1    *Title:*
2    **Genome reduction in an abundant and ubiquitous soil bacterial lineage**
3
4    *Author affiliation:*
5    Tess E Brewer[1, 2], Kim M Handley[3], Paul Carini[1], Jack A Gilbert[4, 5], Noah Fierer[1, 6, *]
6
7    [1]Cooperative Institute for Research in Environmental Sciences, University of
8    Colorado, Boulder, CO 80309; [2]Department of Molecular, Cellular, and
9    Developmental Biology, University of Colorado, Boulder, CO 80309;
10   [3]School of Biological Sciences, The University of Auckland, Auckland 1142, New
11   Zealand; [4]Department of Ecology and Evolution, The University of Chicago,
12   Chicago, IL 60637; [5]Argonne National Laboratory, Institute for Genomic and
13   Systems Biology, Argonne, IL 60439; [6]Department of Ecology and Evolutionary
14   Biology, University of Colorado, Boulder, CO 80309
15
16   *Corresponding author:*
17   Noah Fierer
18   noah.fierer@colorado.edu
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

## **Abstract**

Although bacteria within the *Verrucomicrobia* phylum are pervasive in soils around the world, they are underrepresented in both isolate collections and genomic databases. Here we describe a single verrucomicrobial phylotype within the class *Spartobacteria* that is not closely related to any previously described taxa. We examined >1000 soils and found this spartobacterial phylotype to be ubiquitous and consistently one of the most abundant soil bacterial phylotypes, particularly in grasslands, where it was typically the most abundant phylotype. We reconstructed a nearly complete genome of this phylotype from a soil metagenome for which we propose the provisional name '*Candidatus* Udaeobacter copiosus'. The *Ca.* U. copiosus genome is unusually small for soil bacteria, estimated to be only 2.81 Mbp compared to the predicted effective mean genome size of 4.74 Mbp for soil bacteria. Metabolic reconstruction suggests that *Ca.* U. copiosus is an aerobic chemoorganoheterotroph with numerous amino acid and vitamin auxotrophies. The large population size, relatively small genome and multiple putative auxotrophies characteristic of *Ca.* U. copiosus suggests that it may be undergoing streamlining selection to minimize cellular architecture, a phenomenon previously thought to be restricted to aquatic bacteria. Although many soil bacteria need relatively large, complex genomes to be successful in soil, *Ca.* U. copiosus appears to have identified an alternate strategy, sacrificing metabolic versatility for efficiency to become dominant in the soil environment.

## **Introduction**

Soils harbor massive amounts of undescribed microbial diversity. For example, more than 120,000 unique bacterial and archaeal taxa were found in surface soils of Central Park in New York City, of which only ~15% had 16S rRNA gene sequences matching those contained in reference databases and <1% had representative genome sequence information (1). This undescribed soil microbial diversity is not evenly distributed across the tree of life. For example, *Acidobacteria* and *Verrucomicrobia*, two of the more abundant bacterial phyla found in soil (2, 3), represent only 0.08% and 0.06% of all cultured bacterial isolates in the Ribosomal Database Project (RDP, 4) and only 0.08% and 0.14% of publicly-available bacterial genomes found in Integrated Microbial Genomes (IMG, 5), respectively. Although the ecology and genomic attributes of abundant soil taxa are beginning to be described (6), we still lack basic information on the vast majority of soil microbes. These knowledge gaps highlight that a huge fraction of living biomass in terrestrial systems remains undescribed (7) and that we are only beginning to identify the influence of specific microbes on soil biogeochemistry and fertility.

For this study, we focus our exploration of undescribed microbial diversity on the *Verrucomicrobia* phylum. Although *Verrucomicrobia* are generally recognized

2

90 as being among the most numerically abundant taxa in soil (2, 3), we know very
91 little about the ecological or genomic attributes that contribute to their success.
92 The phylum *Verrucomicrobia* is highly diverse and its members possess a broad
93 range of metabolic capabilities. For example, members of the class
94 *Methylacidiphilae* are nitrogen-fixing acidophiles capable of methane oxidation
95 (8) while *Akkermansia muciniphila* of the class *Verrucomicrobiae* is a mucin-
96 degrading resident of the human gut linked to reduced host obesity (9).
97 However, the dominant *Verrucomicrobia* found in soil typically belong to the
98 class *Spartobacteria*. For example, while *Verrucomicrobia* accounted for >50%
99 of all bacterial 16S rRNA gene sequences in tallgrass prairie soils in the United
100 States, >75% of these sequences were assigned to the class *Spartobacteria*
101 (10). Currently, the class *Spartobacteria* contains only a single described and
102 sequenced isolate, *Chthoniobacter flavus*, a slow-growing aerobic
103 chemoorganoheterotroph capable of using common components of plant
104 biomass for growth (11, 12). While *Spartobacteria* are prevalent in soils, they
105 have also been observed in marine systems ('*Spartobacteria baltica*', 13) and as
106 nematode symbionts (genus *Xiphinematobacter*, 14).
107
108 Here we report the distribution of a dominant *Spartobacteria* lineage, compiling
109 data from both amplicon and shotgun metagenomic 16S rRNA gene surveys to
110 quantify its relative abundance across >1000 unique soils. We assembled a
111 near-complete genome of this lineage from a single soil where it was
112 exceptionally abundant. These results provide our first glimpse into the
113 phylogeny, ecology, and potential physiological traits of a dominant soil
114 Verrucomicrobia and suggest that members of this group are efficient at
115 growing and persisting in the low resource conditions common in many soil
116 microenvironments.
117
118 **Results and Discussion**
119
120 ***Distribution of the dominant Verrucomicrobia in soil***
121 A single spartobacterial clade dominates bacterial communities found in a wide
122 range of soil types across the globe. One phylotype from this group of
123 *Spartobacteria* represented up to 31% of total 16S rRNA gene sequences
124 recovered from prairie soils (10). This phylotype shares 99% 16S rRNA gene
125 sequence identity with a ribosomal clone named 'Da101', first described in 1998
126 as a particularly abundant 16S rRNA sequence recovered from a grassland soil
127 in the Netherlands (15). To determine if the Da101 phylotype (termed 'Da101'
128 herein) is abundant in other soils, we re-analyzed amplicon 16S rRNA gene
129 sequence data obtained from >1000 soils representing a wide range of soil and
130 site characteristics (Table S1). We found that Da101 was on average ranked
131 within the top two most abundant bacterial phylotypes in each study (Fig. 1). In
132 over 70% of the soils analyzed Da101 was within the top ten most abundant
133 phylotypes. Interestingly, phylotypes belonging to the same family as Da101

134    (Chthoniobacteraceae) were also found within the top 5 most abundant
135    phylotypes of several studies (Fig. 1).
136

137    As some 16S rRNA gene PCR primer sets can misestimate the relative
138    abundance of *Verrucomicrobia* (3, 16), we investigated whether the apparent
139    numerical dominance of Da101 in the amplicon datasets was a product of PCR
140    primer biases. To do so, we quantified the abundance of Da101 16S rRNA
141    genes within 75 previously published soil shotgun metagenomes (17, 18). The
142    relative abundance of Da101 in amplicon data was reasonably well correlated
143    with the relative abundance of Da101 determined from shotgun metagenomic
144    data ($P<0.0001$, ρ=0.50). Confirming the amplicon-based results (Fig. 1), we
145    found that Da101 was among the most abundant phylotypes observed in the
146    soil bacterial communities characterized via shotgun metagenomic sequencing
147    (Fig. S1). Thus, we conclude that the numerical dominance of Da101 in soils is
148    not simply a product of primer biases.
149

150    Despite Da101 being one of the most abundant phylotypes found in soil, its
151    proportional abundance can vary significantly across soil types (Fig. 1 and S1).
152    We used metadata associated with each soil sample to determine which of the
153    measured soil and site characteristics best predicted the relative abundance of
154    Da101. We found that Da101 was significantly more abundant in grassland soils
155    than in forest soils *($P<0.0001$, two tailed *t* test, Fig. S2); on average, Da101 is six
156    times more abundant in grassland soils. These findings suggest that the soils in
157    which Da101 excels do not overlap with those forest soils dominated by non-
158    symbiotic *Bradyrhizobium* taxa (6). Across the grassland soils included in our
159    meta-analysis, the relative abundance of Da101 was positively correlated with
160    both soil microbial biomass ($P<0.0001$, ρ=0.57, Fig. S3), and aboveground plant
161    biomass ($P<0.0001$, ρ=0.47, Fig. S3). Together, these results indicate that
162    Da101 prefers soils receiving elevated amounts of labile carbon inputs. We did
163    not identify any consistent correlations between the abundance of Da101 and
164    other prokaryotic or eukaryotic taxa, suggesting that Da101 is unlikely to be a
165    part of an obligate pathogenic or symbiotic relationship.
166

167    ***Diversity of soil Verrucomicrobia***
168    We determined the phylogenetic placement of Da101 and other soil
169    *Verrucomicrobia* by assembling near full-length 16S rRNA gene sequences from
170    six distinct grassland soils collected from multiple continents (Fig. 2, Table S2)
171    using EMIRGE (19). Although we were able to assemble representative 16S
172    rRNA gene sequences from all verrucomicrobial classes except
173    *Methylacidiphilae*, 93% of verrucomicrobial sequences fell within the
174    *Spartobacteria* class and 87% of these fell within the Da101 clade. These
175    phylogenetic analyses confirm that Da101 belongs to the *Spartobacteria* class
176    (Fig. 2). However, within the Spartobacteria class, the Da101 clade is clearly
177    distinct from the clade containing *Chthoniobacter flavus* (11, 12), as Da101
178    shares only 92%16S rRNA gene sequence identity with *C. flavus*. These findings

179     indicate that Da101 is likely a representative of a new verrucomicrobial genus.
180     We propose the candidate genus name 'Candidatus Udaeobacter' for the Da101
181     clade; the proposed name combines Udaeus ('of the earth', Greek) with bacter
182     (rod or staff, Greek), and like *Chthoniobacter* refers to one of the Spartoi of the
183     Cadmus myth. We recommend the provisional name 'Candidatus Udaeobacter
184     copiosus' for the Da101 phylotype, which refers to its numerical dominance in
185     soil.
186

### 187     Draft genome of 'Candidatus Udaeobacter copiosus' recovered from
### 188     metagenomic data

189     Despite their ubiquity and abundance in soil, there is no genomic data currently
190     available for any representative of the 'Candidatus Udaeobacter' clade.
191     Typically, soil hyper-diversity confounds the assembly of genomes from
192     metagenomes (20), requiring single-cell analysis or laboratory isolation to
193     produce an assembled genome. However, we leveraged the sheer abundance
194     of *Ca.* U. copiosus in an individual soil to obtain a nearly complete genome from
195     shotgun metagenomic data. We deeply sequenced a soil where *Ca.* U. copiosus
196     accounted for >30% of 16S rRNA gene sequences and assembled a draft
197     genome from the resulting metagenome. We used GC content, coverage,
198     tetranucleotide frequencies, and the phylogenetic affiliation of predicted proteins
199     to bin assembled contigs, resulting in a draft *Ca.* U. copiosus genome with 238
200     contigs. The draft genome is 2.65 Mbp in size, has a GC content of 54%, and
201     encodes for 3,042 predicted proteins, 67% of which could be assigned to Pfam
202     protein families (21) by the IMG annotation pipeline (5). We estimate that the full
203     *Ca*. U copiosus genome is 2.81 Mbp in length based on the recovery of 94% of
204     single copy housekeeping genes (34 of 36 genes, Table S3) that are commonly
205     used to estimate genome completion (22). The *Ca.* U. copiosus genome shares
206     only 69.3% average nucleotide identity (23) with the genome of its closest
207     sequenced relative *C. flavus*, further supporting its proposed placement as the
208     distinct genus 'Candidatus Udaeobacter'. Additionally, the *Ca.* U copiosus 16S
209     rRNA gene has 100% identity to the Da101 phylotype sequence, indicating that
210     this genome is indeed a representative of the aforementioned dominant Da101
211     clade.
212

213     *Ca.* U. copiosus has a particularly small genome size compared to *C. flavus*
214     (2.81 Mbp to 8.80 Mbp, predicted genome sizes). To see how the genome of
215     *Ca.* U. copiosus compares to other soil bacteria, we compiled data from 378 soil
216     bacteria with finished or permanent draft genome sequences in IMG whose 16S
217     rRNA gene sequences matched the 16S rRNA gene amplicon sequences
218     obtained by Leff et al. (2015) with at least 99% identity. Nearly all (99%) of these
219     378 genomes came from cultivated taxa. We estimated the genome
220     completeness for each of the 378 taxa using the same method as for *Ca*. U.
221     copiosus and found the mean estimated genome size of these taxa to be 5.28 ±
222     2.15 Mbp (mean ± SD), which is nearly identical to metagenomic based
223     estimates of mean genome size for soil microbes (24). Strikingly, the 2.81 Mbp

224 genome of *Ca*. U. copiosus is ~50% smaller than the mean genome size of
225 these 378 taxa and only 13% of these genomes were smaller than the genome
226 of *Ca*. U. copiosus.
227
228 Although soil bacteria with larger genomes tend to be more common in soil, *Ca*.
229 U. copiosus is a notable exception to this pattern. We linked the genome size of
230 each of the matched IMG bacterial genomes with the average abundance of
231 their corresponding amplicon sequence from Leff et al. (2015) and found that
232 genome size is positively correlated with average relative abundance (*P* <0.001,
233 ρ=0.37, Fig. 3). That is, bacteria with large genomes tend to comprise a
234 significantly larger proportion of soil bacterial communities. On average, the
235 genomes of soil prokaryotes are larger than those inhabiting aquatic
236 ecosystems (25) or the human gut (26). These relatively large genomes are
237 thought to provide soil-dwelling bacteria with a more diverse genetic inventory
238 to enhance survival in conditions where resources are diverse, but sparse (27,
239 28). However, the *Ca*. U. copiosus genome has a conspicuously reduced
240 genome given its numerical abundance (Fig. 3). This suggests that *Ca.* U.
241 copiosus occupies a niche space that does not require expansive functional
242 diversity and points to an alternative route to success for soil bacteria. These
243 results also suggest that abundant, uncultivated soil bacteria may have smaller
244 genomes than the cultivated taxa that represent the vast majority of available
245 genomic data. A similar pattern has been observed in aquatic systems where
246 uncultivated taxa often have smaller genomes than cultivated taxa (29). Because
247 the majority of available genomic information is derived from cultivated bacterial
248 taxa, the lack of genomic information from bacteria with reduced genomes likely
249 stems from challenges associated with culturing taxa with reduced genomes
250 (25).
251
252 Metabolic reconstruction of the *Ca*. U. copiosus genome points to an aerobic
253 chemoorganoheterotrophic lifestyle with the capacity to use a limited range of
254 carbon substrates for growth including glucose, pyruvate, and chitobiose.
255 Glycogen/starch synthesis and utilization genes were identified (*glgABCP* and
256 *amyA*), suggesting that *Ca*. U. copiosus has the capacity to store surplus carbon
257 as glycogen or starch. Glycogen metabolism has been demonstrated in other
258 *Verrucomicrobia* (30). Genes encoding for the complete biosynthesis of vitamins
259 $B_2$, $B_3$, $B_5$ (from valine) and $B_6$ were recovered. Full biosynthetic pathways for *de*
260 *novo* synthesis of alanine, aspartate, asparginine, glutamate, glutamine and
261 proline were also present. Nearly complete pathways were recovered for
262 glycine, threonine and methionine biosynthesis. Genes encoding for the
263 conversion of methionine to cysteine were present as the only apparent route to
264 cysteine biosynthesis. The only putative serine biosynthesis pathway is via the
265 transamination of pyruvate. Genes indicative of autotrophic metabolism (for
266 example, RuBisCO, ATP citrate lyase) were not identified. Additionally, genes
267 indicative of methanotrophy (*pmo*), methylotrophy (*mxaF* or *xoxF*), ammonia
268 (*amo*) or nitrite oxidation (*nxr*) were absent.

269

270  Genes encoding for the biosynthesis of all branched-chain (isoleucine, leucine
271  and valine) and aromatic (tryptophan, tyrosine and phenylalanine) amino acids
272  were conspicuously absent from the *Ca*. U copiosus genome. Additionally, the
273  biosynthetic pathways for arginine and histidine were also absent. These eight
274  amino acids are among the most energetically expensive to make (Fig S4, 31),
275  suggesting that their acquisition from the environment offers *Ca*. U copiosus an
276  energetic savings relative to taxa that synthesize them *de novo*. In contrast, Ca.
277  U. copiosus does have the complete suite of genes for the biosynthesis of those
278  amino acids that are energetically less expensive to make (including alanine,
279  aspartate, and glutamate, Fig. S4). The absence of branched-chain amino acid
280  synthesis pathways in the *Ca*. U. copiosus genome is consistent with previous
281  observations that these genes are underrepresented in natural populations of
282  soil *Verrucomicrobia* (10). Although *Ca*. U. copiosus lacks many amino acid
283  synthesis pathways, numerous genes encoding for peptide transport,
284  degradation and recycling were identified. For example, when scaled for
285  genome size, *Ca*. U. copiosus encodes four times as many putative peptide and
286  amino acid transporters as *C. flavus* (1.5% of genome to 0.37%) and twice as
287  many predicted proteases (6.5% of genome versus 3.2%). *Ca*. U. copiosus also
288  encodes for all components of the bacterial proteasome. Proteasomal
289  degradation is critical for amino acid recycling under starvation conditions in
290  mycobacteria (32). The enrichment of peptide transport and degradation
291  systems in the *Ca*. U. copiosus genome suggest that at least some of the amino
292  acids *Ca*. U. copiosus cannot synthesize are available directly from the soil
293  environment or by indirect associations with other soil biota.
294  In addition to these likely amino acid auxotrophies, *Ca*. U. copiosus has several
295  putative B-vitamin or B-vitamin precursor requirements. For example, the entire
296  vitamin $B_{12}$ synthesis pathway was absent in *Ca*. U. copiosus, despite the
297  presence of three genes encoding vitamin $B_{12}$-dependent proteins (methionine
298  synthase, ribonucleotide reductase, and methylmalonyl-CoA mutase). Vitamin
299  $B_{12}$ auxotrophies are relatively common in soil (33), making it likely that
300  exogenous vitamin $B_{12}$ is generally available to many soil bacteria. Genes
301  encoding for the complete vitamin $B_1$ biosynthetic pathway was complete
302  except for the 4-amino-3-hydroxymethyl-2 methylpyrmidine (HMP) synthase
303  (*thiC*), encoding the first step in $B_1$ biosynthesis. The absence of this gene, in the
304  presence of the remainder of the $B_1$ synthesis pathway, was recently linked to
305  an obligate HMP requirement in marine bacteria (34).
306

307  The abundance and cosmopolitan distribution of *Ca*. U. copiosus (Fig. 1),
308  together with its small genome size relative to other soil microbes (Fig. 3),
309  suggest that it is undergoing streamlining selection to minimize genome size.
310  The genome-streamlining hypothesis proposes that, in large bacterial
311  populations, reduced genome complexity is a trait under natural selection,
312  especially in environments where nutrients are sparse and can periodically limit
313  growth (25). All contemporary free-living organisms with streamlined genomes

314   inhabit aquatic environments (25, 35). However, compared to these aquatic
315   environments, soils are more heterogeneous (36), have higher overall microbial
316   diversity (37), and slower carbon turnover (38). Therefore, the functional
317   complexity required by soil microbes to succeed within a given niche is likely
318   large relative to that required by aquatic microbes. This means that the effects
319   of genome streamlining are likely to be most evident (i.e., result in smaller
320   genomes) in aquatic environments and that we might expect genome reduction
321   to be relatively uncommon across soil taxa. This expectation is consistent with
322   the fact that, on average, the genomes of aquatic microbes are smaller than
323   their terrestrial counterparts (25). However, the small genome and numerous
324   putative auxotrophies of *Ca.* U. copiosus show that genome streamlining is not
325   unique to aquatic organisms and that genome streamlining may also confer a
326   selective growth advantage in the soil environment.
327
328   Genome streamlining in *Ca*. U. copiosus has resulted in reduced catabolic and
329   biosynthetic capacity, and thus a loss of metabolic versatility. The absence of
330   multiple costly amino acid and vitamin biosynthetic pathways from the *Ca*. U.
331   copiosus genome implies that these compounds can be acquired from the soil
332   environment. Several studies have shown that free amino acids are present in
333   soil (39, 40), although oligopeptides are reported to be more abundant in
334   grasslands and may be assimilated with kinetics similar to free amino acids (41).
335   The enrichment of proteases and amino acid and peptide importers in the *Ca.* U.
336   copiosus genome suggests that it is well equipped to assimilate this fraction of
337   soil organic matter. Dispensing the capacity to synthesize costly amino acids
338   and vitamins likely provides *Ca*. U copiosus a growth advantage in resource
339   limiting conditions when competition for labile carbon is high. Alternatively,
340   many of the putative amino acid auxotrophies described here are involved in
341   synergistic growth (42) and may be supplied by other microbes as common
342   community goods (43). Based on the few spartobacterial isolates that have been
343   cultivated (11), culture-independent studies (10, 44), and the genomic data
344   presented here, we speculate that *Ca*. U. copiosus is a small, oligotrophic soil
345   bacterium that reduces its requirement for soil organic carbon by acquiring
346   costly amino acids and vitamins from the environment.
347
348   ***Conclusions***
349   Whereas successful soil microbes are predicted to have large genomes (27, 28,
350   Fig. 3), *Ca*. U. copiosus has a small genome, indicating that, similar to aquatic
351   microbes, minimization of cellular architecture can also represent a successful
352   strategy for soil microbes. We do not know if other uncultivated abundant soil
353   taxa also contain streamlined genomes because pre-existing genome databases
354   are preferentially biased towards cultivated isolates. For example, only 4.5% of
355   bacterial genomes in IMG are from uncultivated taxa (accessed April 2016).
356   Bacteria encoding for greater metabolic versatility likely have larger genomes
357   and therefore may be easier to culture in the laboratory (29). On the other hand,
358   specific and combinatorial nutrient requirements such as those described for

359   *Ca*. U. copiosus, present a complex problem for researchers attempting to
360   cultivate microbes with reduced genomes (45). Although *Ca.* U. copiosus has
361   not yet been grown in the laboratory, cultivation is clearly a crucial next step to
362   describing this organism, using the information described here to 'tailor' a
363   growth medium specifically for *Ca*. U. copiosus and related microbes. Such an
364   approach will undoubtedly improve our ability to describe and study the majority
365   of soil microbes, even dominant soil microbes like *Ca.* U. copiosus, which
366   remain difficult to cultivate under laboratory conditions.
367
368   **Materials and Methods**
369
370   ***Estimating the abundances and distributions of Verrucomicrobia in soil***
371   While five abundant *Verrucomicrobia* phylotypes were described in Fierer et al.
372   (2013), a single phylotype with 99% identity to the clone Da101 (15) was clearly
373   dominant. We searched previously published soil datasets for representative
374   sequences with 100% identity to this Da101 phylotype, including 31 soils from
375   United States native tallgrass prairies (10), 64 soils from matched forest and
376   grassland sites across North America (46), 595 soils collected from Central Park
377   in New York City (1), 367 grassland soils collected from North America, Europe,
378   Australia, and Africa (17), and a cross-biome collection of 11 desert and non-
379   desert soils from across the globe (18).  We also included a dataset from a
380   grassland terrace near Boulder, Colorado (105.23W, 40.12N, Table Mountain)
381   where 30 soils were collected from a depth of 25 cm within a $100m^2$ area on
382   28Jan15. Amplicon sequences and associated metadata from this study are
383   available at https://dx.doi.org/10.6084/m9.figshare.3363505.v3. Collectively
384   these datasets represent 1097 unique soil samples collected from a wide range
385   of ecosystem and soil types.
386
387   For all samples, DNA was extracted with the MoBio PowerSoil kit and the V4
388   region of the 16S rRNA gene was amplified in triplicate with the 515f/806r primer
389   pair. After normalization to equimolar concentrations, amplicons were
390   sequenced on an Illumina MiSeq (151bp paired end) at the University of
391   Colorado BioFrontiers Institute Next-Gen Sequencing Core Facility. Sequences
392   were processed as described previously (17). In brief, we used a combination of
393   QIIME (47) and UPARSE (48) to quality-filter, remove singletons, and merge
394   paired reads. Sequences were assembled into phylotypes at the 97% identity
395   level using UCLUST (49). Taxonomy was assigned using the Greengenes 13_8
396   database (50) and the Ribosomal Database Project classifier (4) and each
397   dataset was rarefied independently (Table S1).
398   As PCR primer biases can misestimate the relative abundances of
399   *Verrucomicrobia* (3,16), we also estimated the abundances of the Da101
400   phylotype directly from shotgun metagenomic data. We used Metaxa2 with
401   default settings (51) to extract bacterial 16S rRNA gene sequences from shotgun
402   metagenomic data compiled from previous analyses of 75 different soils (data
403   from 17, 18). Extracted 16S rRNA gene fragments were matched to GreenGenes

404    full-length sequences at 99% ID using the usearch7 command usearch_global.
405    The matched Greengenes sequences were then clustered and assigned
406    taxonomy as described above.
407

408    ***Describing the phylogenetic diversity of soil Verrucomicrobia***
409    We reconstructed near-full length 16S rRNA gene sequences to construct a
410    phylogeny of soil *Verrucomicrobia* from six soil samples (see Table S2) that were
411    selected to represent geographically distinct grasslands with a range of
412    verrucomicrobial abundances. We extracted DNA from each of these soils as
413    described previously (17) and used the 27f/1392r primer pair (52) to amplify near
414    full-length 16S rRNA genes as described in (19). The amplicons were sheared
415    using the Covaris M220 (Covaris, Woburn, MA) and the 16S rRNA gene libraries
416    were prepared using TruSeq DNA LT library preparation kits (Illumina, San
417    Diego, CA). Samples were pooled and sequenced on an Illumina MiSeq
418    (2x300bp) at the University of Colorado Next Generation Sequencing Facility.
419

420    After quality filtering of sequences, near full length SSU sequences were
421    reconstructed using EMIRGE (19). After 40 iterations, sequences were merged
422    into phylotypes with ≥97% similarity. Reconstructed sequences were trimmed
423    to 1200 bp and all sequences were further clustered at 95% identity due to gaps
424    in some assemblies. Full-length 16S rRNA sequences from named
425    verrucomicrobial isolates were aligned along with the reconstructed sequences
426    using PyNAST (53). A UPGMA tree was constructed using the R packages
427    seqnir and phangorn and visualized with GraPhlAn (R 3.2.2, version 0.9.7).
428

429    ***Assembly and annotation of a genome from the dominant soil***
430    ***Verrucomicrobial phylotype***
431    We assembled the genome of '*Candidatus* Udaeobacter copiosus' from a
432    metagenome of a U.S. prairie soil sample (NTP21, Hayden, IA) estimated to have
433    particularly high abundances of bacteria within the Da101 clade (10).
434    Fragmented DNA extracted from this soil was prepared for sequencing using
435    WaferGen's PrepX ILM DNA library Kit (WaferGen Biosystems Inc, Fremont, CA)
436    and the Apollo 324 Automated Library Prep System for library generation. The
437    library was sequenced on one Illumina HiSeq2000 lane (2×101 bp), yielding 17
438    Gb of sequence with an average paired-end insert size of 345 bp. Low quality
439    reads were trimmed using Sickle v. 1.29 with a quality score threshold of Q=3,
440    or removed if trimmed to <80 bp long (https://github.com/najoshi/sickle). The
441    sequences were assembled using IDBA_ud v. 1.1.0 (54) with a kmer range of 40
442    to 70 and step size of 15. To improve recovery of the most abundant
443    *Verrucomicrobia*, the genome was selectively re-assembled using Velvet with a
444    kmer size of 59, and expected kmer coverage of 11.5 (range 7.5 to 15.5). To bin
445    contigs ≥2 kb long, genes and protein sequences were predicted using Prodigal
446    v. 2.60 in metagenomics mode (55). For each contig, we determined the GC
447    content, coverage, and the phylogenetic affiliation based on the best hit for each
448    predicted protein in the Uniref90 database (Sept-2013, 56) following ublast

449 searches (49). We also constructed emergent self-organizing maps (ESOM)
450 using tetranucleotide frequencies of 5 kb DNA fragments (57). A combination of
451 these approaches were used to identify the genome. The draft genome was
452 uploaded to IMG for annotation under the taxon ID 2651869889.
453
454 No rRNA genes were annotated by IMG, so we used Metaxa with default
455 settings on the unassembled sequences to extract any 16S rRNA genes. Metaxa
456 recovered two ~500 bp 16S rRNA gene fragments at 23-29× coverage which
457 aligned to separate regions of the full-length 16S rRNA gene from the closest
458 related verrucomicrobial genome (*C. flavus*). Because these two rRNA gene
459 fragments have the same coverage as the genome and align to separate regions
460 of one 16S rRNA gene, it is likely that the sequences encode a single rRNA
461 operon.
462

467
468
469 ***References***
470

471 1. Ramirez KS et al. (2014) Biogeographic patterns in below-ground diversity in
472 New York City's Central Park are similar to those observed globally. *Proc R Soc*
473 *B* 281(1795):20141988.
474
475 2. Janssen PH. (2006) Identifying the dominant soil bacterial taxa in libraries of
476 16S rRNA and 16S rRNA genes. *Appl Environ Microbiol* 72(3):1719–28.
477
478 3. Bergmann GT et al. (2011) The under-recognized dominance of
479 Verrucomicrobia in soil bacterial communities. *Soil Biol Biochem* 43(7):1450-5.
480
481 4. Wang, Q, Garrity GM, Tiedje JM, and Cole JR. (2007) Naïve Bayesian
482 Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial
483 Taxonomy. *Appl Environ Microbiol* 73(16):5261-7.
484
485 5. Markowitz VM et al. (2014) IMG 4 version of the Integrated Microbial
486 Genomes comparative analysis system. *Nucleic Acids Res* 42:D560–7.
487
488 6. VanInsberghe, D et al. (2015) Non-symbiotic Bradyrhizobium ecotypes
489 dominate North American forest soils. *ISME J* 9(11):2435-41.
490
491 7. Fierer, N., M.S. Strickland, D. Liptzin, M.A. Bradford, C.C. Cleveland. (2009)
492 Global patterns in belowground communities. *Ecol Lett* 12(11):1238-49.
493

494  8. Dunfield PF et al. (2007) Methane oxidation by an extremely acidophilic
495  bacterium of the phylum Verrucomicrobia. *Nature* 450(7171):879–82.
496
497  9. Everard A et al. (2013) Cross-talk between *Akkermansia muciniphila* and
498  intestinal epithelium controls diet-induced obesity. *Proc Natl Acad Sci USA*
499  110(22):9066-71
500
501  10. Fierer N et al. (2013) Reconstructing the microbial diversity and function of
502  pre-agricultural tallgrass prairie soils in the United States. *Science* 342(6158):
503  621-4.
504
505  11. Sangwan P, Chen X, Hugenholtz P and Janssen PH. (2004) *Chthoniobacter*
506  *flavus* gen. nov., sp. nov., the First Pure-Culture Representative of Subdivision
507  Two, *Spartobacteria* classis nov., of the Phylum *Verrucomicrobia*. *Appl Environ*
508  *Microbiol* 70(10):5875–81.
509
510  12. Kant R, van Passel MWJ, Palva A, et al. (2011) Genome Sequence of
511  Chthoniobacter flavus Ellin428, an Aerobic Heterotrophic Soil Bacterium. *J*
512  *Bacteriol* 193(11):2902-3.
513
514  13. Herlemann DPR et al. (2013) Metagenomic De Novo Assembly of an Aquatic
515  Representative of the Verrucomicrobial Class Spartobacteria. *mBio* 4(3):e00569-
516  12
517
518  14. Vandekerckhove TT, Willems A, Gillis M, Coomans A. (2000) Occurrence of
519  novel verrucomicrobial species, endosymbiotic and associated with
520  parthenogenesis in Xiphinema americanum-group species (Nematoda,
521  Longidoridae). *Int J Syst Evol Microbiol* 50:2197–205.
522
523  15. Felske A, Akkermans ADL. (1998) Prominent occurrence of ribosomes from
524  an uncultured bacterium of the Verrucomicrobiales cluster in grassland
525  soils. *Lett Appl Microbiol* 26(3):219–23.
526
527  16. Guo J, Cole JR, Zhang Q, Brown CT, Tiedje JM. (2016) Microbial community
528  analysis with ribosomal gene fragments from shotgun metagenomes. *Appl*
529  *Environ Microbiol* 82(1):157–66.
530
531  17. Leff, JW et al. (2015) Consistent responses of soil microbial communities to
532  elevated nutrient inputs in grasslands across the globe. *Proc Natl Acad Sci USA*
533  112(35):10967-72.
534
535  18. Fierer N et al. (2012) Cross-biome metagenomic analyses of soil microbial
536  communities and their functional attributes. *Proc Natl Acad Sci USA*
537  109(52):21390–5.
538

539 19. Miller CS, Baker BJ, Thomas BC, Singer SW, Banfield JF. (2011) EMIRGE:
540 reconstruction of full-length ribosomal genes from microbial community short
541 read sequencing data. *Genome Biol* 12:R44.
542
543 20. Howe AC et al. (2014) Tackling soil diversity with the assembly of large,
544 complex metagenomes. *Proc Natl Acad Sci USA* 111(13):4904–9.
545
546 21. Finn RD et al. (2016) The Pfam protein families database: towards a more
547 sustainable future. *Nucleic Acids Res* 44(D1):D279-85.
548
549 22. Ciccarelli FD et al. (2006) Toward Automatic Reconstruction of a Highly
550 Resolved Tree of Life. *Science* 311(5765):1283-7.
551
552 23. Varghese NJ et al. (2015) Microbial species delineation using whole genome
553 sequences. *Nucleic Acids Res* 10.1093.
554
555 24. Raes J, Korbel JO, Lercher MJ, von Mering C, Bork P. (2007). Prediction of
556 effective genome size in metagenomic samples. *Genome Biol* 8(1):R10.
557
558 25. Giovannoni SJ, Thrash JC, Temperton B. (2014) Implications of streamlining
559 theory for microbial ecology. *ISME J* 8:1553-65.
560
561 26. Nayfach S, Pollard KS. (2015) Average genome size estimation improves
562 comparative metagenomics and sheds light on the functional ecology of the
563 human microbiome. *Genome Biol* 16(1):51.
564
565 27. Konstantinidis KT, Tiedje JM. (2003) Trends between gene content and
566 genome size in prokaryotic species with larger genomes. *Proc Natl Acad Sci
567 USA* 101(9):3160-5.
568
569 28. Barberán A et al. (2014) Why are some microbes more ubiquitous than
570 others? Predicting the habitat breadth of soil bacteria. *Ecol Lett* 17(7):794–802.
571
572 29. Button DK, Robertson, BR. (2001) Determination of DNA content of aquatic
573 bacteria by flow cytometry. *Appl Environ Microbiol* 67(4):1636-45.
574
575 30. Khadem AF et al. (2012) Genomic and Physiological Analysis of Carbon
576 Storage in the Verrucomicrobial Methanotroph "*Ca*. Methylacidiphilum
577 Fumariolicum" SolV. *Front Microbiol* 3:345.
578
579 31. Akashi H, and Gojobori T (2002) Metabolic efficiency and amino acid
580 composition in the proteomes of *Escherichia coli* and *Bacillus subtilis*. *Proc Natl
581 Acad Sci USA* 99(6):3695-3700.
582

13

583  32. Elharar Y et al. (2014) Survival of mycobacteria depends on proteasome-
584  mediated amino acid recycling under nutrient limitation. *EMBO J* 33(16):1802-
585  14.
586
587  33. Lochhead, AG. (1958) Soil bacteria and growth promoting substances.
588  *Bacteriol Rev* 22(3):145-53.
589
590  34. Carini P et al. (2014) Discovery of a SAR11 growth requirement for thiamin's
591  pyrimidine precursor and its distribution in the Sargasso Sea. *ISME J* 8:1727–38.
592
593  35. Kantor RS et al. (2013) Small Genomes and Sparse Metabolisms of
594  Sediment-Associated Bacteria from Four Candidate Phyla. *mBio* 4(5):e00708-
595  13.
596
597  36. Vos M, Wolf AB, Jennings SJ, Kowalchuk GA. (2013) Micro-scale
598  determinants of bacterial diversity in soil. *FEMS Microbiol Rev* 37(6):936–54.
599
600  37. Fierer N, Lennon JT. (2011) The generation and maintenance of diversity in
601  microbial communities. *Am J Bot* 98(3):439–48.
602
603  38. Giovannoni SJ, Vergin VL. (2012) Seasonality in Ocean Microbial
604  Communities. *Science* 335(6069):671-6.
605
606  39. Friedel JK, Scheller E. (2002) Composition of hydrolysable amino acids in
607  soil organic matter and soil microbial biomass. *Soil Bio Biochem* 34(3):315–25.
608
609  40. Sauheitl L, Glaser B, Dippold M, Lieber K, Weigelt A (2010) Amino acid
610  fingerprint of a grassland soil reflect changes in plant species richness. *Plant
611  Soil* 334(1):353–63.
612
613  41. Farrell M et al. (2013) Oligopeptides Represent a Preferred Source of
614  Organic N Uptake: A Global Phenomenon? *Ecosystems* 16(1):133-45.
615
616  42. Mee MT, Collins JJ, Church GM, Wang HH (2014) Syntrophic exchange in
617  synthetic microbial communities. *Proc Natl Acad Sci USA* 111(20):E2149–
618  E2156.
619
620  43. Morris JJ, Lenski RE, Zinser ER (2012) The Black Queen Hypothesis:
621  Evolution of dependencies through adaptive gene loss. *mBio* 3(2):e00036-12.
622
623  44. Portillo, M.C., J.W. Leff, C.L. Lauber, N. Fierer. (2013) Cell size distributions
624  of soil bacterial and archaeal taxa. *Appl Environ Microbiol* 79(24):7610-7617.
625

14

626   45. Carini P, et al. (2013) Nutrient requirments for growth of the extreme
627   oligotroph 'Candidatus Pelagibacter ubique' HTCC1062 on a defined medium.
628   ISME J 7(3):592-602.
629

630   46. Crowther, TW et al. (2014) Predicting the responsiveness of soil biodiversity
631   to deforestation: a cross-biome study. Glob Chang Biol 20(9):2983-94.
632

633   47. Caporaso JG et al. (2010) QIIME allows analysis of high-throughput
634   community sequencing data. Nat Methods 7(5):335–336.
635

636   48. Edgar RC (2013) UPARSE: highly accurate phylotype sequences from
637   microbial amplicon reads. Nat Methods 10(10):996-8
638

639   49. Edgar, RC (2010) Search and clustering orders of magnitude faster than
640   BLAST. Bioinformatics 26(19):2460-1.
641

642   50. DeSantis, T.Z et al. (2006) Greengenes, a Chimera-Checked 16S rRNA Gene
643   Database and Workbench Compatible with ARB. Appl Environ Microbiol
644   72(7):5069-72.
645

646   51. Bengtsson-Palme J et al. (2015) Metaxa2: improved identification and
647   taxonomic classification of small and large subunit rRNA in metagenomic data.
648   Mol Ecol Resour 15(6):1403-14.
649

650   52. Miller CS et al. (2013) Short-Read Assembly of Full-Length 16S Amplicons
651   Reveals Bacterial Diversity in Subsurface Sediments. PLoS ONE 8(2):e56018.
652

653   53. Caporaso JG et al. (2010) PyNAST: a flexible tool for aligning sequences to a
654   template alignment. Bioinformatics 26(2):266-7
655

656   54. Peng Y, Leung HCM, Yiu SM, Chin FYL. (2012) IDBA-UD: a de novo
657   assembler for single-cell and metagenomic sequencing data with highly uneven
658   depth. Bioinformatics 28(11):1420–8.
659

660   55. Hyatt D, et al. (2010) Prodigal: prokaryotic gene recognition and translation
661   initiation site identification. BMC Bioinformatics 11:119.
662

663   56. Suzek BE, Huang H, McGarvey P, Mazumder R, Wu CH (2007) UniRef:
664   comprehensive and non-redundant UniProt reference clusters. Bioinformatics
665   23(10):1282–8.
666

667   57. Dick GJ, et al. (2009) Community-wide analysis of microbial genome
668   sequence signatures. Genome Biol 10(8):R85.
669

670

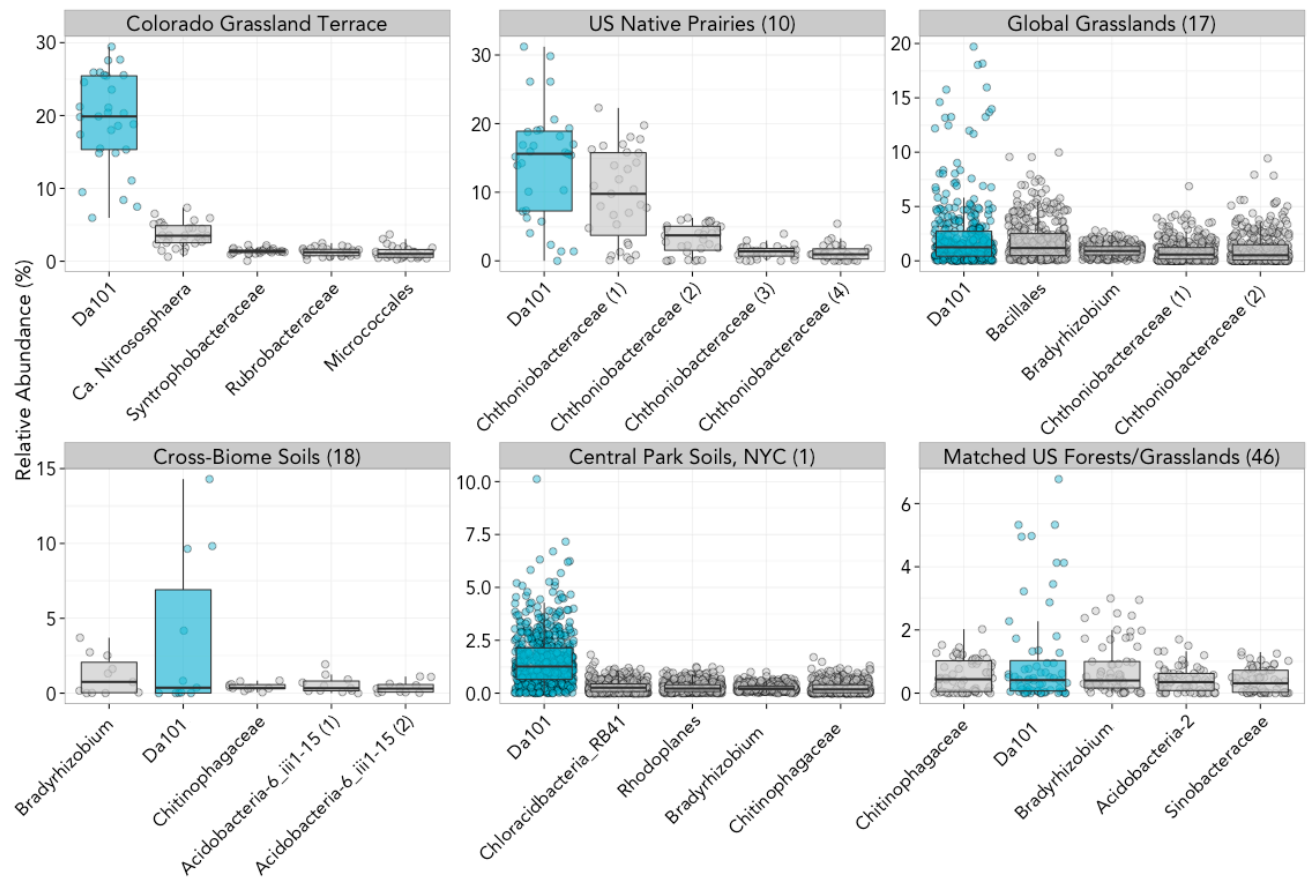671 ***Figures***

672



673
674 Fig. 1. Da101 is, on average, one of the most abundant bacterial phylotypes
675 found across >1000 soils collected from a wide range of soil and ecosystem
676 types across the globe. The Da101 phylotype is indicated in blue while other
677 abundant taxa are indicated in grey. Taxa are listed on the x-axis in order of
678 their rank abundance (taxa on the left are the most abundant). Data comes from
679 (10) Fierer et al. 2013, (17) Leff et al. 2015, (18) Fierer et al. 2012, (1) Ramirez et
680 al. 2014, and (46) Crowther et al. 2014. Further details on each of these studies
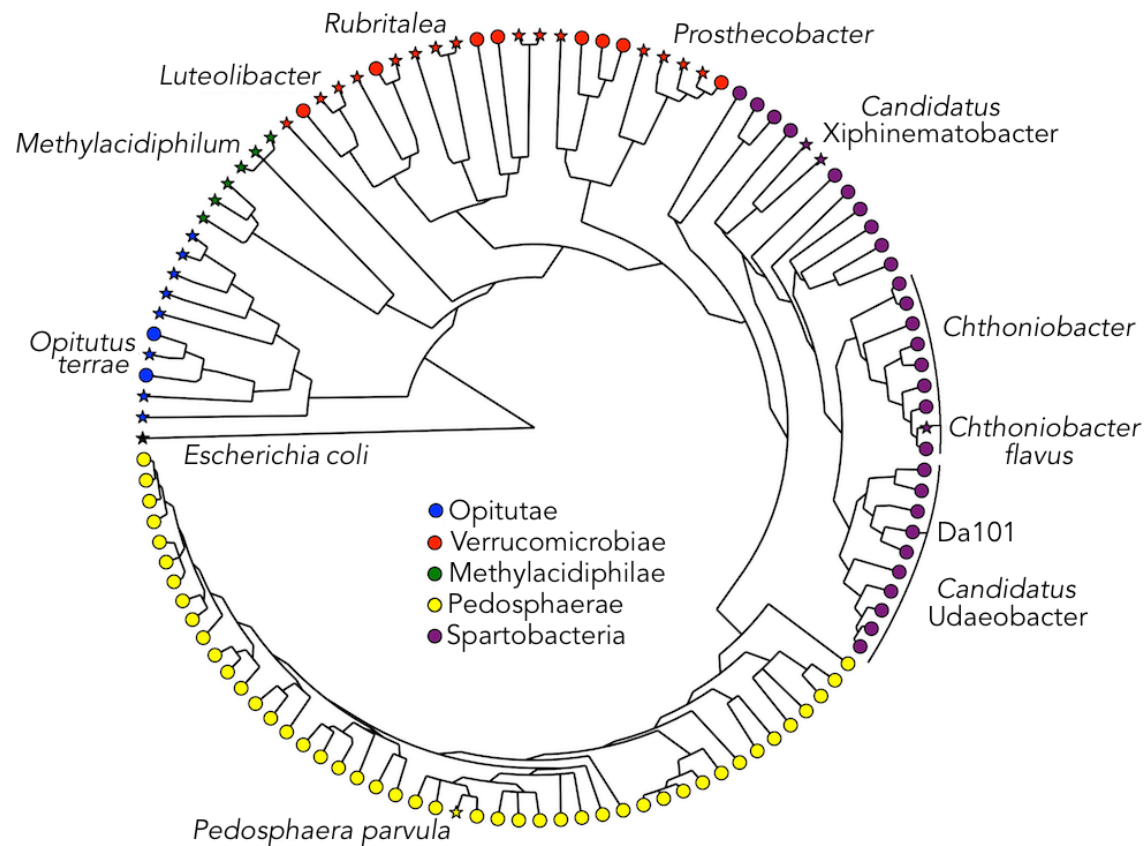681 are provided in Table S1.
682

Fig. 2. Phylogenetic analyses of soil Verrucomicrobia. Stars denote 16S rRNA gene sequences of named isolates while circles represent environmental 16S rRNA gene sequences assembled from 6 soils using EMIRGE (Table S2). The uncultivated verrucomicrobial phylotype Da101 falls within a cluster distinct from cultivated *Spartobacteria*. The UPGMA phylogenetic tree was constructed using 1200bp verrucomicrobial 16S rRNA gene sequences and is rooted with a 16S rRNA gene sequence from *Escherichia coli* K-12. Notable verrucomicrobial isolates and genera are labeled. Colors indicate verrucomicrobial classes.
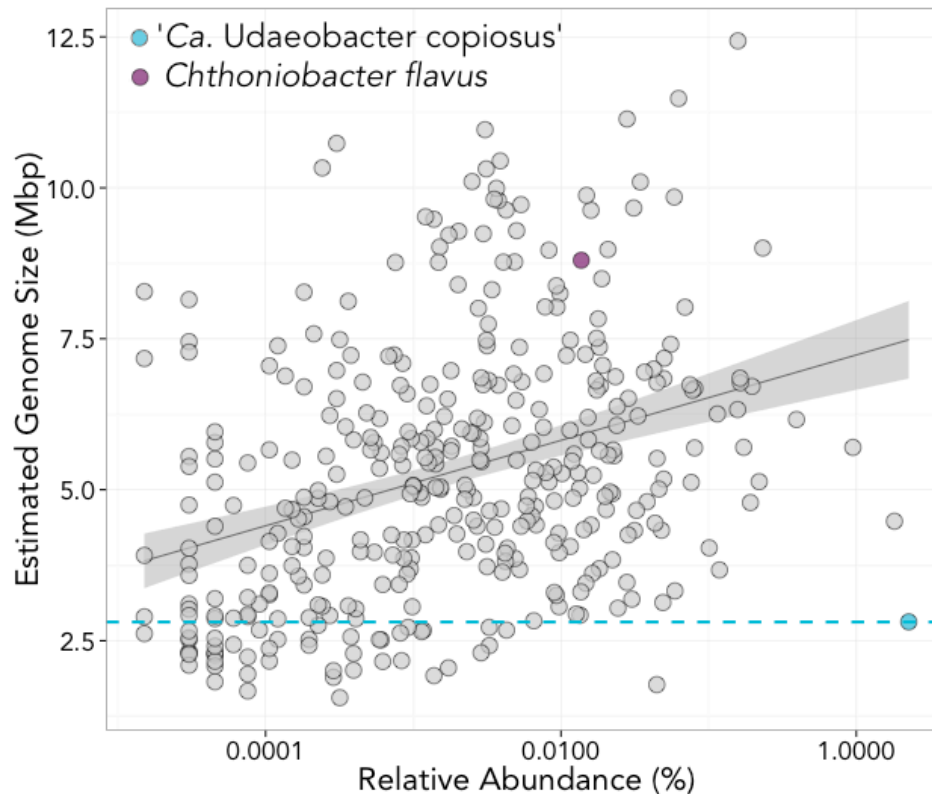
Fig. 3. '*Ca.* Udaeobacter copiosus' has a reduced genome size compared to other abundant soil bacteria. Points represent the estimated genome size and relative abundances of 378 bacterial genomes obtained by matching 16S rRNA gene sequences from (17) to 16S rRNA gene sequences extracted from genomes in IMG at 99% sequence identity. Only genomes classified as 'permanent draft' or 'finished' status were used. Bacteria with larger genomes tend to be more abundant (p<0.0001, ρ=0.368, Spearman correlation) with *Ca*. U. copiosus (indicated in blue) being a notable exception to this pattern, as it has a high relative abundance (2.26% of 16S rRNA sequences) but has a relatively small genome.
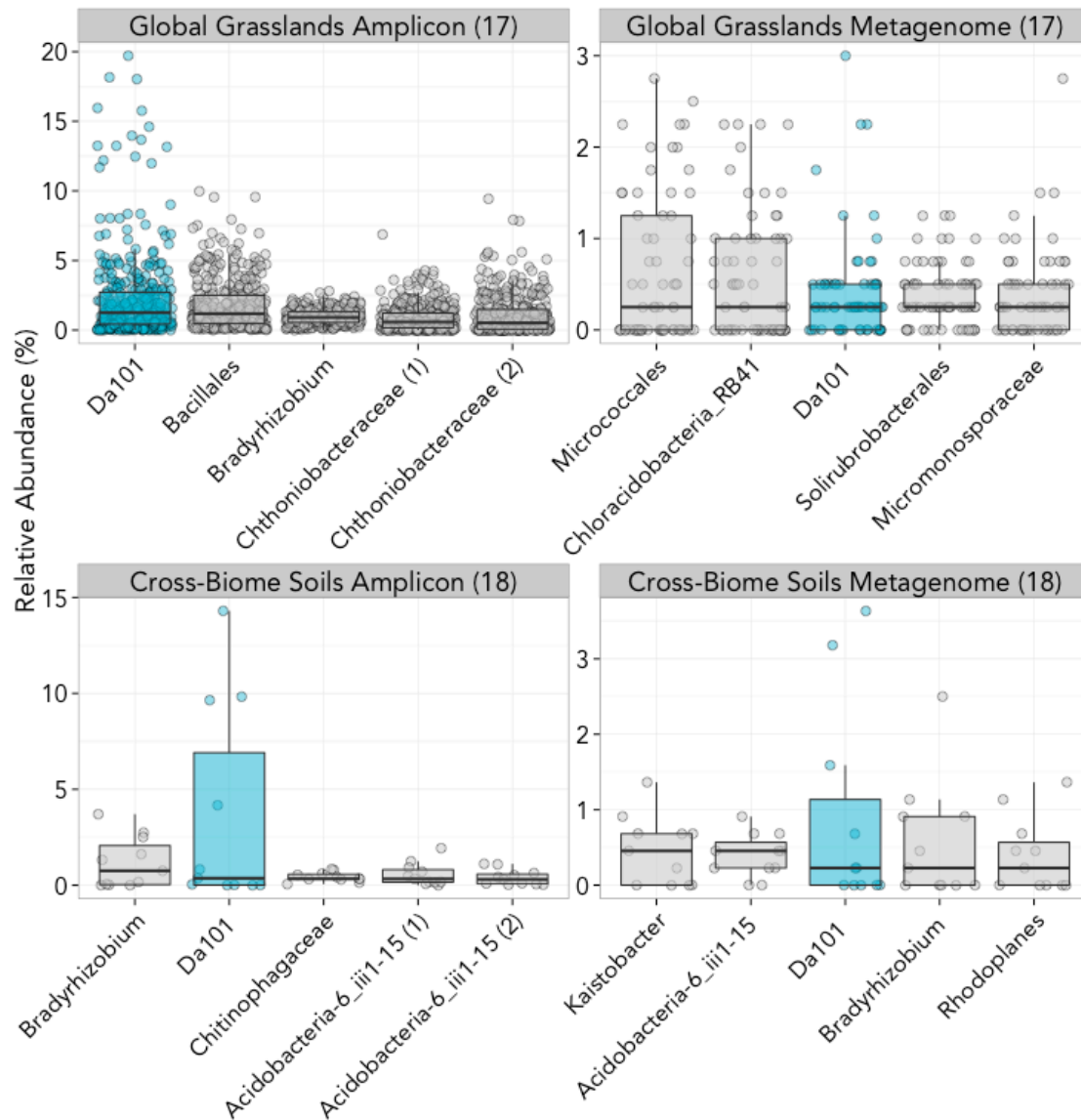
18

Fig. S1: Da101 rank is similar in amplicon and metagenomic data. The top 5 phylotypes from two matched amplicon and metagenomic datasets (Leff et al. 2015 (17), Fierer et al. 2012 (18)) are shown in order of decreasing median rank. Each point represents one sample within the corresponding dataset. Da101's position is highlighted with blue while all other phylotypes are grey.
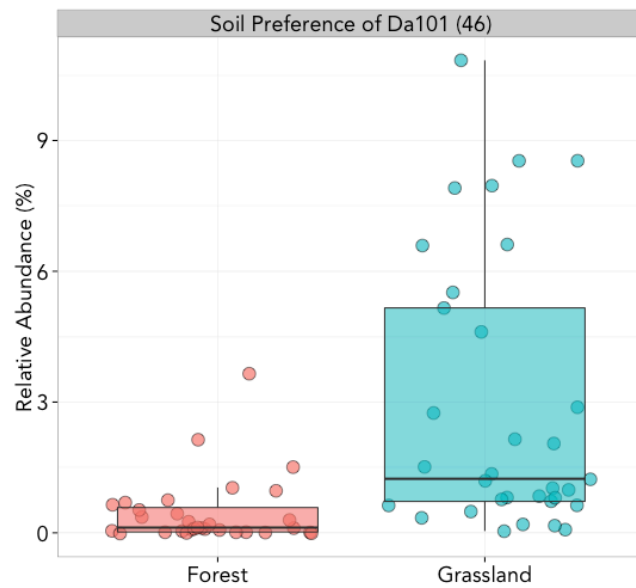
712
713  Fig. S2. Phylotype Da101 is more abundant in grasslands than forests;
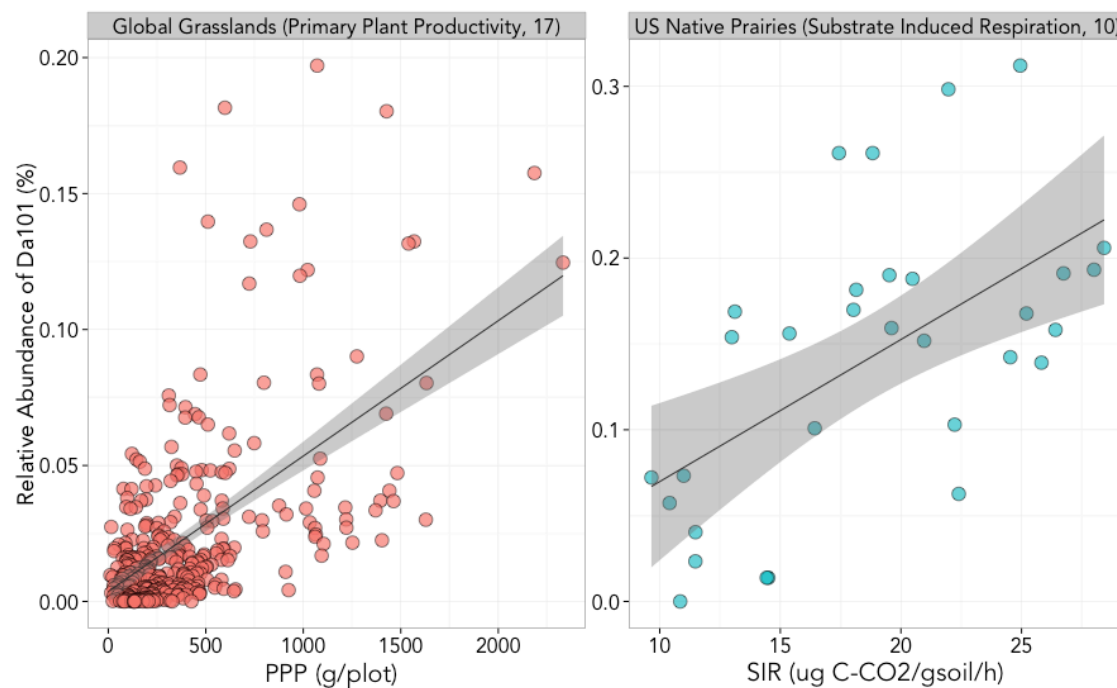714  p<0.0001, two tailed t-test. Data is from Crowther et al. 2014 (46).
715
716



717
718  Fig. S3. The abundance of phylotype Da101 correlates with measures of
719  microbial and plant biomass in two separate studies: SIR p<0.001 ρ=0.57, PPP
720  p<0.0001 ρ=0.47, Spearman correlations. Data is from Fierer et al. 2013 (10) and
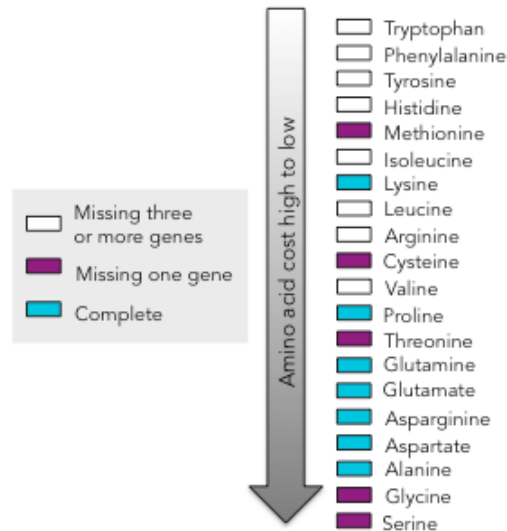721  Leff et al. 2015 (17).
722

20

723
724
725 Fig. S4. '*Ca.* Udaeobacter copiosus' lacks pathways to synthesize several
726 expensive amino acids, yet has partial or complete pathways for relatively
727 inexpensive amino acids. Cost of each of amino acid was estimated in *E. coli* by
728 number of high-energy phosphate bonds hydrolyzed (31).
729
730
731

**Table S1: Descriptions of Datasets Used in this Study**

|  | Fierer et al. 2012 | Fierer et al. 2013 | Crowther et al. 2014 | Ramirez et al. 2014 | Leff et al. 2015 | Table Mountain, CO |
|---|---|---|---|---|---|---|
| **REGION DESCRIPTION** | Global cross-biome | U.S. native prairies | U.S. matched forest/grasslands | Central park, NYC | Global grasslands | CO Terrace |
| **# SAMPLES** | 11 | 31 | 64 | 595 | 367 | 29 |
| **# SITES** | 11 | 31 | 11 | 1 | 25 | 1 |
| **METAGENOMIC DATA** | 11 | - | - | - | 87 | - |
| **ITS DATA** | × | × | ✓ | ✓ | ✓ | × |
| **18S DATA** | × | × | × | ✓ | × | × |
| **pH RANGE** | 4.1-9.5 | 5.8-7.9 | 4.0-8.1 | 3.9-8.4 | 4.4-8.2 | - |
| **MAT RANGE (°C)** | -19-25 | 3.7-18.9 | -3.2-22.8 | 13 | 0.3-18.4 | 11 |
| **MAP RANGE (mm)** | 100-4000 | 503-1148 | 287-3460 | 1016 | 262-1898 | 525 |
| **RAREFACTION DEPTH** | 15000 | 942 | 4000 | 5000 | 18000 | 11000 |

21

**Table S2: Samples used to construct EMIRGE phylogeny**

| SAMPLE | DATASET | LOCATION | LATITUDE | LONGITUDE | DESCRIPTION |
|--------|---------|----------|----------|-----------|-------------|
| NTP21 | Fierer et al. 2013 | Hayden, IA | 43.26 | -92.23 | Native prairie |
| NTP28 | Fierer et al. 2013 | Glynn Prairie, MN | 44.15 | -95.41 | Native prairie |
| NN1182 | Leff et al. 2015 | Val Mustair, Switzerland | 46.63 | 10.37 | Alpine grassland |
| NN772 | Leff et al. 2015 | Msunduzi Municipality, South Africa | -29.67 | 30.40 | Mesic grassland |
| TM25 | New data set | Table Mountain, CO | 40.01 | -105.5 | Grassland terrace |
| GG14 | New data set | Gordon Gulch, CO | 40.12 | -105.2 | Meadow |

732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762

**Table S3 'Candidatus Udaeobacter copiosus' genome has 34/36 single copy housekeeping genes**

| IMG GENE ID | COG/PFAM PRODUCT NAME | SEQUENCE LENGTH (BP) | COG/PFAM |
|---|---|---|---|
| 2653240560 | arginyl-tRNA synthetase | 1737 | COG0018 |
| 2653239845 | DNA-directed RNA polymerase subunit alpha | 1044 | COG0202 |
| 2653239819 | DNA-directed RNA polymerase subunit beta | 2910 | COG0085 |
| 2653240331 | histidyl-tRNA synthetase | 1239 | COG0124 |
| 2653239579 | Ribosome-binding ATPase GTP1/OBG family | 1125 | COG0012 |
| 2653241119 | isoleucyl-tRNA synthetase | 2748 | COG0060 |
| 2653239815 | large subunit ribosomal protein L1 | 702 | COG0081 |
| 2653239810 | large subunit ribosomal protein L11 | 429 | COG0080 |
| 2653239936 | large subunit ribosomal protein L13 | 438 | COG0102 |
| 2653241207 | large subunit ribosomal protein L14 | 366 | COG0093 |
| 2653241203 | large subunit ribosomal protein L16 | 426 | COG0197 |
| 2653241213 | large subunit ribosomal protein L18 | 375 | COG0256 |
| 2653241195 | large subunit ribosomal protein L3 | 726 | COG0087 |
| 2653241210 | large subunit ribosomal protein L5 | 573 | COG0094 |
| 2653241212 | large subunit ribosomal protein L6 | 540 | COG0097 |
| 2653239619 | leucyl-tRNA synthetase | 2412 | COG0495 |
| 2653242017 | N6-L-threonylcarbamoyladenine synthase | 1086 | COG0533 |
| 2653241883 | phenylalanyl-tRNA synthetase alpha chain | 282 | pfam01409 |
| 2653241216 | preprotein translocase subunit SecY | 1509 | COG0201 |
| 2653241215 | ribosomal protein L15 | 636 | COG0200 |
| 2653241201 | ribosomal protein L22 | 480 | COG0091 |
| 2653241222 | small subunit ribosomal protein S11 | 630 | COG0100 |
| 2653241191 | small subunit ribosomal protein S12 | 435 | COG0048 |
| 2653241221 | small subunit ribosomal protein S13 | 396 | COG0099 |
| 2653241789 | small subunit ribosomal protein S15 | 177 | pfam00312 |
| 2653241206 | small subunit ribosomal protein S17 | 297 | COG0186 |
| 2653239969 | small subunit ribosomal protein S2 | 702 | COG0052 |
| 2653241202 | small subunit ribosomal protein S3 | 690 | COG0092 |
| 2653241223 | small subunit ribosomal protein S4 | 612 | COG0522 |
| 2653241214 | small subunit ribosomal protein S5 | 564 | COG0098 |
| 2653241192 | small subunit ribosomal protein S7 | 474 | COG0049 |
| 2653241211 | small subunit ribosomal protein S8 | 399 | COG0096 |
| 2653239935 | small subunit ribosomal protein S9 | 396 | COG0103 |
| 2653241990 | valyl-tRNA synthetase | 816 | pfam00133 |