communities. PeerJ. 2015;3:e1165.

815    16.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence
816         alignment/map format and SAMtools. Bioinformatics. 2009;25(16):2078-9.

817    17.  Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.
818         bioinformatics. 2009;25(14):1754-60.

37

819   18.  Brown CT, Irber L. sourmash: a library for MinHash sketching of DNA. J Open Source
820         Software. 2016;1(5):27.

821   19.  Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, et al. Mash: fast
822         genome and metagenome distance estimation using MinHash. Genome Biol.
823         2016;17(1):132.

824   20.  Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
825         alignments. Genome Biol. 2014;15(3):R46.

826   21.  Gu Z, Eils R, Schlesner MJB. Complex heatmaps reveal patterns and correlations in
827         multidimensional genomic data. 2016;32(18):2847-9.

828   22.  Foster ZS, Sharpton TJ, Grünwald NJ. Metacoder: An R package for visualization and
829         manipulation of community taxonomic diversity data. PLoS computational biology.
830         2017;13(2):e1005404.

831   23.  Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome
832         biology. 2019;20(1):257.

833   24.  Pruitt KD, Tatusova T, Maglott DR. NCBI reference sequences (RefSeq): a curated non-
834         redundant sequence database of genomes, transcripts and proteins. Nucleic acids research.
835         2007;35(suppl_1):D61-D5.

836   25.  Foster ZS, Sharpton TJ, Grunwald NJ. Metacoder: An R package for visualization and
837         manipulation of community taxonomic diversity data. PLoS Comput Biol.
838         2017;13(2):e1005404.

839   26.  Lu J, Breitwieser FP, Thielen P, Salzberg SL. Bracken: estimating species abundance in
840         metagenomics data. Peerj Computer Science. 2017;3:e104.

841   27.  McMurdie PJ, Holmes S. phyloseq: an R package for reproducible interactive analysis and
842         graphics of microbiome census data. PLoS One. 2013;8(4):e61217.

843    28. Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics.
844         2014;30(14):2068-9.

845    29. Eddy SR. HMMER: Profile hidden Markov models for biological sequence analysis. 2001.

846    30. Haft DH, Loftus BJ, Richardson DL, Yang F, Eisen JA, Paulsen IT, et al. TIGRFAMs: a
847         protein family resource for the functional identification of proteins. Nucleic acids research.
848         2001;29(1):41-3.

849    31. Bateman A, Coin L, Durbin R, Finn RD, Hollich V, Griffiths‑Jones S, et al. The Pfam
850         protein families database. Nucleic acids research. 2004;32(suppl_1):D138-D41.

851    32. Wickham H, Francois R. dplyr: A grammar of data manipulation. R package version 04.
852         2015;1:20.

853    33. Zhang H, Yohe T, Huang L, Entwistle S, Wu P, Yang Z, et al. dbCAN2: a meta server for
854         automated carbohydrate-active enzyme annotation. Nucleic Acids Res. 2018;46(W1):W95-
855         W101.

856    34. Wickham H, Averick M, Bryan J, Chang W, McGowan L, François R, et al. Welcome to the
857         Tidyverse. Journal of Open Source Software. 2019;4(43):1686.

858    35. Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. ggtree: an r package for visualization and
859         annotation of phylogenetic trees with their covariates and other associated data. Methods in
860         Ecology and Evolution. 2017;8(1):28-36.

861    36. Seemann T. Prokka: rapid prokaryotic genome annotation. Bioinformatics.
862         2014;30(14):2068-9.

863    37. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, et al. Roary: rapid large-
864         scale prokaryote pan genome analysis. Bioinformatics. 2015;31(22):3691-3.

865    38. Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, Harris SR. Phandango: an
866         interactive viewer for bacterial population genomics. Bioinformatics. 2017;34(2):292-3.

867    39.  Carvalho TL, Balsemao-Pires E, Saraiva RM, Ferreira PC, Hemerly AS. Nitrogen signalling
868          in plant interactions with associative and endophytic diazotrophic bacteria. J Exp Bot.
869          2014;65(19):5631-42.

870    40.  Amicucci MJ, Galermo AG, Guerrero A, Treves G, Nandita E, Kailemia MJ, et al. Strategy
871          for Structural Elucidation of Polysaccharides: Elucidation of a Maize Mucilage that Harbors
872          Diazotrophic Bacteria. Anal Chem. 2019;91(11):7254-65.

873    41.  Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. The
874          Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics.
875          Nucleic Acids Res. 2009;37(Database issue):D233-8.

876    42.  Saier MH, Jr., Reddy VS, Tsu BV, Ahmed MS, Li C, Moreno-Hagelsieb G. The Transporter
877          Classification Database (TCDB): recent advances. Nucleic acids research.
878          2016;44(D1):D372-D9.

879    43.  Mus F, Alleman AB, Pence N, Seefeldt LC, Peters JW. Exploring the alternatives of
880          biological nitrogen fixation. Metallomics. 2018;10(4):523-38.

881    44.  Arnold W, Rump A, Klipp W, Priefer UB, Pühler AJJomb. Nucleotide sequence of a
882          24,206-base-pair DNA fragment carrying the entire nitrogen fixation gene cluster of
883          Klebsiella pneumoniae. 1988;203(3):715-38.

884    45.  Milenkov M, Thummer R, Glöer J, Grötzinger J, Jung S, Schmitz RA. Insights into
885          membrane association of Klebsiella pneumoniae NifL under nitrogen-fixing conditions from
886          mutational analysis. Journal of bacteriology. 2011;193(3):695-705.

887    46.  Chaumeil P-A, Mussig AJ, Hugenholtz P, Parks DH. GTDB-Tk: a toolkit to classify
888          genomes with the Genome Taxonomy Database. Bioinformatics. 2019.

889    47.  Varghese NJ, Mukherjee S, Ivanova N, Konstantinidis KT, Mavrommatis K, Kyrpides NC,
890          et al. Microbial species delineation using whole genome sequences. Nucleic acids research.
891          2015:gkv657.

892    48.    Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI
893            analysis of 90K prokaryotic genomes reveals clear species boundaries. Nature
894            Communications. 2018;9(1):5114.

895    49.    Haiminen N, Edlund S, Chambliss D, Kunitomi M, Weimer BC, Ganesan B, et al. Food
896            authentication from shotgun sequencing reads with an application on high protein powders.
897            NPJ science of food. 2019;3(1):1-11.

898    50.    Bulgarelli D, Schlaeppi K, Spaepen S, Van Themaat EVL, Schulze-Lefert P. Structure and
899            functions of the bacterial microbiota of plants. Annual review of plant biology.
900            2013;64:807-38.

901    51.    Hardoim P, Nissinen R, van Elsas JD. Ecology of bacterial endophytes in sustainable
902            agriculture.  Bacteria in Agrobiology: Plant Probiotics: Springer; 2012. p. 97-126.

903    52.    Levy A, Gonzalez IS, Mittelviefhaus M, Clingenpeel S, Paredes SH, Miao JM, et al.
904            Genomic features of bacterial adaptation to plants. Nature Genetics. 2018;50(1):138-150.

905    53.    Raymond J, Siefert JL, Staples CR, Blankenship RE. The natural history of nitrogen
906            fixation. Molecular biology and evolution. 2004;21(3):541-54.

907    54.    Amicucci MJ, Galermo AG, Guerrero A, Treves G, Nandita E, Kailemia MJ, et al. Strategy
908            for Structural Elucidation of Polysaccharides: Elucidation of a Maize Mucilage that Harbors
909            Diazotrophic Bacteria. Analytical Chemistry. 2019.

910    55.    Pozzo T, Higdon SM, Pattathil S, Hahn MG, Bennett AB. Characterization of novel
911            glycosyl hydrolases discovered by cell wall glycan directed monoclonal antibody screening
912            and metagenome analysis of maize aerial root mucilage. PloS one. 2018;13(9):e0204525.

913    56.    Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. The
914            Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics.
915            Nucleic Acids Research. 2008;37(suppl_1):D233-D8.

916    57.    Saier Jr MH. Families of transmembrane sugar transport proteins: MicroReview. Molecular
917            microbiology. 2000;35(4):699-710.

918    58.    Milenkov M, Thummer R, Gloer J, Grotzinger J, Jung S, Schmitz RA. Insights into
919           Membrane Association of Klebsiella pneumoniae NifL under Nitrogen-Fixing Conditions
920           from Mutational Analysis. Journal of Bacteriology. 2011;193(3):695-705.

921    59.    PAUL W, MERRICK M. The roles of the nifW, nifZ and nifM genes of Klebsiella
922           pneumoniae in nitrogenase biosynthesis. European journal of biochemistry.
923           1989;178(3):675-82.

924    60.    Chen M, Zhu B, Lin L, Yang L, Li Y, An Q. Complete genome sequence of Kosakonia
925           sacchari type strain SP1(T.). Standards in genomic sciences. 2014;9(3):1311-8.

926    61.    Govindarajan M, Kwon S-W, Weon H-Y. Isolation, molecular characterization and growth-
927           promoting activities of endophytic sugarcane diazotroph Klebsiella sp. GR9. World Journal
928           of Microbiology and Biotechnology. 2007;23(7):997-1006.

929    62.    Andreozzi A, Prieto P, Mercado-Blanco J, Monaco S, Zampieri E, Romano S, et al.
930           Efficient colonization of the endophytes Herbaspirillum huttiense RCA24 and Enterobacter
931           cloacae RCA25 influences the physiological parameters of Oryza sativa L. cv. Baldo rice.
932           Environmental Microbiology. 2019;21(9):3489-504.

933    63.    Kandel SL, Joubert PM, Doty SL. Bacterial Endophyte Colonization and Distribution within
934           Plants. Microorganisms. 2017;5(4):77.

935    64.    Luo T, Ou‐Yang XQ, Yang LT, Li YR, Song XP, Zhang GM, et al. Raoultella sp. strain
936           L03 fixes N2 in association with micropropagated sugarcane plants. Journal of basic
937           microbiology. 2016;56(8):934-40.

938    65.    Fox AR, Soto G, Valverde C, Russo D, Lagares Jr A, Zorreguieta Á, et al. Major cereal
939           crops benefit from biological nitrogen fixation when inoculated with the nitrogen-fixing
940           bacterium Pseudomonas protegens Pf-5 X940. Environmental Microbiology.
941           2016;18(10):3522-34.

942    66. Ke X, Feng S, Wang J, Lu W, Zhang W, Chen M, et al. Effect of inoculation with nitrogen-
943          fixing bacterium Pseudomonas stutzeri A1501 on maize plant growth and the microbiome
944          indigenous to the rhizosphere. Syst Appl Microbiol. 2019;42(2):248-60.

945    67. Hardy RWF, Burns RC, Holsten RD. Applications of the acetylene-ethylene assay for
946          measurement of nitrogen fixation. Soil Biology and Biochemistry. 1973;5(1):47-81.

947    68. Shand CA, Williams BL, Coutts G. Determination of N-species in soil extracts using
948          microplate techniques. Talanta. 2008;74(4):648-54.

949    69. Song AA-L, In LLA, Lim SHE, Rahim RA. A review on Lactococcus lactis: from food to
950          factory. Microbial cell factories. 2017;16(1):55-.

951    70. Viehweger A, Krautwurst S, König B, Marz M. Distributed representations of protein
952          domains and genomes and their compositionality. bioRxiv. 2019:524280.

953    71. Brynildsrud O, Bohlin J, Scheffer L, Eldholm V. Rapid scoring of genes in microbial pan-
954          genome-wide association studies with Scoary. Genome Biology. 2016;17(1):238.

955    72. Bennett AB, Pankievicz VCS, Ané J-M. A Model for Nitrogen Fixation in Cereal Crops.
956          Trends in Plant Science. 2020.

957    73. Hadfield J, Croucher NJ, Goater RJ, Abudahab K, Aanensen DM, Harris SR. Phandango: an
958          interactive viewer for bacterial population genomics. Bioinformatics. 2018;34(2):292-3.

959

960

961  **List of Tables**

962  **Table 1. Summary of BNF Assay Results.** Isolates were grouped using defined ranges of

963  $^{15}N/^{14}N$ ratio values. $^{15}N/^{14}N$ ratios were computed by summing the peak intensities of all N-

964  containing bio-markers common to both enriched and control cultures that had q-values less than

965  or equal to 0.05 after analyzing metabolite data using Metaboanalyst (11).

| BNF Group | N Isolates | $^{15}N/^{14}N$ Ratio |
|-----------|-----------|-----------------------|
| A | 4 | $x > 4$ |
| B | 10 | $4 > x > 3$ |
| C | 14 | $3 > x > 2$ |
| D | 461 | $2 > x > 1$ |
| E | 85 | $x < 1$ |
| F | 14 | Not determined |

966

967

968     **List of Figures**

969     **Fig 1. Comparative analysis of draft genome assemblies from Sierra Mixe bacterial isolates.**

970     All-by-all comparison of MinHash sketches of draft genome assemblies from 588 bacterial

971     isolates using Sourmash (18) . MinHash sketches of each draft genome assembly used in the

972     comparison had a k-mer size of 31. Genus classification from MinHash sketches for each isolate

973     genome is presented as a color-coded sidebar alongside the matrix. Results from genome binning

974     analysis with Metabat (15) is included as a second color-coded sidebar. The Jaccard Index scale

975     represents the Jaccard Similarity Index (JSI) value computed for each pairwise comparison of

976     isolate genome MinHash sketches. Darker coloring indicates higher genome similarity and

977     lighter coloring indicates lower similarity.

978

979     **Fig 2. Glycosyl hydrolase and sugar transporter genome profiles of diazotrophic isolates.**

980     Analysis using dbCAN2 (33) was done to query total predicted coding sequences in each

981     genome. Gene sequences encoding CAZymes and sugar transporters with substrate specificities

982     that correspond to monosaccharide residues of the Sierra Mixe aerial root mucilage

983     polysaccharide were selected from query results by generating a manually curated list of CAZy

984     HMMs and TCDB accession IDs. Predicted gene sequence-HMM matching pairs were reported

985     after filtering total hits from each genome to select all records with > 85% model coverage and

986     an e-value $\leq 1e^{-09}$. A) Glycosyl Hydrolase family HMM hits with designated sugar residue

987     specificities: Ara – Arabinose, Gal – Galactose, GlcA – Glucuronic Acid, Fuc – Fucose, Man –

988     Mannose, Xyl – Xylose. B) Sugar Transporter HMM-Gene hits with designated sugar residue

989     transporter activity.

990　**Fig 3. Canonical *nif* gene profiles of diazotrophic isolate genomes.** Total predicted protein

991　sequences of each pure isolate genome were queried against Hidden Markov Models (HMMs)

992　for genes of the *K. pneumoniae* NIF regulon – including the six essential *nif* genes of the Dos

993　Santos (DS) Model. Pure isolate genomes were clustered based on their relative MinHash

994　genomic distances followed by heatmap visualization of their associated *nif* gene profiles. Three

995　groups of pure diazotrophic isolates were formed based on the detected presence of homologous

996　protein coding sequences to *nif*HDKENB: DS-Positive (DSP; red), Semi-DS (SDS; green) and

997　DS-Negative (DSN; blue). Predicted amino acid sequence queries for each genome were

998　considered as matches if *nif* gene HMM coverage was greater than or equal to 75% along with e-
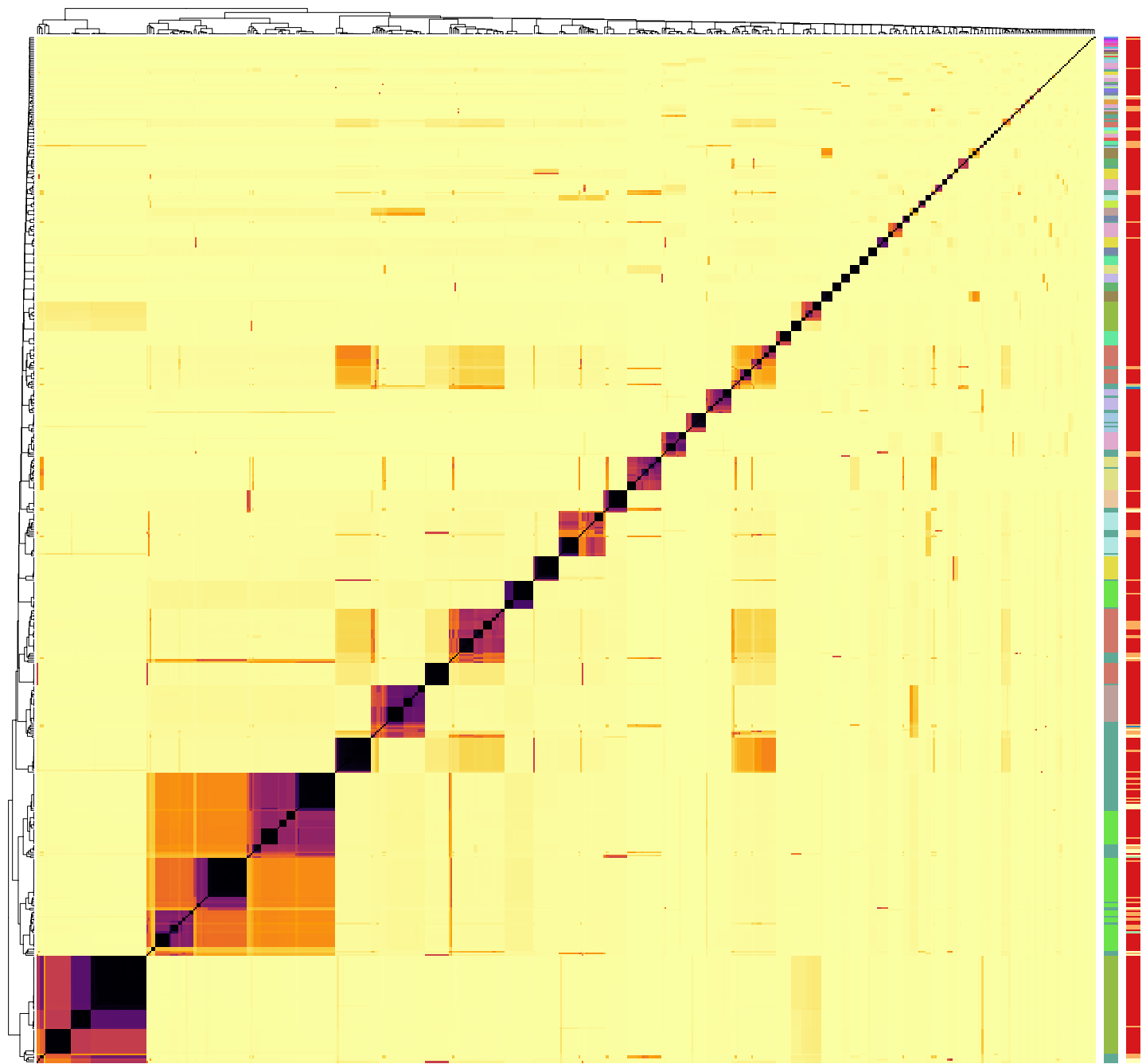
999　values $\leq 1e^{-9}$.

1000

1001　**Fig 4. BNF ratios of mucilage diazotrophs from atmospheric $^{15}N_2$ incorporation assay.** As a

1002　means to connect each diazotroph's *nif* gene profile with its corresponding BNF phenotype, BNF

1003　ratios are presented in heatmaps that accompany dendrograms clustered by MinHash genome

1004　distance under the context of the three NIF Groups determined from the genome mining analysis

1005　(Fig 3). Heatmap annotations indicate the $^{15}N/^{14}N$ ratios (BNF ratios) that represent the

1006　summation of peak intensities for all N-containing metabolites used as biomarkers in the assay.

1007　A) Dos-Santos Positive (DSP) isolates; B) Semi-Dos Santos (SDS) isolates; C) Dos-Santos

1008　Negative (DSN) isolates. Grey bars on the BNF ratio heatmap indicate values that were not

1009　determined.

1010

1011　**Fig 5. Mucilage bacterial isolates exhibit alternative nitrogenase genes.** A) The presence of

1012　predicted protein sequences in diazotrophic isolate genomes was detected using TIGRFAM

1013 HMMs corresponding to the Fe-Fe and Vn-Fe alternative *nif* genes (*anf*D, *anf*K, *anf*G, *vnf*D,

1014 *vnf*K, *vnf*G) along with HMMs for *nif*HDK. Genomes with detected presence of the targeted

1015 genes were compared and quantified using a Venn Diagram to determine the list of diazotrophs

1016 with genes resembling Vn-Fe, Fe-Only, Mo-Fe Type nitrogenases and *nif*H. B) Genomes with

1017 alternative *nif* genes were compared using Sourmash and visualized as a composite matrix that

1018 included annotation of genus level classification. C) Pan-genome analysis of diazotrophs with

1019 alternative *nif* genes was conducted using Roary (37) and data for gene presence and absence

1020 was visualized using Phandango (73) along with genus classification data from Sourmash LCA.

1021 Orange annotations indicate genomes classified as *Metakosakonia*, green annotations indicate

1022 *Raoultella* isolates, and purple annotations indicate "unassigned" classification at the genus

1023 level.
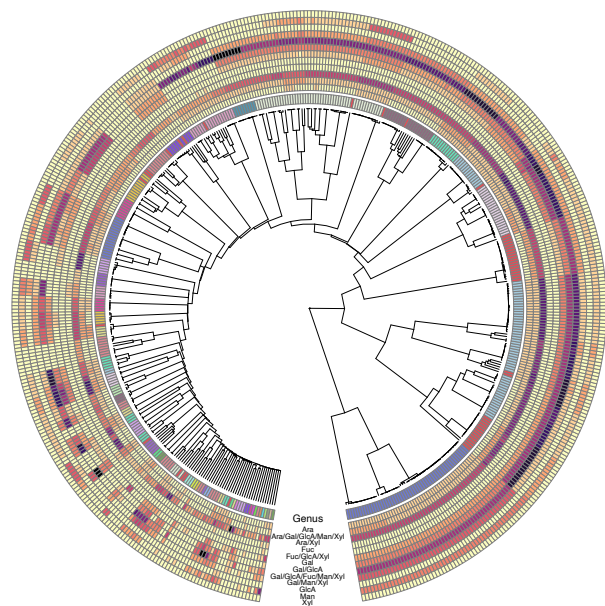
Jaccard Index

0    0.2   0.4   0.6   0.8    1

**Genus**

| | | | | | |
|---|---|---|---|---|---|
| Achromobacter | Atlantibacter | Curtobacterium | Escherichia | Kosakonia | Lelliottia |
| Acidovorax | Bacillus | Enterobacter | Hafnia | Lactococcus | Metakosakonia |
| Acinetobacter | Burkholderia | Enterococcus | Herbaspirillum | Leclercia | Microbacterium |
| Agrobacterium | Citrobacter | Erwinia | Klebsiella | Leifsonia | Micrococcus |

| | | |
|---|---|---|
| Morganella | Raoultella | Stenotrophomonas |
| Pantoea | Rhodococcus | unassigned |
| Pseudomonas | Serratia | |
| Rahnella | Staphylococcus | |

**N Bins**

1
2
3
4
5

**A**

Genus
Ara
Ara/Gal/GlcA/Man/Xyl
Ara/Xyl
Fuc
Fuc/GlcA/Xyl
Gal
Gal/GlcA
Gal/GlcA/Fuc/Man/Xyl
Gal/Man/Xyl
GlcA
Man
Xyl

**B**

Genus
Arabinose
Fucose
Galactose
Glucuronate
Mannose
Xylose

**Glycosyl Hydrolase HMM Hits**

0 1 2 3 4 5 6 7 8 9 10 11

**Genus**

Acidovorax          Citrobacter          Hafnia             Leifsonia           Morganella          Rhodococcus
Acinetobacter       Curtobacterium       Herbaspirillum     Lelliottia          Pantoea             Serratia
Agrobacterium       Enterobacter         Klebsiella         Metakosakonia       Pseudomonas         Staphylococcus
Atlantibacter       Erwinia              Kosakonia          Microbacterium      Rahnella            Stenotrophomonas
Bacillus            Escherichia          Lactococcus        Micrococcus         Raoultella          unassigned

**TCDB Diamond Hits**

0 1 2 3 4 5 6 7 8 9 10

**A**

**B**

**C**

**Genus**

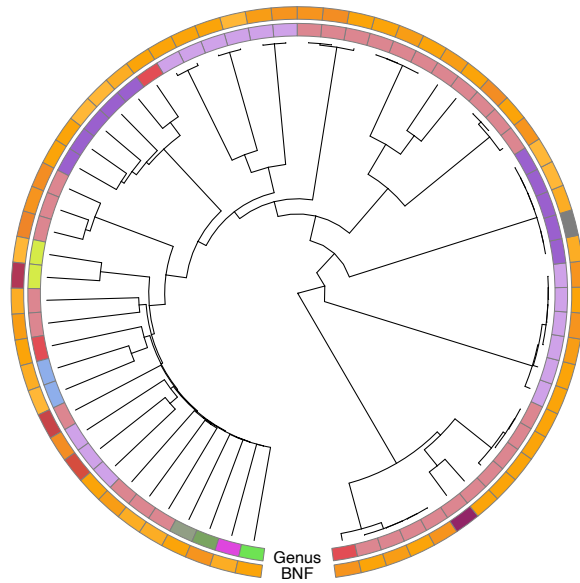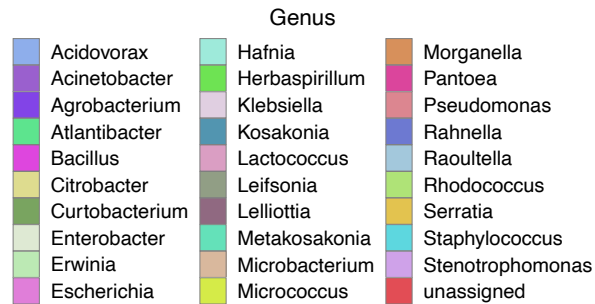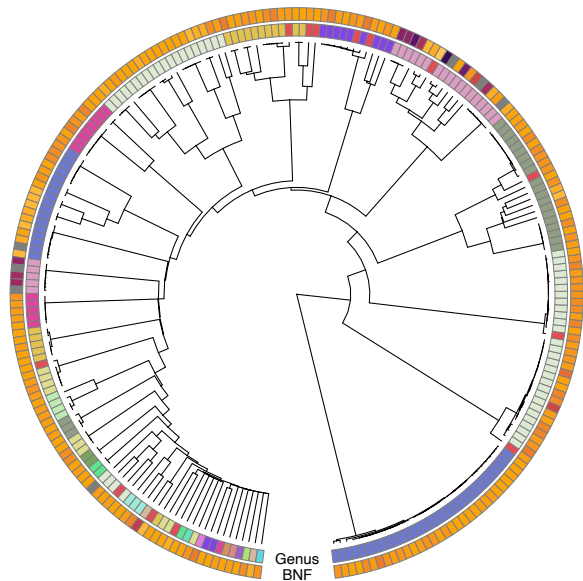| Genus | | |
|---|---|---|
| ■ Acidovorax | ■ Hafnia | ■ Morganella |
| ■ Acinetobacter | ■ Herbaspirillum | ■ Pantoea |
| ■ Agrobacterium | ■ Klebsiella | ■ Pseudomonas |
| ■ Atlantibacter | ■ Kosakonia | ■ Rahnella |
| ■ Bacillus | ■ Lactococcus | ■ Raoultella |
| ■ Citrobacter | ■ Leifsonia | ■ Rhodococcus |
| ■ Curtobacterium | ■ Lelliottia | ■ Serratia |
| ■ Enterobacter | ■ Metakosakonia | ■ Staphylococcus |
| ■ Erwinia | ■ Microbacterium | ■ Stenotrophomonas |
| ■ Escherichia | ■ Micrococcus | ■ unassigned |

**BNF (15N/14N Ratio)**

0    1    2    3    4    5