

1 **Title:**

2 **The ecological relevance of flagellar motility in soil bacterial**  
3 **communities**

4

5 **Authors:**

6 **Josep Ramoneda**, Cooperative Institute for Research in Environmental Sciences,  
7 University of Colorado, Boulder, Colorado, USA

8 **Kunkun Fan**, Laboratory of Soil and Sustainable Agriculture, Institute of Soil  
9 Science, Chinese Academy of Sciences, Nanjing, China

10 **Jane M. Lucas**, Cary Institute of Ecosystem Studies, Millbrook, New York, USA

11 **Haiyan Chu**, Laboratory of Soil and Sustainable Agriculture, Institute of Soil  
12 Science, Chinese Academy of Sciences, Nanjing, China

13 University of Chinese Academy of Sciences, Beijing, China

14 **Andrew Bissett**, CSIRO, Hobart, TAS7000, Australia

15 **Michael S. Strickland**, Department of Soil and Water Systems, University of Idaho,  
16 Moscow, Idaho, USA

17 **Noah Fierer**, Cooperative Institute for Research in Environmental Sciences,  
18 University of Colorado, Boulder, Colorado, USA

19 Department of Ecology and Evolutionary Biology, University of Colorado, Boulder,  
20 Colorado, USA

21

22 **Correspondence:**

23 Josep Ramoneda, [ramoneda.massague@gmail.com](mailto:ramoneda.massague@gmail.com)

24 Noah Fierer, [noah.fierer@colorado.edu](mailto:noah.fierer@colorado.edu)

25

26

## 27 **Abstract**

28 Flagellar motility is a key bacterial trait as it allows bacteria to navigate their  
29 immediate surroundings. Not all bacteria are capable of flagellar motility, and the  
30 distribution of this trait, its ecological associations, and the life history strategies of  
31 flagellated taxa remain poorly characterized. We developed and validated a genome-  
32 based approach to infer the potential for flagellar motility across 12 bacterial phyla  
33 (26,192 genomes in total). The capacity for flagellar motility was associated with a  
34 higher prevalence of genes for carbohydrate metabolism and higher maximum  
35 potential growth rates, suggesting that flagellar motility is more prevalent in resource-  
36 rich environments due to the energetic costs associated with this trait. To test this  
37 hypothesis, we focused on soil bacterial communities, where flagellar motility is  
38 expected to be particularly important given the heterogeneous nature of the soil  
39 environment. We applied a method to infer the prevalence of flagellar motility in  
40 whole bacterial communities from metagenomic data, and quantified the prevalence  
41 of flagellar motility across 4 independent field studies that each captured putative  
42 gradients in soil carbon availability (148 metagenomes). As expected, we observed a  
43 positive relationship between the prevalence of bacterial flagellar motility and soil  
44 carbon availability in each of these datasets. Given that soil carbon availability is  
45 often correlated with other factors that could influence the prevalence of flagellar  
46 motility, we validated these observations using metagenomic data acquired from a  
47 soil incubation experiment where carbon availability was directly manipulated with  
48 glucose amendments, confirming that the prevalence of bacterial flagellar motility is  
49 consistently associated with soil carbon availability over other potential confounding  
50 factors. Flagellar motility is a fundamental phenotypic trait for bacterial adaptation to  
51 soil, defining life history strategies primarily associated with resource availability.  
52 More generally, this work highlights the value of combining genomic and  
53 metagenomic approaches to expand our understanding of microbial phenotypic traits  
54 and reveal their general environmental associations.

55

## 56 Introduction

57 Microorganisms navigate their environment by responding to gradients in nutrients,  
58 toxins, and environmental conditions, in a process called chemotaxis [1, 2]. Flagellar  
59 motility is a widespread adaptation that allows bacteria to colonize new micro-  
60 environments by facilitating access to space and nutrients [3, 4], and enables escape  
61 from unfavorable conditions [5] and predators [6]. For example, moving towards  
62 environmental cues is an effective mechanism by which most pathogens [7, 8] and  
63 symbionts [9, 10] colonize their hosts. Despite the recognition that swimming and  
64 swarming (the two main modes of flagellar motility) are widely used to navigate  
65 microbial environments, empirical knowledge on the environmental conditions where  
66 bacterial flagellar motility can be beneficial remains rather limited, as most  
67 knowledge derives from laboratory-based studies using model organisms.

68 In laboratory conditions, bacteria have been widely investigated for their ability to  
69 swim towards resources [11, 12], display quorum sensing [13], or swim away from  
70 toxins [14]. Several experimental studies show that the hydration level of surfaces  
71 generally predicts how easily bacteria can colonize a given surface [15], and that  
72 flagellar motility also predicts the temporal persistence of bacterial pathogens in host  
73 microbiomes [16]. The high energetic cost of powering the flagellar machinery is  
74 tightly linked to regulatory systems that control flagellar expression depending on the  
75 spatial proximity and quality of available resources (i.e., optimal foraging based on  
76 energetic constraints; [17–19]). Notably, different flagellar systems have evolved in  
77 response to distinct environmental conditions, as exemplified by the case of the  
78 *Vibrio* genus, which use different flagellar systems depending on the spatial  
79 complexity of their surroundings [20]. This broad body of knowledge leads to the  
80 expectation that flagellar motility should display general ecological associations, but  
81 such patterns have not been comprehensively explored.

82 Research on bacterial flagellar motility predates modern microbiology, and  
83 laboratory-based approaches have enabled the discovery of the genes involved in  
84 flagellar assembly [21, 22]. The genes encoding for the flagellar machinery are  
85 reasonably well-known and generally conserved across a broad diversity of bacterial  
86 groups [23, 24]. Because the production of flagella requires a well-defined gene  
87 repertoire, the prediction of flagellar motility across taxa is likely feasible [25]. Yet,  
88 the proportion of bacterial taxa for which flagellar motility could theoretically be  
89 inferred from genomic information contrasts with the relatively limited number of  
90 strains for which flagellar expression has been empirically determined. A comparison  
91 of the number of strains with known motility information in bacterial phenotypic trait  
92 databases [26] to the total number of genomes contained in the Genome Taxonomy  
93 Database (GTDB r207, [27]) highlights that we have information on whether taxa are  
94 flagellated or not for only ~10% of the bacterial strains with available whole genome  
95 information. If we could infer the capacity for flagellar motility across a broad diversity

96 of microbial taxa, we could determine the set of traits that generally characterize  
97 flagellated taxa (so-called life history strategies; [28, 29]). Previous studies have  
98 linked flagellar motility to a fast growth (copiotroph) strategy [30], and flagellar  
99 motility is expected to be associated with a life history strategy for rapid nutrient  
100 acquisition [31]. However, one of the main challenges with identifying the life history  
101 strategies of bacteria remains the quantification of phenotypic traits. Thus,  
102 developing methods to infer flagellar motility across single bacterial genomes and  
103 metagenomes can help us identify the main ecological and life history associations  
104 of this important trait.

105 Flagellar motility is likely common for bacteria living in many environments –  
106 including host-associated and aquatic environments [2]. However, we are particularly  
107 interested in the prevalence of flagellar motility in soil environments because soil is a  
108 heterogeneous environment where resources are patchily distributed, and access to  
109 resources is a key factor structuring soil bacterial communities [32, 33]. We expect a  
110 high degree of variability in the prevalence of motile bacteria in soil as motility  
111 requires continuous water films [34, 35], and the high energetic cost associated with  
112 flagellar motility may be disadvantageous in the resource-limited conditions often  
113 common in soil [17, 36]. Since organic carbon compounds are likely the main  
114 sources of energy for soil bacteria, soil C availability is likely a key factor determining  
115 the selective advantage of flagellar motility in soil. Indeed, several studies have  
116 found a higher prevalence of bacterial flagellar motility in soil environments that  
117 generally have higher C availability. For example, plant rhizospheres usually contain  
118 elevated levels of available C compared to adjacent bulk soil environments due to  
119 plant-derived organic carbon inputs [37], and generally harbor a higher prevalence of  
120 flagellar genes [38, 39]. In arid environments, several studies have detected a  
121 negative relationship between the prevalence of flagellar motility genes and aridity  
122 [40, 41], which could be due to both lower C availability or to lower moisture. Given  
123 the spatial heterogeneity of soil, and the fitness advantage theoretically gained from  
124 flagellar motility in conditions where energy-rich resources are patchily distributed  
125 [18], we hypothesize that bacterial flagellar motility should exhibit a general positive  
126 relationship with soil C availability.

127 We had three objectives with this study. First, we wanted to build genome and  
128 metagenome-based models to accurately infer the potential for flagellar motility  
129 across bacterial taxa and whole bacterial communities. Second, we sought to identify  
130 the general life history strategies associated with flagellar motility in bacteria. Third,  
131 we aimed to determine the prevalence of flagellar motility in soil bacterial  
132 communities and to test the hypothesis that flagellar motility is more prevalent in  
133 soils with higher C availability. To this end, we estimated the potential for flagellar  
134 motility across 26,192 bacterial taxa with available genomic information based on a  
135 machine learning model trained on empirical information for this trait, and explored  
136 whether flagellar motility is associated with broader life history strategies. We then  
137 applied a method to estimate flagellar motility as a community-aggregated trait  
138 directly from metagenomes. We used this method to investigate the prevalence of

139 flagellar motility across four independent sample sets that we would expect to  
140 capture gradients in soil C availability, and confirmed our findings using  
141 metagenomes from a soil incubation experiment where C availability was  
142 experimentally manipulated via glucose amendments.

143

## 144 **Results and Discussion**

### 145 **Development of a genomic model to predict flagellar motility in bacteria**

146 Since the genes involved in flagellum assembly are well-described and conserved  
147 across bacterial groups [23], we were able to use information on the  
148 presence/absence of flagellar genes to predict the capacity for flagellar motility from  
149 genomic information alone. We used genomic data for 1225 bacterial strains known  
150 to be motile or non-motile using strain description information compiled in [26] as  
151 training data for a boosted regression machine learning model to predict the capacity  
152 for flagellar motility in bacteria (388 unique strains with known flagellar motility and  
153 837 unique strains with no flagellar motility; Supplementary Data 1). We note that  
154 being non-flagellated does not mean taxa are non-motile. For example, within the  
155 Bacteroidota, which had only 2 flagellated members in the training data, the majority  
156 of aquatic and terrestrial members display gliding motility [42, 43]. Of the initial set of  
157 35 genes we identified from the literature as being associated with flagellum  
158 assembly (Supplementary Data 2), we found that 14 out of these 35 genes were  
159 either not frequently found in the genomes of taxa with experimentally validated  
160 flagellar motility, or occurred in >50% of the genomes of non-flagellated taxa. As  
161 these 14 genes were not useful for predictive purposes, the final model was based  
162 on the presence/absence of 21 genes that were sufficiently prevalent across  
163 bacterial genomes and less frequently found in non-motile taxa (Supplementary Data  
164 2). These genes encode different structural parts of the flagellar apparatus, including  
165 the basal body (*FlaE*, *FliL*, *Flg\_bbr\_C*, *Flg\_bb\_rod*), the flagellar rotor (*FliG\_C*), the  
166 flagellar hook (*FlgD*, *Flg\_hook*, *FliD\_C*, *FliE*), or the M-ring (*YscJ\_FliF\_C*), as well as  
167 multiple proteins for protein export and the flagellins required for flagellar assembly  
168 (Supplementary Data 2). We verified that the presence/absence of this set of genes  
169 could effectively distinguish taxa with flagellar motility from non-flagellated taxa using  
170 Principal Components Analysis (PCA) (Supplementary Figure 1A).

171 Our model inferred that taxa were able to display flagellar motility correctly in all taxa  
172 with experimentally verified flagellar motility, and inferred that taxa were non-  
173 flagellated correctly in 94.5% of the cases (Supplementary Figure 1B). We verified  
174 that many of the genomes that the model incorrectly predicted as having flagellar  
175 motility belonged to strains whose genomes contain the majority of flagellar motility  
176 genes and have sister taxa that do display flagellar motility (Supplementary Data 3).  
177 We also recognize that a number of strains might express flagella under certain

178 environmental conditions that would not be captured with the specific *in vitro*  
179 conditions used for strain isolation and phenotyping. Unsurprisingly, the phylum  
180 Proteobacteria was overrepresented in the phenotypic trait database  
181 (Supplementary Figure 2), but our predictions of flagellar motility for this phylum  
182 were not necessarily more accurate than predictions for other phyla (Supplementary  
183 Figure 3), as it contained numerous taxa considered to be non-motile with flagellated  
184 sister taxa (Supplementary Data 3). We recognize that our dataset is over-  
185 represented by taxa (particularly those within the Proteobacteria, Actinobacteria, and  
186 Firmicutes) that are readily cultivated *in vitro* as those are the only taxa for which  
187 phenotypic information on flagellar motility is available. However, given that the  
188 genes associated with flagellar motility are generally well-conserved across a broad  
189 diversity of bacteria [23], and given that our model was robust across multiple phyla  
190 (Supplementary Figure 3), we expect that our genome-based model can also  
191 effectively predict flagellar motility for taxa with no available phenotypic information  
192 on motility, including taxa not included in our test set.

### 193 **Prevalence of flagellar motility across a broad diversity of bacteria**

194 We next used our validated genome-based model (based on the presence/absence  
195 of 21 genes) to determine how the potential for flagellar motility is distributed across  
196 a broad diversity of bacteria, including a wide range of taxa for which no phenotypic  
197 information on motility is currently available. We did so to assess the degree to which  
198 flagellar motility is predictable based on taxonomic or phylogenetic information, and  
199 to investigate the genomic attributes that are generally associated with flagellar  
200 motility. We predicted the capacity for flagellar motility for 26,192 bacterial genomes  
201 spanning 12 major phyla (covering all high-quality genomes in GTDB r207 [27],  
202 belonging to the main bacterial phyla; see Methods). The predicted prevalence of  
203 flagellar motility was highly variable among phyla, ranging from the phylum  
204 Spirochaetota, which had the highest proportion of flagellated taxa (93.2%) to the  
205 Deinococcota and Mycoplasmatota which our model suggests do not have any  
206 flagellated members (Figure 1A). Among the phyla with the largest number of  
207 genomes, we found that the Proteobacteria are predominantly flagellated (78.3%),  
208 with lower proportions for the Firmicutes (54.6%), and very low proportions of  
209 flagellated taxa in the phyla Actinobacteriota (15.9%) and Bacteroidota (0.7%)  
210 (Figure 1A).

211 The majority of bacterial phyla contain numerous families with both flagellated and  
212 non-flagellated members (Supplementary Figure 4). This means that family-level  
213 taxonomic information alone cannot necessarily provide robust inferences of flagellar  
214 motility, stressing the need for alternative approaches to evaluate the prevalence of  
215 this trait across microbial communities. However, at the genus level, most taxa are  
216 either flagellated or non-flagellated, indicating that this trait is typically conserved at  
217 this level of taxonomic resolution (Supplementary Figure 5).

218 Consistent with our taxonomic analyses, the phylogenetic analyses also highlight  
219 that the prevalence of flagellar motility is highly variable at broad taxonomic levels,  
220 which was reflected in a weak phylogenetic signal (phylogenetic  $D = -0.077$ ,  $P <$   
221  $0.001$ ; Figure 1B). Higher-resolution phylogenetic and taxonomic information can  
222 often be useful for inferring flagellar motility, particularly for those groups that are  
223 well-characterized (i.e., where information on flagellar motility, or lack thereof, is  
224 available for closely related taxa). However, phenotypic information is often  
225 unavailable for the broad diversity of taxa found in environmental samples,  
226 highlighting the utility of the genome-based predictive approach described here that  
227 makes it feasible to leverage the rapidly expanding databases of bacterial genomes  
228 to comprehensively investigate the prevalence of this trait in microbial communities.

### 229 **Is there a general life history strategy associated with flagellar motility in** 230 **bacteria?**

231 We expect that bacteria with the capacity for flagellar motility should have distinct  
232 ecologies from non-flagellated taxa. In particular, we expect that flagellated taxa  
233 should be capable of more rapid growth and a greater capacity for carbohydrate  
234 degradation than non-flagellated taxa (i.e. a ‘resource-acquisition’ life history  
235 strategy; [31]). By analyzing the 26,192 genomes for which we had inferred the  
236 capacity for flagellar motility, we were able to identify genomic attributes that were  
237 consistently associated with flagellar motility, conducting these analyses separately  
238 for each of the 6 phyla which had sufficient representation of both flagellated and  
239 non-flagellated taxa (Figure 2A; see Methods). Besides the expected  
240 overrepresentation of genes for motility and extracellular structures (Figure 2A), the  
241 two gene categories that were consistently over-represented in taxa with the  
242 capacity for flagellar motility were signal transduction mechanisms (linked to  
243 chemotaxis) and carbohydrate transport and metabolism (Figure 2A). The latter  
244 observation is consistent with our general expectation that flagellar motility should be  
245 associated with a ‘resource-acquisition’ life history strategy (sensu [31]). However,  
246 we note that this pattern was only evident in 4 of the 6 phyla examined (Figure 2A).

247 We also found that 51.3% of genomes obtained from cultured isolates were  
248 predicted to be flagellated, compared to only 35.2% of genomes of assembled origin  
249 (metagenome-assembled and single cell-assembled genomes, MAGs and SAGs,  
250 respectively; Figure 2B). As culture collections are generally biased towards faster  
251 growing bacterial taxa with adaptations for rapid substrate uptake [44], these results  
252 provide additional support for the hypothesis that flagellar motility is often indicative  
253 of a ‘resource acquisition’ life history strategy [31].

254 To complement these analyses, we also determined the total number of 16S rRNA  
255 gene copies per genome as a proxy for maximum potential growth rate in bacteria  
256 [45]. The number of 16S rRNA gene copies was significantly higher in genomes of  
257 taxa inferred to have the capacity for flagellar motility (Mann-Whitney U  $P < 0.001$ ;

258 Figure 2C), but this pattern was only significant in two phyla (Firmicutes and the  
259 Proteobacteria, Supplementary Figure 6A). Genome size was also significantly  
260 larger in taxa predicted to display flagellar motility (Mann-Whitney U  $P < 0.001$ ;  
261 Figure 2D), a pattern that was consistent across all phyla except for the  
262 Actinobacteriota (Supplementary Figure 6B), and agrees with previous work [29]. We  
263 additionally verified that flagellated taxa harboured a significantly higher number of  
264 genes for chemotaxis than taxa predicted to be non-flagellated (Figure 2E), as we  
265 would expect (Figure 2A; [4]).

266 Together, our genomic analyses suggest that bacteria with flagellar motility tend to  
267 be capable of more rapid growth and the rapid acquisition of organic C substrates,  
268 but this pattern is variable across phyla. Consistent with our findings, a recent global  
269 classification of life history strategies in bacteria found flagellar motility to be  
270 associated with elevated genomic capacity for carbohydrate metabolism, higher 16S  
271 rRNA gene copy numbers, and larger genomes [29]. Recent studies focusing on soil  
272 bacterial communities have had similar findings [38, 46], and in aquatic  
273 environments flagellar motility is considered a signature of copiotrophic lifestyles  
274 [30]. Overall, our findings suggest that flagellar motility is often part of a general life  
275 history strategy for rapid organic carbon metabolism and high maximum potential  
276 growth [31], recognizing that these analyses are based on a biased subset of  
277 bacterial diversity [44] given that most of the genomes included in this analysis  
278 (83%) were derived from cultivated isolates.

## 279 **Application of a metagenome-based approach to quantify the prevalence of** 280 **flagellar motility in bacterial communities**

281 We next extended our genome-based method so it could be used to infer the  
282 prevalence of flagellar motility in whole communities. As the prevalence of flagellar  
283 motility is difficult to reliably infer from taxonomic information alone (see above), and  
284 because neither genomic data nor phenotypic information is available for many  
285 environmental bacteria, we used a metagenome-based method to quantify flagellar  
286 motility as a community-aggregated trait [47]. This method is based on calculating  
287 the ratio between the 21 genes identified as being indicative of flagellar motility and  
288 single-copy marker genes detected per metagenome (see Methods and overview  
289 provided in Figure 3A). We first validated this metagenomic approach using  
290 simulated metagenomic data (see Methods). The simulated data were derived by  
291 mixing different proportions of genomes from taxa with experimentally-verified  
292 flagellar motility capabilities, creating a gradient of metagenomes containing between  
293 0% and 100% flagellated taxa (Figure 3B). This allowed us to obtain a linear  
294 equation to predict the prevalence of flagellated bacteria in any given metagenome  
295 based on the summed abundances of the 21 genes indicative of flagellar motility (as  
296 determined from the genomic analyses above) to the summed abundances of single-  
297 copy genes shared across nearly all bacteria (using a similar approach to [48];  
298 Figure 3A). With these simulated metagenomes, the ratio between the median gene



299 length-corrected reads per kilobase assigned to flagellar and single-copy marker  
300 genes was strongly correlated with the proportion of flagellated taxa in bacterial  
301 communities assembled *in silico* (Pearson's correlation  $r = 0.99$ ,  $P < 0.0001$ ; Figure  
302 3B). We further validated the approach with metagenomic data obtained by  
303 sequencing a DNA mixture from the commercial ZymoBIOMICS microbial  
304 community standard, which contains known amounts of genomic DNA from different  
305 bacterial taxa whose flagellar motility capabilities are known *a priori* (see Methods).  
306 We found that this method accurately inferred the proportion of taxa that were  
307 flagellated based on metagenomic information alone (estimated proportion of  
308 flagellated taxa = 52.0%, expected proportion of flagellated taxa = 48.2%; Figure  
309 3B). We also verified that our estimates using only the forward reads did not differ  
310 from those using the reverse or merged reads (Figure 3B). Together, these results  
311 highlight that we can accurately infer the community-level prevalence of bacterial  
312 flagellar motility in any metagenome of interest simply by calculating the ratio  
313 between the sum of the 21 flagellar genes and the sum of single-copy bacterial  
314 marker genes.

### 315 **Prevalence of bacterial flagellar motility across gradients in soil carbon** 316 **availability**

317 We used our metagenome-based approach to further test our hypothesis that  
318 flagellar motility is most likely to be associated with taxa adapted for fast resource  
319 acquisition under resource-rich conditions (Figure 2A-C). If this hypothesis is valid,  
320 we would expect the community-wide prevalence of bacterial flagellar motility to be  
321 higher in soils with greater amounts of available organic C. Since it is challenging to  
322 directly quantify the amount of C in soil that is available to fuel microbial activities, we  
323 selected 4 independent metagenomic datasets that we expect to effectively capture  
324 gradients in soil C availability, and the results from the analyses of these datasets  
325 are described below.

326 Soil C availability is expected to decrease with soil depth [49, 50]. Across the 9 soil  
327 depth profiles analyzed [51], we consistently observed a higher prevalence of  
328 flagellar motility in the surface (top 20cm) compared to deeper soil horizons (20-  
329 90cm, linear mixed effects model,  $\text{Estimate}_{\text{Surface}} = 11.88 \pm 1.30$  (mean  $\pm$  SD),  
330  $\text{Estimate}_{\text{Subsurface}} = 8.64 \pm 1.34$ ,  $P = 0.005$ ,  $N = 66$ ; Figure 4A; Supplementary Figure  
331 7). We recognize that soil C availability is not the only factor that is likely to change  
332 appreciably with soil depth. For example, soil water and nutrient availability can also  
333 vary with depth [51], so we cannot conclude that soil C availability is the only factor  
334 responsible for the elevated prevalence of flagellar motility in surface soil  
335 communities.

336 To further test our hypothesis, we analyzed 38 surface soils collected from across  
337 Australia. For this sample set, we assume that net primary productivity (NPP) is a  
338 reasonable proxy for soil C availability, as higher NPP leads to increased plant-

339 derived organic matter inputs to soil [52]. We found that across these varied soils,  
340 the prevalence of flagellar motility in bacterial communities was strongly correlated  
341 with NPP (Pearson's  $r = 0.619$ ,  $P < 0.001$ ; Figure 4B). As with the 'soil depth'  
342 analyses, these findings also support our hypothesis that flagellar motility is more  
343 prevalent in soils with higher C availability. However, as other factors likely co-vary  
344 with NPP (including mean annual precipitation), these findings on their own are not  
345 sufficient to confidently support a general association between flagellar motility and  
346 soil C availability.

347 As both the 'soil depth' and the 'Australian surface soil' datasets indicate an  
348 association between inferred soil C availability and flagellar motility, we then sought  
349 to determine the prevalence of flagellar motility in bacteria from rhizospheres and  
350 associated bulk soils. While many factors differ between rhizosphere and bulk soils,  
351 we would expect that soil C availability is one of the more prominent factors differing  
352 between these two soil habitats. The rhizosphere receives abundant inputs of  
353 available plant-derived C via root exudation [53], and rhizosphere soils generally  
354 support higher microbial respiration rates than adjacent bulk soils [37, 54]. We  
355 analyzed two independent metagenomic datasets that compared bacterial  
356 communities in rhizospheres and adjacent bulk soils. One dataset contained paired  
357 rhizosphere and bulk soil samples across the globe from diverse citrus species ([55],  
358  $N = 20$ ), and the other dataset contained samples from a controlled pot experiment  
359 with wheat plants ([56],  $N = 24$ ). We found that in both datasets, rhizosphere  
360 bacterial communities consistently had a higher prevalence of flagellar motility  
361 compared to their adjacent bulk soils (Figures 4C,D). Across citrus species, the  
362 prevalence of flagellar motility was on average a 11.5% higher in rhizospheres than  
363 in bulk soils (one-sample t-test  $P = 0.012$ ; Figure 4C; Supplementary Figure 8), and  
364 was higher in the rhizosphere than in the paired bulk soil in 9 out of the 10 sites  
365 analyzed. In wheat plants, we also found that rhizosphere bacterial communities  
366 contained a higher prevalence of flagellar motility ( $\text{Estimate}_{\text{Rhizosphere}} = 23.7 \pm 2.6$ )  
367 than bulk soils ( $\text{Estimate}_{\text{Bulk soil}} = 11.3 \pm 8.0$ ; Welch two-sample t-test  $P = 0.0002$ ;  
368 Figure 4D). While we recognize that other factors could contribute to the elevated  
369 prevalence of flagellar motility in rhizosphere communities, these results provide  
370 further support for our hypothesis that flagellar motility is favored under conditions of  
371 higher soil C availability, as also indicated by the analyses of the 'soil depth' and the  
372 'Australian surface soil' datasets.

### 373 **Experimental verification that bacterial flagellar motility is associated with soil** 374 **carbon availability**

375 To more conclusively test whether soil C availability is associated with the  
376 prevalence of bacterial flagellar motility, we generated new metagenomic data from a  
377 117-day soil incubation experiment where C availability was directly manipulated via  
378 regular glucose amendments (see Methods; [57]). This experiment was performed in  
379 the absence of a growing plant and under uniform moisture conditions, thus

380 minimizing the impact of these potential confounding factors. The prevalence of  
381 flagellar motility in the bacterial communities amended with glucose ( $15.82 \pm 0.89\%$ )  
382 was higher than in the soils that did not receive glucose ( $13.19 \pm 1.56\%$ ) (Welch two-  
383 sample t-test  $P = 0.017$ ; Figure 5A). This pattern is in line with the results from the  
384 field studies (Figure 4) and supports our central hypothesis that the prevalence of  
385 flagellar motility is positively associated with soil C availability. The rather small size  
386 of these effects is likely due to the fact that relatively few bacterial taxa responded to  
387 the glucose addition. While glucose addition shifted the overall community  
388 composition (Figure 5B), only 28 bacterial taxa (ASVs) out of the total 1203 ASVs  
389 detected were significantly more abundant in the glucose-amended soils. These taxa  
390 that significantly responded to glucose addition belonged to 7 different bacterial  
391 phyla (Supplementary Figure 9) .

## 392 **Conclusions**

393 We have shown that flagellar motility is a key trait linking C dynamics and microbial  
394 communities in soil. Consistent with expectations [31], our genomic analyses reveal  
395 that flagellated taxa tend to be associated with a 'resource-acquisition' life history  
396 strategy. This observation was supported by our metagenomic analyses which  
397 revealed a positive relationship between the prevalence of flagellar motility in  
398 bacterial communities and soil C availability across multiple, independent datasets.  
399 This relationship between flagellar motility and soil C availability can be explained  
400 based on fundamental energetic constraints, which make flagellar motility a  
401 beneficial trait in environments where C availability is elevated, particularly in  
402 spatially structured environments like soil where available C can be patchily  
403 distributed [35].

404 The methods to predict microbial traits from genomic information presented here are  
405 particularly relevant for traits that are difficult to quantify *in situ* or for those that  
406 require isolation and culturing [58, 59]. Our metagenome-based approach to infer the  
407 proportion of a microbial community harboring any given phenotypic trait would be  
408 very useful for this purpose (Figure 3A). This method can also be applied to  
409 investigate processes where flagellar motility is expected to play an important role,  
410 such as microbial colonization and persistence in host-associated microbiomes [60,  
411 61]. In efforts to improve microbiome management, a better quantification of the  
412 prevalence of flagellar motility in these systems could help identify microbiomes that  
413 are likely to be more persistent in the host or more likely to deliver beneficial  
414 functions [62]. These methods could also be used to explore the prevalence of  
415 motility and its associated traits across gradients in C availability in other  
416 environments of interest, such as freshwater systems. Overall, genome-based  
417 predictive approaches offer opportunities for expanding our trait-based  
418 understanding of microbial communities beyond cultivated taxa, and help us  
419 understand microbial community patterns across environmental gradients.

420

## 421 **Materials and Methods**

### 422 *Genome selection and annotation*

423 We compiled genomic data from ~62,000 unique bacterial taxa ('species clusters')  
424 available in the Genome Taxonomy Database (GTDB) (release 207; [27]). We  
425 restricted our analyses to bacterial phyla with more than 100 representative  
426 genomes available in GTDB and only included genomes estimated to be >95%  
427 complete based on CheckM (v1.1.6) [63]. We also removed all genomes that lacked  
428 a 16S rRNA gene, as well as those with signals of chimerism based on GUNC  
429 (Genome Unclutterer; [64]), yielding 26,192 genomes in total belonging to 12  
430 different phyla.

431 The coding sequences of the 26,192 genomes were identified using Prodigal (v2.6.3;  
432 [65]). We then aligned the predicted coding sequences for each genome to the Pfam  
433 database (v35.0; [66]) using HMMER (v3; [67]) to obtain information on all potential  
434 domains and genes present in those genomes. All matches with a bit score lower  
435 than 10 were discarded. We then binarized all copy numbers of genes and domains  
436 in each genome to presence/absence for further analyses. We selected a set of 21  
437 genes out of a total of 35 genes involved in flagellar assembly in Pfam based on their  
438 prevalence among strains with empirical information on flagellar motility  
439 (Supplementary Data 2). Specifically, this subset of genes was chosen based on the  
440 following criteria: 1) genes were present in >80% of taxa with experimentally  
441 demonstrated flagellar motility, and 2) genes were not present in >50% of taxa  
442 classified as non-motile based on available phenotypic information (see below). This  
443 step was necessary as many non-motile taxa conserve genes for flagellar motility  
444 (Supplementary Data 3), and some of the flagellar genes are not well represented in  
445 Pfam. We used the information on the presence/absence of these 21 genes across  
446 genomes to build a predictive model of flagellar motility in bacteria (Supplementary  
447 Data 2).

### 448 *Genome-based prediction of flagellar motility in bacteria*

449 We compiled all information on whether bacterial taxa displayed flagellar motility or  
450 not from the bacterial phenotypic trait database compiled in [26]. This database  
451 contains information on motility traits for 13,481 unique bacterial strains [26]. We first  
452 selected only the subsets categorized as having flagellar motility or as being non-  
453 motile (8191 unique strains). To obtain representative genomes for these strains, we  
454 matched the National Center for Biotechnology Information (NCBI) taxon id of each  
455 of these strains to their corresponding genome accession in GTDB. To ensure  
456 maximal reliability of the genomic information used for model training, we only kept

457 those genomes that were 100% complete, and applied the same quality filters  
458 mentioned above. This led to a final subset of 1225 high quality genomes (388  
459 categorized as having flagellar motility, 837 non-motile) that we used for model  
460 training (Supplementary Data 1). We note that these 1225 genomes included taxa  
461 from 18 unique phyla, with the proportions of motile taxa per phylum ranging from 0-  
462 100% (Supplementary Figure 2).

463 Since some of the genes involved in flagellar assembly are often present in several  
464 non-motile taxa ([68]; Supplementary Data 3), we were not able to use standard  
465 statistical approaches to build a predictive model of flagellar motility based on the  
466 presence/absence of 21 flagellar genes. We thus used gradient boosted regression  
467 decision trees that could accommodate the complexity of having 21 predictive  
468 features using the *xgboost* package in R (v1.7.5; [69]). To this end, we first built a  
469 training and a test set (70:30, randomly selected) of the matrix containing the  
470 presence/absence of the flagellar genes for each of the representative bacterial  
471 genomes with experimental information on flagellar motility using the *xgb.DMatrix*  
472 function of *xgboost*. We then applied Bayesian hyperparameter optimization to select  
473 the best parameters for the regressor model using the *bayesOpt* function of the  
474 *ParBayesianOptimization* R package (v1.2.6; [70]), specifying the objective function  
475 as a binary logistic regression. We ran k-fold cross-validation using the *xgb.cv*  
476 function in *xgboost* to identify the optimal number of iterations of model improvement  
477 for the final model training function. We built the boosted regression model using the  
478 optimized parameters and iterations calculated above using the *xgboost* function of  
479 package *xgboost*. We used the function *xgb.importance* from *xgboost* to compute the  
480 predictive importance of the different genes in the final model, which identified 14  
481 flagellar genes that were most useful for predicting flagellar motility (Supplementary  
482 Data 4), even though the full set of 21 genes was needed for accurate prediction  
483 (see Results and Discussion for details on the performance of the final selected  
484 model). We evaluated model performance using the accuracy index.

#### 485 *Phylogenetic analysis*

486 To investigate the phylogenetic distribution of flagellar motility in bacteria, we first  
487 randomly selected a single genome from each family within the 12 predominant  
488 phyla investigated (485 genomes in total). Since we had already predicted the  
489 potential for flagellar motility across GTDB genomes, we simply subsetted the tree  
490 provided by GTDB with the selected genomes, which can be found in  
491 [https://data.gtdb.ecogenomic.org/releases/release207/207.0/bac120\\_r207.tree](https://data.gtdb.ecogenomic.org/releases/release207/207.0/bac120_r207.tree). This  
492 tree is based on the alignment of 120 single-copy marker genes and is therefore  
493 more robust than a conventional maximum likelihood tree based on the alignment of  
494 full 16S rRNA gene fragments. We visualized and edited the trees using iTOL (v5;  
495 [71]). We tested whether flagellar motility had a phylogenetic signal by calculating  
496 the phylogenetic *D* index for binary traits [72], where values (positive or negative)  
497 closer to 0 indicate phylogenetic conservatism, and values closer to 1 indicate a

498 random phylogenetic pattern. This phylogenetic analysis was conducted using the R  
499 package ape (v5.7-1; [73]). We additionally explored the degree of conservatism of  
500 flagellar motility across different levels of phylogenetic resolution by measuring the  
501 standard deviation (SD) of the flagellar motility status (flagellated, 1; non-flagellated,  
502 0) across taxa from different taxonomic ranks (phyla, classes, orders, families, and  
503 genera). For this analysis, we only included those taxa that were represented by  
504 more than one genome.

#### 505 *Analysis of bacterial life history strategies associated with flagellar motility*

506 We investigated associations between flagellar motility and broad functional gene  
507 categories by testing the prevalence of Clusters of Orthologous Genes (COGs) in the  
508 genomes of taxa predicted to be flagellated or non-flagellated [74]. We excluded the  
509 phyla Bacteroidota, Chloroflexota, Cyanobacteria, Spirochaetota, and  
510 Mycoplasmatota from this analysis as these phyla had either too high (>90%) or too  
511 low (<15%) proportions of flagellated taxa to perform robust statistical comparisons.  
512 We annotated genomes (N = 21,551) into COG categories using eggNOG-mapper  
513 v2 [75], and calculated the genome size-corrected prevalence of each COG category  
514 per genome. We also investigated general genomic features such as genome size  
515 and the 16S rRNA gene copy number for each of the genomes to compare these  
516 genomic attributes between motile and non-motile taxa within each phylum. We  
517 included 16S rRNA gene copy number as it is considered a proxy for maximal  
518 potential growth rates in bacteria [45]. We compiled and identified the genes involved  
519 in chemotaxis (Supplementary Data 5) across the genomes of flagellated and non-  
520 flagellated taxa as a validation given the chemotaxis signaling pathway is an  
521 activator of the flagellar motor system [4].

#### 522 *Estimation of the prevalence of flagellar motility in microbial communities using* 523 *metagenomic information*

524 We applied a method to estimate the prevalence of flagellar motility as a community-  
525 aggregated trait using metagenomic information from bacterial communities (Figure  
526 3A). To this end, we first assembled 'mock' metagenomes containing different  
527 proportions of genomes from flagellated and non-flagellated taxa from the subset we  
528 originally used for boosted regression model training. We selected 20 genomes of  
529 taxa with empirically verified flagellar motility capabilities spanning the phyla  
530 Proteobacteria, Firmicutes, and Actinobacteria as these are ubiquitous taxa and are  
531 well-represented in our training data. We used ART (a next-generation sequencing  
532 simulator; [76]) to simulate short (150bp) shotgun sequencing reads at a coverage of  
533 50% of these genome mixtures. We did not choose higher coverage as soil  
534 metagenomic datasets do not usually exceed 50% community coverage [77]. We  
535 then constructed a DIAMOND (v2.0.7; [78]) database containing the protein variants  
536 for each of the 21 selected genes (7-633 variants per gene) identified from the  
537 genomic analyses (see above) that were determined to be robust predictors of

538 flagellar motility, as well as the variants contained in GTDB for the 120 single-copy  
539 marker genes that constitute the taxonomic basis of GTDB [27]. We annotated the  
540 simulated metagenomes using blastx (v2.13.0; [79]) on this custom protein  
541 database. In this way, we obtained a reads-per-kilobase (RPK) index for both the  
542 flagellar gene and the single-copy marker gene sets by taking the median gene  
543 length-corrected number of hits of each protein across the 21 and 120 unique  
544 proteins, respectively. We finally built a ‘flagellar motility index’ based on the ratio  
545 between the flagellar gene RPK and the single-copy marker gene RPK (see  
546 overview in Figure 3A). The use of single-copy marker genes in this manner offers a  
547 general normalization of the flagellar gene read count – which can vary due to  
548 differences in library size, coverage, or diversity – as these single copy genes are  
549 assumed to be present in every bacterial genome [80]. We then determined the  
550 linear relationship between the ‘flagellar motility index’ and the proportions of  
551 genomes that were able to produce flagella across mock metagenomes, following a  
552 similar approach to [48]. We used this standard curve to estimate the proportion of  
553 genomes in a given metagenome that are able to produce flagella based on the  
554 ‘flagellar motility index’, as expressed in equation (1):

555 (1)  $\% \text{ of bacteria with flagellar motility} = 3650 \times \text{Flagellar motility index} - 0.321$

556 where the ‘flagellar motility index’ is the ratio between the median RPK of the 21  
557 flagellar motility genes over the median RPK of the 120 single-copy marker genes  
558 (Figure 3A). This method allows the estimation of the prevalence of flagellar motility  
559 in any given bacterial metagenome based on the assumption that flagellar genes are  
560 usually found in single copies among bacterial genomes [23].

### 561 *Testing associations between bacterial flagellar motility and soil carbon availability*

562 We selected metagenomic datasets that covered expected gradients in soil C  
563 availability, which we hypothesized to be positively associated with bacterial flagellar  
564 motility. Soil C availability is challenging to measure *in situ* and direct measurements  
565 of soil C availability (which is not equivalent to total C concentrations) are rarely  
566 compiled along with metagenomic data. We thus selected datasets that we expect  
567 based on published research to span gradients in C availability, recognizing that C  
568 availability is often correlated with other soil variables. The datasets included are the  
569 following: 1) soils from across the USA spanning gradients in soil depth (surface, 0-  
570 20; subsurface, 20-90cm, N = 66), where total organic C decreases with depth [51];  
571 2) a net primary productivity (NPP) gradient across Australia (N = 38, [81]), where  
572 higher NPP is expected to be associated with higher soil organic C availability [52];  
573 3) a global comparison of rhizosphere and bulk soils associated with citrus plants (N  
574 = 20, [55]), where we would expect C availability to be higher in rhizosphere soils  
575 than in bulk soils [82]; and 4) a pot experiment with controlled water inputs  
576 comparing the rhizosphere and adjacent bulk soil of wheat plants (N = 24, [56]).

577 Since all these datasets contain factors that likely covary with soil C availability, we  
578 additionally obtained metagenomic data from soils that were incubated with or  
579 without glucose amendments over a 117-d incubation period in a previous study [57].  
580 In this experiment, glucose was added weekly to sub-samples of a single soil at a  
581 rate of 260  $\mu\text{g C g dry wt soil}^{-1} \text{ day}^{-1}$  (see [57] for full details). Since this experiment  
582 was performed under constant moisture conditions and in the absence of plants [57],  
583 the glucose amendments should lead to an increase in C availability with minimal  
584 direct effects on other soil attributes. The addition of glucose in this experiment led to  
585 a 7.9-fold increase in the microbial  $\text{CO}_2$  respiration rates [57], confirming that the C  
586 available to soil microbes increased in the soils amended with glucose compared to  
587 the controls (i.e. soils that received only an equivalent amount of water).

588 We then generated metagenomic data from the 9 soil samples harvested from the  
589 glucose amendment experiment. For each soil sample (4 with added glucose, 5  
590 without glucose), we used 0.25g of soil for DNA extraction using the DNeasy  
591 PowerSoil Pro tube kit (Qiagen). The shotgun sequencing library was prepared using  
592 Illumina's DNA Prep kit and Unique Dual Indexes (Illumina, CA). Samples were  
593 quantified using Qubit and pooled at equimolar concentrations. The library was run  
594 on a NovaSeq 6000 (Illumina, CA) at the Texas A&M AgriLife Genomics &  
595 Bioinformatics Service (USA) using a 2x150 cycle flow cell. Sequence cluster  
596 identification, quality prefiltering, base calling and uncertainty assessment were done  
597 in real time using Illumina's NCS 1.0.2 and RFV 1.0.2 software (Illumina, CA) with  
598 default parameter settings. We also analyzed the 16S rRNA gene sequencing  
599 information on the same soil communities (see [57] for details on how this data was  
600 generated and processed).

#### 601 *Processing of shotgun metagenomic sequencing reads from datasets covering* 602 *gradients in soil C availability*

603 To process the metagenomic data from all of the datasets described above (157  
604 metagenomes in total), we first downloaded the sequences from the Sequence Read  
605 Archive (SRA) of NCBI when applicable, and ran trimmomatic (v0.39; [83]) to remove  
606 adapters and low quality base pairs using a phred score of 33 as a threshold, only  
607 keeping reads above 100bp after trimming. We used blastx on the custom  
608 DIAMOND database we created to annotate the metagenomic reads. We filtered out  
609 reads that had <50% bit score, <60% identity to the reference protein, and an e-  
610 value higher than 0.001. We finally measured the flagellar motility index based on  
611 the ratio between the median reads-per-kilobase (RPK) of the flagellar genes and  
612 the median RPK of the 120 single-copy marker genes as described above, and fitted  
613 equation (1) to quantify the prevalence of flagellar motility in any given metagenome  
614 using the method outlined in Figure 3A.

#### 615 *Statistical analysis*



616 All statistical analyses were conducted in R (v4.1.3; [84]). We used principal  
617 components analysis (PCA) to visualize how well the presence/absence of the  
618 selected flagellar motility genes was able to discriminate between the genomes of  
619 flagellated and non-flagellated taxa. To identify potential differences in the life history  
620 strategies of flagellated and non-flagellated taxa, we used multiple Mann-Whitney U  
621 tests with Bonferroni correction for multiple comparisons to investigate whether  
622 particular COG categories were overrepresented in genomes from flagellated versus  
623 non-flagellated taxa. The results were presented as the log<sub>2</sub>-fold ratio. We used  
624 Mann-Whitney U tests to investigate associations between flagellar motility and the  
625 16S rRNA gene copy number and the number of chemotaxis genes in any given  
626 genome due to non-normality of the data. We compared differences in genome size  
627 between flagellated and non-flagellated taxa using Welch two-sample t-tests.

628 To test for differences in the prevalence of flagellar motility between surface and  
629 subsurface soils we used a mixed effects linear model with location coded as  
630 random factor, and for the test between rhizosphere and bulk soils in wheat we used  
631 Welch two-sample t-tests. Since we only had a single rhizosphere and bulk soil  
632 observation per site, in the global citrus rhizosphere dataset we first calculated the  
633 difference in the prevalence of flagellar motility in the rhizosphere over bulk soil at  
634 each site, and then tested whether these differences were significantly different from  
635 zero using a one-sample t-test. These tests were implemented using different  
636 arguments of the *t.test* function in base R [84]. We used Pearson's correlations to  
637 evaluate relationships between the prevalence of flagellar motility and NPP, and  
638 used linear regression to represent the standard curve to quantify the prevalence of  
639 flagellar motility in metagenomes.

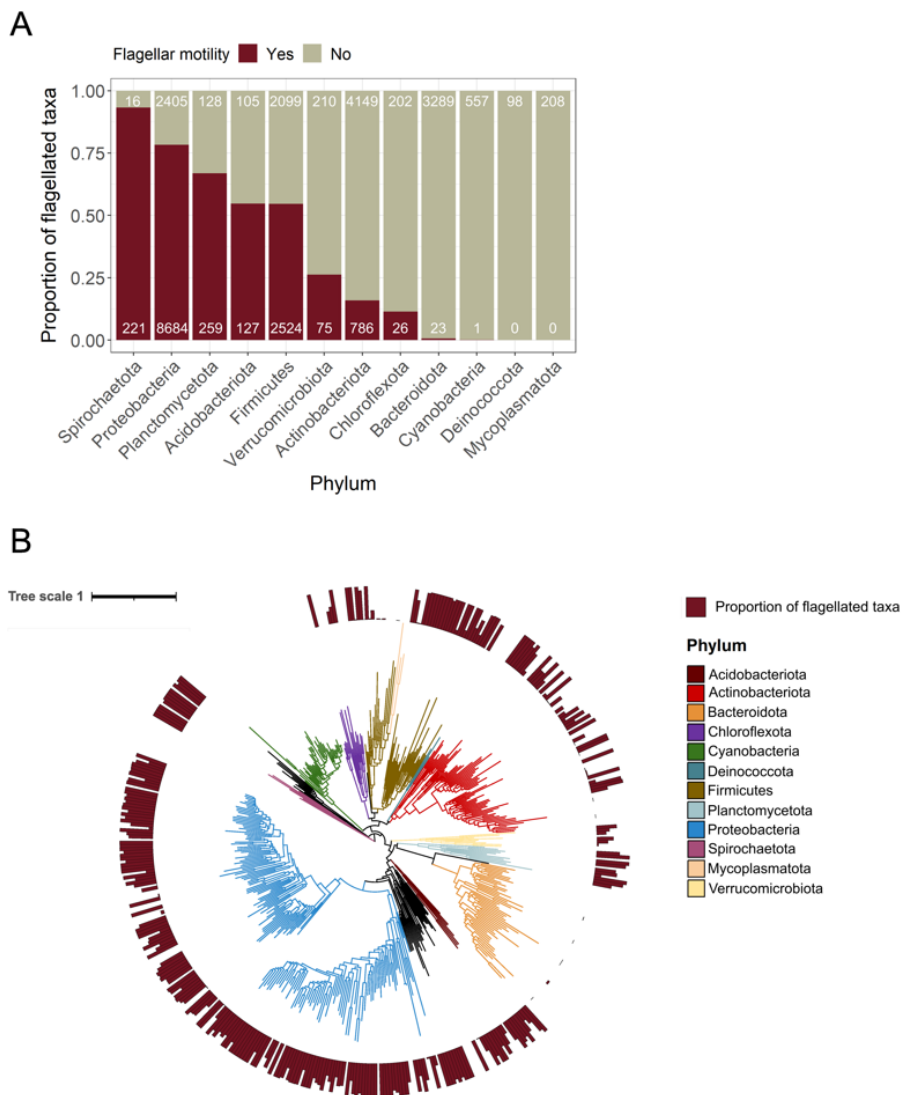
640 Finally, we used 16S rRNA gene sequencing information from samples in the  
641 glucose amendment experiment [57] to investigate the shifts in the taxonomic  
642 composition of soil bacterial communities upon glucose addition. Specifically, we  
643 investigated which bacterial Amplicon Sequence Variants (ASVs) responded to  
644 glucose addition using ANCOM-BC [85]. The taxonomic composition of these  
645 bacterial communities was investigated using the phyloseq R package (v1.38.0;  
646 [86]), and we tested the effect of glucose amendment on the prevalence of flagellar  
647 motility assessed using our metagenome-based method using the Welch two-sample  
648 t-test.

649

650

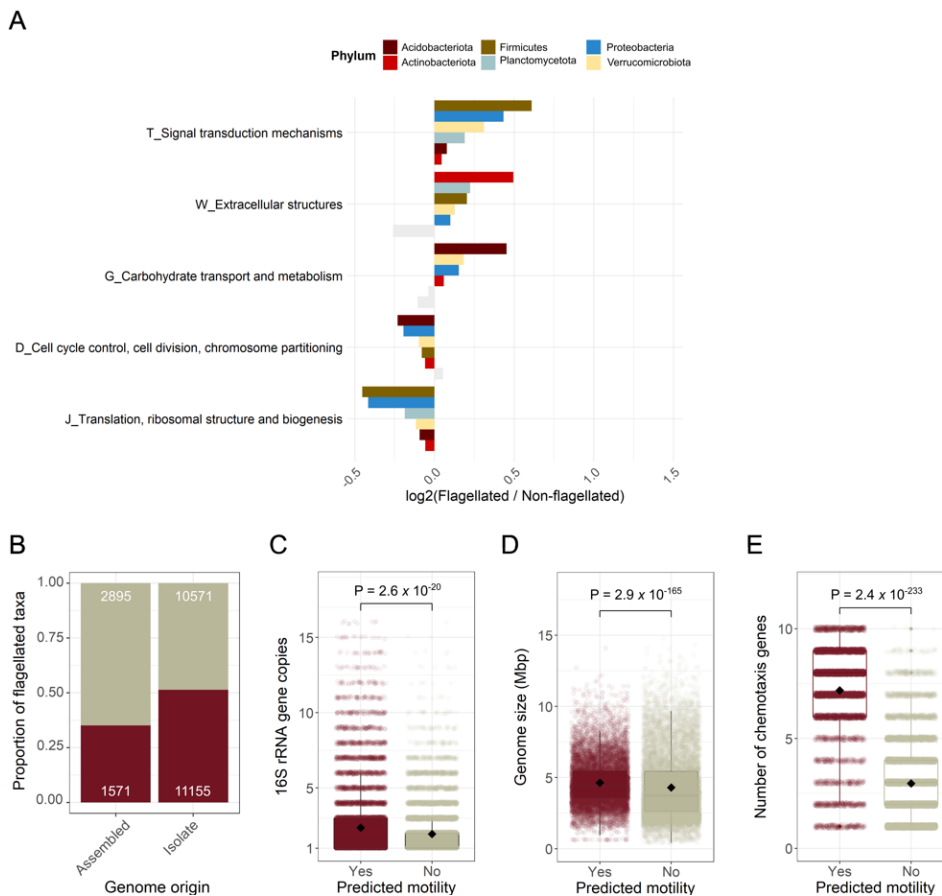
## 651 Figures

652 **Figure 1. Taxonomic and phylogenetic distribution of flagellar motility in**  
 653 **bacteria.** A. Prevalence of flagellar motility in bacterial taxa from the 12 phyla best-  
 654 represented phyla in a curated database of reference genomes (N = 26,192  
 655 genomes). B. Phylogenetic distribution of flagellar motility across the 12 bacterial  
 656 phyla. To construct the tree, we randomly selected a single genome representative  
 657 of each family found in each phylum, and predicted the capacity for flagellar motility  
 658 in these genomes. Higher bars indicate a greater proportion of genomes within that  
 659 family that are inferred to have the capacity for flagellar motility (based on our  
 660 genome-based model, see Methods). The tree was constructed from the Genome  
 661 Taxonomy Database phylogeny (GTDB r207; [27]).



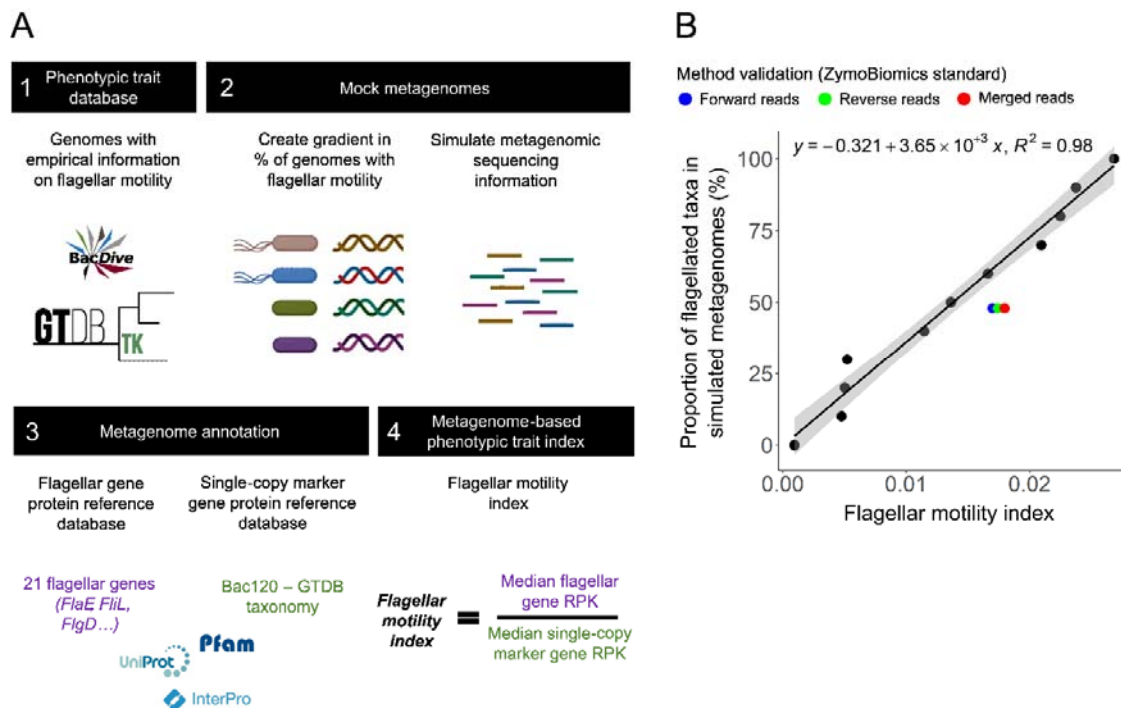
663

664 **Figure 2. Genomic attributes associated with bacteria inferred to have the**  
665 **capacity for flagellar motility.** A. Functional gene categories that are consistently  
666 overrepresented in genomes from taxa predicted to be flagellated or non-flagellated  
667 across the 6 most dominant phyla that contain >15% flagellated taxa. Functional  
668 categories were defined as Clusters of Orthologous Genes (COGs). We indicated  
669 those gene categories that were not statistically different with grey shading based on  
670 Mann-Whitney U tests ( $P > 0.01$ ). B. Prevalence of flagellar motility in genomes  
671 derived from environmental metagenomes (MAGs) or single cells (SAGs)  
672 ('Assembled'), and in genomes obtained from bacterial isolates ('Isolate'). Numbers  
673 on the upper and lower ends of the plot indicate the number of genomes predicted to  
674 be non-flagellated and flagellated, respectively. C. Number of 16S rRNA gene copies  
675 in genomes of taxa predicted to be flagellated ( $N = 12,726$ ) versus non-flagellated ( $N = 13,236$ ).  
676 D. Genome size of taxa predicted to be flagellated and non-flagellated. E.  
677 Number of genes involved in chemotaxis identified in the genomes of taxa that are  
678 flagellated and non-flagellated. In panel A,  $N_{\text{Acidobacteriota}} = 232$ ;  $N_{\text{Actinobacteriota}} = 4935$ ;  
679  $N_{\text{Firmicutes}} = 4623$ ;  $N_{\text{Planctomycetota}} = 387$ ;  $N_{\text{Proteobacteria}} = 11,089$ ;  $N_{\text{Verrucomicrobiota}} = 285$ . In  
680 panels C and E the P-value was obtained from Mann-Whitney U tests due to non-  
681 normality of the data. In panel D, the P-value was obtained from a Welch two-sample  
682 t-test, ( $P < 0.05$ );  $N = 26,192$  genomes.

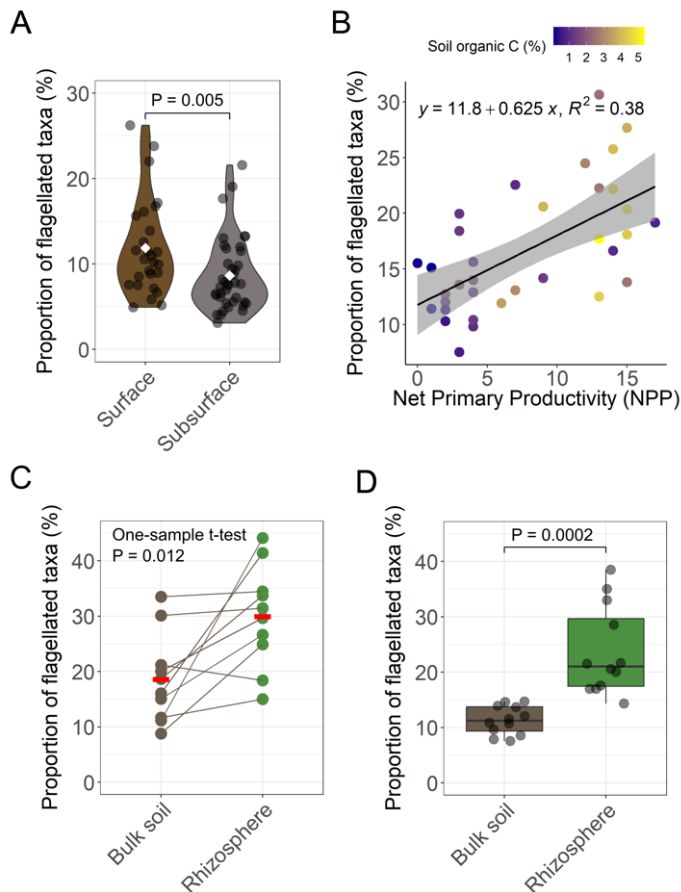


683 **Figure 3. Developing a metagenome-based approach to quantify the**  
 684 **prevalence of flagellar motility in bacterial communities.** A. Method overview.  
 685 We first collected whole-genome data for bacterial taxa directly observed to have  
 686 flagellar motility *in vitro* (1). We then make combinations of genomes with and  
 687 without the capacity for flagellar motility to create a gradient of the prevalence of  
 688 flagellar motility in ‘mock’ metagenomes (2). These ‘mock’ metagenomes are created  
 689 by simulating shotgun metagenomic reads from the whole genomes (see Methods).  
 690 We annotate the metagenomes to identify the 21 flagellar genes determined from the  
 691 genomic analyses to be indicative of flagellar motility along with a set of 120 single-  
 692 copy marker genes that are found in nearly all bacteria (see Methods) (3). Finally, we  
 693 calculate the gene length-corrected reads-per-kilobase (RPK) of each of these gene  
 694 sets and calculate a ‘flagellar motility index’ using the ratio between these indices  
 695 (4). B. Linear relationship between the ‘flagellar motility index’ calculated as shown in  
 696 panel A (4) and the proportion of genomes of taxa with flagellar motility in simulated  
 697 metagenomes (panel A, 2) (N = 14). The y-axis shows a gradient of bacterial  
 698 metagenomes created by combining different proportions of genomes from bacteria  
 699 known to be flagellated or non-flagellated spanning the phyla Proteobacteria,  
 700 Actinobacteria, and Firmicutes. The linear equation resulting from this association  
 701 can be used to quantify the prevalence of flagellar motility in any bacterial  
 702 metagenome. Colored dots indicate the known proportion of flagellated taxa in the  
 703 metagenome of the ZymoBiomics microbial community standard.

704

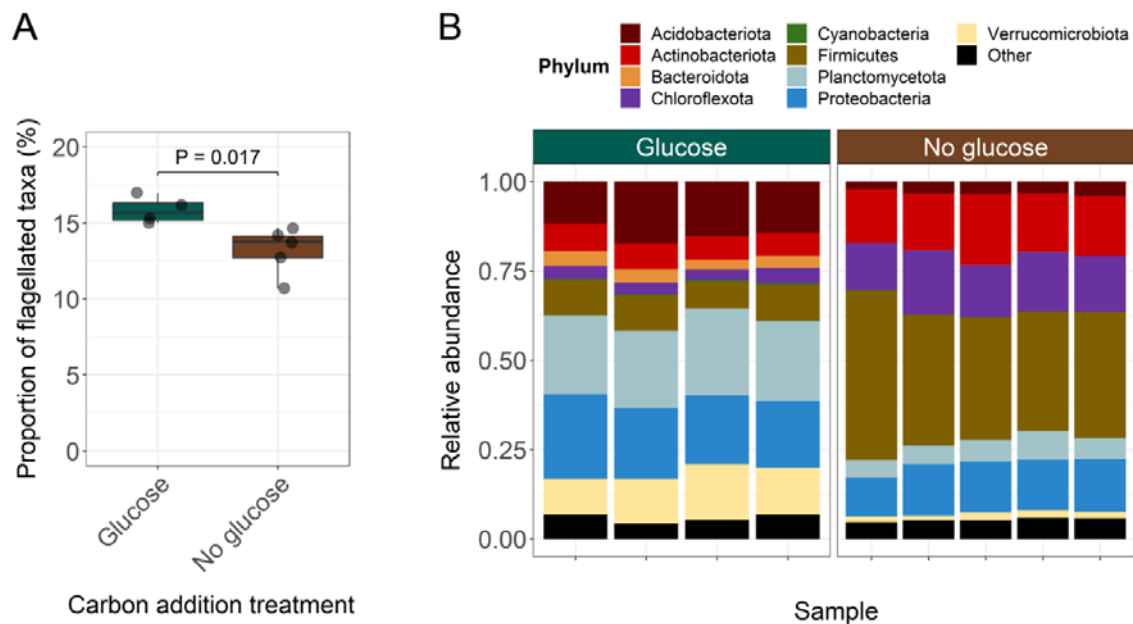


705 **Figure 4. Prevalence of flagellar motility in bacterial communities spanning**  
706 **putative gradients in soil carbon (C) availability.** A. Estimated prevalence of  
707 flagellar motility in bacterial communities found across soil profiles (Surface, 0-20cm;  
708 Subsurface, 20-90cm; N = 66). These soil profiles were sampled from sites that  
709 covered diverse climatic regions across the USA [51]. Group means are shown as  
710 white diamonds, and the P-value was obtained from a linear model with site coded  
711 as a random factor. B. Relationship between the estimated prevalence of flagellar  
712 motility and net primary productivity (NPP) across Australia (N = 38; [81]). The  
713 shaded area depicts the standard error around the mean. C. Comparison of the  
714 prevalence of flagellar motility in bulk soils and rhizospheres of citrus trees found at  
715 10 sites across the globe (N = 20; [55]). Since each site contained a single bulk soil  
716 and a single rhizosphere sample, we indicate which samples come from the same  
717 site using connecting lines. To obtain the P-value, we calculated the difference  
718 between the prevalence of flagellar motility in the rhizosphere and bulk soil at each  
719 site, and then made a comparison against zero using a one-sample t-test. Means are  
720 shown as horizontal red lines. D. Comparison of the prevalence of flagellar motility in  
721 bulk soils and rhizospheres of wheat plants from a controlled pot experiment (N = 24;  
722 [56]). The P-value was obtained using a Welch two-sample t-test. Statistical  
723 significance is set at  $P < 0.05$ .



724 **Figure 5. Prevalence of flagellar motility in bacterial communities from a 117-d**  
725 **soil incubation experiment where soil carbon (C) availability was directly**  
726 **manipulated via addition of glucose.** A. Estimated prevalence of flagellar motility  
727 in bacterial communities found in soils amended (N = 4) and not amended (N = 5)  
728 with glucose as a way to directly manipulate soil C availability [57]. The P-value was  
729 obtained using a Welch two-sample t-test with significance at  $P < 0.05$ . B.  
730 Taxonomic composition of bacterial communities from soils amended or non-  
731 amended with glucose over 117 days of incubation. The taxonomic composition of  
732 the bacterial communities was determined via amplicon sequencing of the 16S rRNA  
733 gene (see Methods).

734



735

736 **Acknowledgements**

737 We thank Michael Hoffert and Thomas B. N. Jensen for assistance with the  
738 bioinformatical analyses, and Jessica Henley, Caihong Vanderburgh, and Jordan  
739 Galletta for generating the metagenomic sequence data.

740 **Author contributions**

741 JR and NF conceived and designed the study. JR performed the data analyses. KF,  
742 JML, HC, AB, and MSS contributed data to the study. JR and NF wrote the  
743 manuscript, with input from all co-authors.

744 **Funding statement**

745 JR acknowledges funding from the Swiss National Science Foundation (Early  
746 PostDoc Mobility grant P2EZP3\_199849). Funding was also provided by a grant  
747 from the US National Science Foundation awarded to NF and MSS (award #  
748 2131837).

749 **Conflicts of interest**

750 The authors declare no conflicts of interest.

## 751 References

- 752 1. Miyata M, Robinson RC, Uyeda TQP, Fukumori Y, Fukushima S ichi, Haruta  
753 S, et al. Tree of motility – A proposed history of motility systems in the tree of  
754 life. *Genes to Cells* 2020; **25**: 6–21.
- 755 2. Keegstra JM, Carrara F, Stocker R. The ecological roles of bacterial  
756 chemotaxis. *Nature Reviews Microbiology* 2022; **20**: 491–504.
- 757 3. Cremer J, Honda T, Tang Y, Wong-Ng J, Vergassola M, Hwa T. Chemotaxis  
758 as a navigation strategy to boost range expansion. *Nature* 2019; **575**: 658–  
759 663.
- 760 4. Colin R, Ni B, Laganenka L, Sourjik V. Multiple functions of flagellar motility  
761 and chemotaxis in bacterial physiology. *FEMS Microbiol Rev* 2021; **45**: 1–19.
- 762 5. D'Souza GG, Povolo VR, Keegstra JM, Stocker R, Ackermann M. Nutrient  
763 complexity triggers transitions between solitary and colonial growth in bacterial  
764 populations. *ISME J* 2021; **15**: 2614–2626.
- 765 6. Matz C, Jürgens K. High motility reduces grazing mortality of planktonic  
766 bacteria. *Appl Environ Microbiol* 2005; **71**: 921–929.
- 767 7. Matilla MA, Krell T. The effect of bacterial chemotaxis on host infection and  
768 pathogenicity. *FEMS Microbiol Rev* 2018; **42**: 40–67.
- 769 8. Zinicola M, Higgins H, Lima S, Machado V, Guard C, Bicalho R. Shotgun  
770 metagenomic sequencing reveals functional genes and microbiome associated  
771 with bovine digital dermatitis. *PLoS One* 2015; **10**: e0133674.
- 772 9. Raina JB, Fernandez V, Lambert B, Stocker R, Seymour JR. The role of  
773 microbial motility and chemotaxis in symbiosis. *Nature Reviews Microbiology*  
774 2019; **17**: 284–294.
- 775 10. Aschtgen MS, Brennan CA, Nikolakakis K, Cohen S, McFall-Ngai M, Ruby EG.  
776 Insights into flagellar function and mechanism from the squid–vibrio symbiosis.  
777 *npj Biofilms and Microbiomes* 2019; **5**: 1–10.
- 778 11. Neumann S, Grosse K, Sourjik V. Chemotactic signaling via carbohydrate  
779 phosphotransferase systems in *Escherichia coli*. *Proc Natl Acad Sci USA*  
780 2012; **109**: 12159–12164.
- 781 12. Yang Y, Pollard AM, Höfler C, Poschet G, Wirtz M, Hell R, et al. Relation  
782 between chemotaxis and consumption of amino acids in bacteria. *Mol*  
783 *Microbiol* 2015; **96**: 1272–1282.
- 784 13. Yang Q, Defoidt T. Quorum sensing positively regulates flagellar motility in  
785 pathogenic *Vibrio harveyi*. *Environ Microbiol* 2015; **17**: 960–968.



- 786 14. Tso WW, Adler J. Negative Chemotaxis in *Escherichia coli*. *J Bacteriol* 1974;  
787 **118**: 560–576.
- 788 15. Dechesne A, Wang G, Gülez G, Or D, Smets BF. Hydration-controlled  
789 bacterial motility and dispersal on surfaces. *Proc Natl Acad Sci USA* 2010;  
790 **107**: 14369–14372.
- 791 16. Morgan SJ, Chaston JM. Flagellar genes are associated with the colonization  
792 persistence phenotype of the *Drosophila melanogaster* microbiota. *Microbiol*  
793 *Spectr* 2023; **11**: e04585-22.
- 794 17. Schavemaker PE, Lynch M. Flagellar energy costs across the tree of life. *Elife*  
795 2022; **11**: e77266.
- 796 18. Yawata Y, Carrara F, Menolascina F, Stocker R. Constrained optimal foraging  
797 by marine bacterioplankton on particulate organic matter. *Proc Natl Acad Sci*  
798 *USA* 2020; **117**: 25571–25579.
- 799 19. Sathyamoorthy R, Kushmaro Y, Rotem O, Matan O, Kadouri DE, Huppert A, et  
800 al. To hunt or to rest: prey depletion induces a novel starvation survival  
801 strategy in bacterial predators. *ISME J* 2020; **15**: 109–123.
- 802 20. Grognot M, Nam JW, Elson LE, Taute KM. Physiological adaptation in flagellar  
803 architecture improves *Vibrio alginolyticus* chemotaxis in complex  
804 environments. *Proc Natl Acad Sci USA* 2023; **120**: e2301873120.
- 805 21. Barka EA, Vatsa P, Sanchez L, Gaveau-Vaillant N, Jacquard C, Klenk H-P, et  
806 al. Taxonomy, physiology, and natural products of Actinobacteria. *Microbiology*  
807 *and Molecular Biology Reviews* 2016; **80**: 1–43.
- 808 22. Kajikawa A, Midorikawa E, Masuda K, Kondo K, Irisawa T, Igimi S, et al.  
809 Characterization of flagellins isolated from a highly motile strain of  
810 *Lactobacillus agilis*. *BMC Microbiol* 2016; **16**: 1–8.
- 811 23. Liu R, Ochman H. Stepwise formation of the bacterial flagellar system. *Proc*  
812 *Natl Acad Sci USA* 2007; **104**: 7116–7121.
- 813 24. Pallen MJ, Penn CW, Chaudhuri RR. Bacterial flagellar diversity in the post-  
814 genomic era. *Trends Microbiol* 2005; **13**: 143–149.
- 815 25. Girgis HS, Liu Y, Ryu WS, Tavazoie S. A comprehensive genetic  
816 characterization of bacterial motility. *PLoS Genet* 2007; **3**: e154.
- 817 26. Madin JS, Nielsen DA, Brbic M, Corkrey R, Danko D, Edwards K, et al. A  
818 synthesis of bacterial and archaeal phenotypic trait data. *Scientific Data* 2020;  
819 **7**: 1–8.

- 820 27. Parks DH, Chuvochina M, Rinke C, Mussig AJ, Chaumeil PA, Hugenholtz P.  
821 GTDB: an ongoing census of bacterial and archaeal diversity through a  
822 phylogenetically consistent, rank normalized and complete genome-based  
823 taxonomy. *Nucleic Acids Res* 2022; **50**: 785–794.
- 824 28. Grime, J. P. Evidence for the existence of three primary strategies in plants  
825 and its relevance to ecological and evolutionary theory. *The american*  
826 *naturalist* 1977; **111**: 1169-1194.
- 827 29. Piton G, Allison SD, Bahram M, Hildebrand F, Martiny JBH, Treseder KK, et al.  
828 Life history strategies of soil bacterial communities across global terrestrial  
829 biomes. *Nature Microbiology* 2023; **8**: 2093–2102.
- 830 30. Noell SE, Brennan E, Washburn Q, Davis EW, Hellweger FL, Giovannoni SJ.  
831 Differences in the regulatory strategies of marine oligotrophs and copiotrophs  
832 reflect differences in motility. *Environ Microbiol* 2023; **25**: 1265–1280.
- 833 31. Malik AA, Martiny JBH, Brodie EL, Martiny AC, Treseder KK, Allison SD.  
834 Defining trait-based microbial strategies with consequences for soil carbon  
835 cycling under climate change. *ISME J* 2019; **14**: 1–9.
- 836 32. Chesson, P. Mechanisms of maintenance of species diversity. *Annu Rev Ecol*  
837 *Syst* 2000; **31**: 343-366.
- 838 33. Sokol NW, Slessarev E, Marschmann GL, Nicolas A, Blazewicz SJ, Brodie EL,  
839 et al. Life and death in the soil microbiome: how ecological processes  
840 influence biogeochemistry. *Nature Reviews Microbiology* 2022; **20**: 415–430.
- 841 34. Bickel S, Or D. Soil bacterial diversity mediated by microscale aqueous-phase  
842 processes across biomes. *Nature Communications* 2020; **11**: 1–9.
- 843 35. Wang G, Or D. Aqueous films limit bacterial cell motility and colony expansion  
844 on partially saturated rough surfaces. *Environ Microbiol* 2010; **12**: 1363–1373.
- 845 36. Kuzyakov Y, Blagodatskaya E. Microbial hotspots and hot moments in soil:  
846 Concept & review. *Soil Biol Biochem* 2015; **83**: 184–199.
- 847 37. Kuzyakov Y, Friedel JK, Stahr K. Review of mechanisms and quantification of  
848 priming effects. *Soil Biol Biochem* 2000; **32**: 1485–1498.
- 849 38. Wu X, Bei S, Zhou X, Luo Y, He Z, Song C, et al. Metagenomic insights into  
850 genetic factors driving bacterial niche differentiation between bulk and  
851 rhizosphere soils. *Science of the Total Environment* 2023; **891**: 164221.
- 852 39. Feng H, Fu R, Hou X, Lv Y, Zhang N, Liu Y, et al. Chemotaxis of Beneficial  
853 Rhizobacteria to Root Exudates: The first step towards root–microbe  
854 rhizosphere interactions. *International Journal of Molecular Sciences* 2021; **22**:  
855 6655.

- 856 40. Li C, Liao H, Xu L, Wang C, He N, Wang J, et al. The adjustment of life history  
857 strategies drives the ecological adaptations of soil microbiota to aridity. *Mol*  
858 *Ecol* 2022; **31**: 2920–2934.
- 859 41. Chen Y, Neilson JW, Kushwaha P, Maier RM, Barberán A. Life-history  
860 strategies of soil microbial communities in an arid ecosystem. *ISME J* 2021;  
861 **15**: 649–657.
- 862 42. Fernández-Gómez B, Richter M, Schüller M, Pinhassi J, Acinas SG, González  
863 JM, et al. Ecology of marine Bacteroidetes: a comparative genomics approach.  
864 *ISME J* 2013; **7**: 1026–1037.
- 865 43. Mark BM, Zhu Y. Gliding motility and por secretion system genes are  
866 widespread among members of the phylum bacteroidetes. *J Bacteriol* 2013;  
867 **195**: 270–278.
- 868 44. Albright S, Louca S. Trait biases in microbial reference genomes. *Scientific*  
869 *Data* 2023; **10**: 1–17.
- 870 45. Vieira-Silva S, Rocha EPC. The systemic imprint of growth and its uses in  
871 ecological (meta)genomics. *PLoS Genet* 2010; **6**: e1000808.
- 872 46. Barnett SE, Egan R, Foster B, Eloe-Fadrosh EA, Buckley DH. Genomic  
873 features predict bacterial life history strategies in soil, as identified by  
874 metagenomic stable isotope probing. *mBio* 2023; **14**: e03584-22.
- 875 47. Fierer N, Barberán A, Laughlin DC. Seeing the forest for the genes: Using  
876 metagenomics to infer the aggregated traits of microbial communities. *Front*  
877 *Microbiol* 2014; **5**: 119845.
- 878 48. Dar D, Thomashow LS, Weller DM, Newman DK. Global landscape of  
879 phenazine biosynthesis and biodegradation reveals species-specific  
880 colonization patterns in agricultural soils and crop microbiomes. *Elife* 2020; **9**:  
881 1–44.
- 882 49. Spohn M, Klaus K, Wanek W, Richter A. Microbial carbon use efficiency and  
883 biomass turnover times depending on soil depth – Implications for carbon  
884 cycling. *Soil Biol Biochem* 2016; **96**: 74–81.
- 885 50. Fontaine S, Barot S, Barré P, Bdioui N, Mary B, Rumpel C. Stability of organic  
886 carbon in deep soil layers controlled by fresh carbon supply. *Nature* 2007; **450**:  
887 277–280.
- 888 51. Brewer, T. E., Aronson, E. L., Arogyaswamy, K., Billings, S. A., Botthoff, J. K.,  
889 Campbell, A. N., et al. Ecological and genomic attributes of novel bacterial  
890 taxa that thrive in subsurface soil horizons. *mBio* 2019, **10**: 10-1128.

- 891 52. Raich JW, Schlesinger WH. The global carbon dioxide flux in soil respiration  
892 and its relationship to vegetation and climate. *Tellus B* 1992; **44**: 81–99.
- 893 53. Sokol NW, Kuebbing SE, Karlsen-Ayala E, Bradford MA. Evidence for the  
894 primacy of living root inputs, not root or shoot litter, in forming soil organic  
895 carbon. *New Phytologist* 2019; **221**: 233–246.
- 896 54. Cheng W, Johnson DW, Fu S. Rhizosphere effects on decomposition. *Soil  
897 Science Society of America Journal* 2003; **67**: 1418–1427.
- 898 55. Xu J, Zhang Y, Zhang P, Trivedi P, Riera N, Wang Y, et al. The structure and  
899 function of the global citrus rhizosphere microbiome. *Nature Communications*  
900 2018; **9**: 1–10.
- 901 56. Fan K, Holland-Moritz H, Walsh C, Guo X, Wang D, Bai Y, et al. Identification  
902 of the rhizosphere microbes that actively consume plant-derived carbon. *Soil  
903 Biol Biochem* 2022; **166**: 108577.
- 904 57. Lucas JM, McBride SG, Strickland MS. Trophic level mediates soil microbial  
905 community composition and function. *Soil Biol Biochem* 2020; **143**: 107756.
- 906 58. Lennon JT, Aanderud ZT, Lehmkuhl BK, Schoolmaster DR. Mapping the niche  
907 space of soil microorganisms using taxonomy and traits. *Ecology* 2012; **93**:  
908 1867–1879.
- 909 59. Sauer DB, Wang DN. Predicting the optimal growth temperatures of  
910 prokaryotes using only genome derived features. *Bioinformatics* 2019; **35**:  
911 3224–3231.
- 912 60. Haiko J, Westerlund-Wikström B. The role of the bacterial flagellum in  
913 adhesion and virulence. *Biology* 2013; **2**: 1242–1267.
- 914 61. Scharf BE, Hynes MF, Alexandre GM. Chemotaxis signaling systems in model  
915 beneficial plant–bacteria associations. *Plant Mol Biol* 2016; **90**: 549–559.
- 916 62. Lemanceau P, Blouin M, Muller D, Moënne-Loccoz Y. Let the core microbiota  
917 be functional. *Trends Plant Sci* 2017; **22**: 583–595.
- 918 63. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM:  
919 assessing the quality of microbial genomes recovered from isolates, single  
920 cells, and metagenomes. *Genome Res* 2015; **25**: 1043–1055.
- 921 64. Orakov A, Fullam A, Coelho LP, Khedkar S, Szklarczyk D, Mende DR, et al.  
922 GUNC: detection of chimerism and contamination in prokaryotic genomes.  
923 *Genome Biol* 2021; **22**: 1–19.

- 924 65. Hyatt D, Chen GL, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal:  
925 Prokaryotic gene recognition and translation initiation site identification. *BMC*  
926 *Bioinformatics* 2010; **11**: 1–11.
- 927 66. Finn RD, Bateman A, Clements J, Coghill P, Eberhardt RY, Eddy SR, et al.  
928 Pfam: the protein families database. *Nucleic Acids Res* 2014; **42**: 222–230.
- 929 67. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence  
930 similarity searching. *Nucleic Acids Res* 2011; **39**: 29–37.
- 931 68. Coloma-Rivero RF, Flores-Concha M, Molina RE, Soto-Shara R, Cartes Á,  
932 Oñate ÁA. Brucella and its hidden flagellar system. *Microorganisms* 2021; **10**:  
933 83.
- 934 69. Chen T, Guestrin C. XGBoost: a scalable tree boosting system. *Proceedings*  
935 *of the 22nd ACM SIGKDD International Conference on Knowledge Discovery*  
936 *and Data Mining* 2016.
- 937 70. Snoek J, Larochelle H, Adams RP. Practical bayesian optimization of machine  
938 learning algorithms. *Adv Neural Inf Process Syst* 2012; **4**: 2951–2959.
- 939 71. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v5: an online tool for  
940 phylogenetic tree display and annotation. *Nucleic Acids Res* 2021; **49**: 293–  
941 296.
- 942 72. Fritz SA, Purvis A. Selectivity in mammalian extinction risk and threat types: a  
943 new measure of phylogenetic signal strength in binary traits. *Conservation*  
944 *Biology* 2010; **24**: 1042–1051.
- 945 73. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and  
946 Evolution in R language. *Bioinformatics* 2004; **20**: 289–290.
- 947 74. Galperin MY, Wolf YI, Makarova KS, Alvarez RV, Landsman D, Koonin E V.  
948 COG database update: focus on microbial diversity, model organisms, and  
949 widespread pathogens. *Nucleic Acids Res* 2021; **49**: 274–281.
- 950 75. Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J.  
951 eggNOG-mapper v2: Functional annotation, orthology assignments, and  
952 domain prediction at the metagenomic scale. *Mol Biol Evol* 2021; **38**: 5825–  
953 5829.
- 954 76. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read  
955 simulator. *Bioinformatics* 2012; **28**: 593–594.
- 956 77. Rodriguez-R LM, Konstantinidis KT. Estimating coverage in metagenomic data  
957 sets and why it matters. *ISME J* 2014; **8**: 2349–2351.

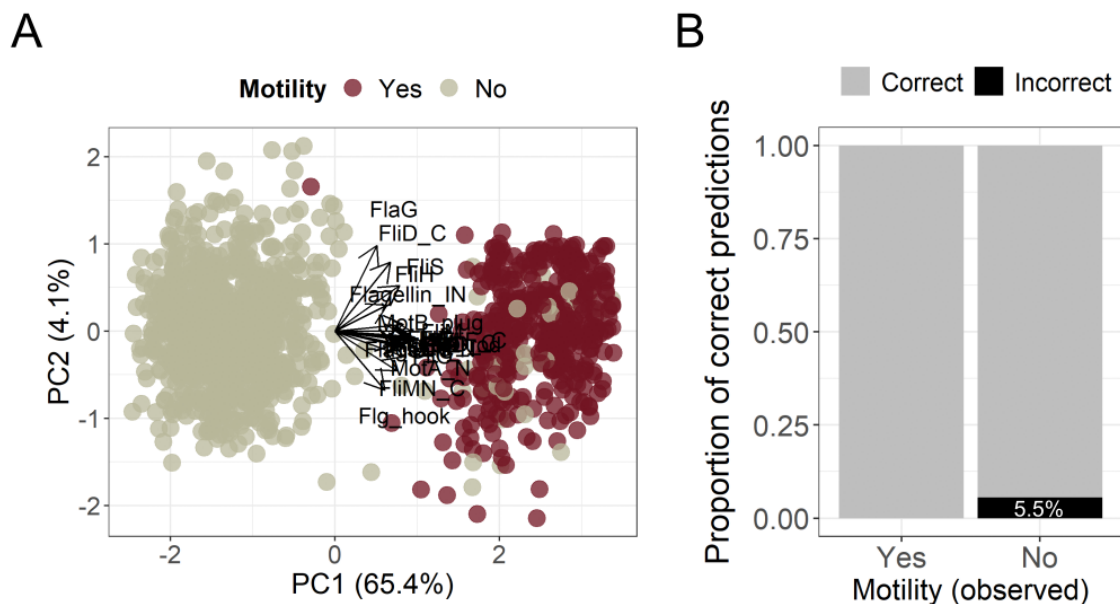
- 958 78. Buchfink B, Reuter K, Drost HG. Sensitive protein alignments at tree-of-life  
959 scale using DIAMOND. *Nature Methods* 2021; **18**: 366–368.
- 960 79. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al.  
961 BLAST+: Architecture and applications. *BMC Bioinformatics* 2009; **10**: 1–9.
- 962 80. Manor O, Borenstein E. MUSiCC: A marker genes based framework for  
963 metagenomic normalization and accurate profiling of gene abundances in the  
964 microbiome. *Genome Biol* 2015; **16**: 1–20.
- 965 81. Holland-Moritz H, Vanni C, Fernandez-Guerra A, Bissett A, Fierer N. An  
966 ecological perspective on microbial genes of unknown function in soil. *bioRxiv*  
967 2021; 2021.12.02.470747.
- 968 82. Mendes R, Garbeva P, Raaijmakers JM. The rhizosphere microbiome:  
969 significance of plant beneficial, plant pathogenic, and human pathogenic  
970 microorganisms. *FEMS Microbiol Rev* 2013; **37**: 634–663.
- 971 83. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina  
972 sequence data. *Bioinformatics* 2014; **30**: 2114–2120.
- 973 84. Team R Core. R: A language and environment for statistical computing.  
974 *Vienna, Austria* 2018; **0**.
- 975 85. Lin H, Peddada S Das. Analysis of compositions of microbiomes with bias  
976 correction. *Nature Communications* 2020; **11**: 1–11.
- 977 86. McMurdie PJ, Holmes S. phyloseq: An R package for reproducible interactive  
978 analysis and graphics of microbiome census data. *PLoS One* 2013; **8**: e61217.
- 979
- 980
- 981
- 982
- 983

984 **Supplementary Material**

985

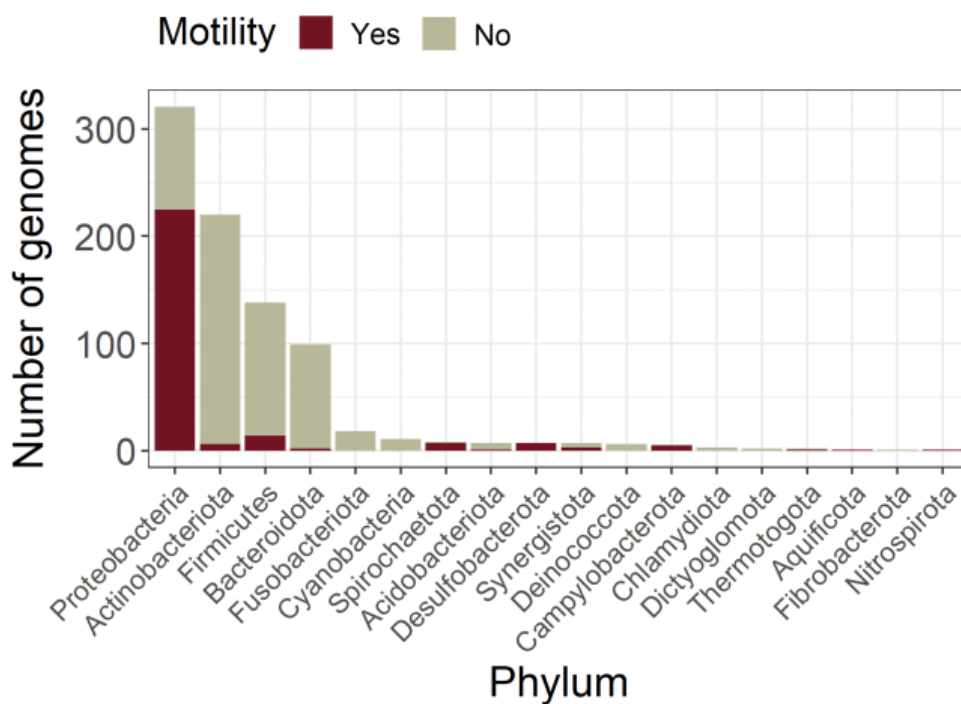
986 **Supplementary Figure 1. Prediction of the capacity to display flagellar motility**  
987 **in bacterial taxa based on presence/absence information of 21 genes involved**  
988 **in flagellar assembly.** A. Principal Components Analysis (PCA) based on the  
989 presence/absence of 21 flagellar genes in genomes of taxa that have been  
990 empirically found to be flagellated (N = 388 genomes) or non-flagellated (N = 837  
991 genomes). Empirical information on flagellar motility was obtained from the bacterial  
992 phenotypic trait data compiled in [26]. We only included genomes that were 100%  
993 complete, contained an assembled 16S rRNA gene, and showed no signs of  
994 chimerism. B. Accuracy of a boosted regression machine learning model trained on  
995 the genomes shown in panel A for the prediction of the capacity for flagellar motility  
996 in any given bacterial genome based on the presence/absence of 21 flagellar genes.  
997 Accuracy was tested on 30% of the original genome set (116 genomes from  
998 flagellated taxa and 251 genomes from non-flagellated taxa). Genomes were  
999 obtained from the Genome Taxonomy Database (GTDB r207; [27]).

1000



1001

1002 **Supplementary Figure 2. Taxonomic distribution of genomes with empirically**  
1003 **determined capacity for flagellar motility that were used as training data for a**  
1004 **boosted regression machine learning model to predict the capacity for**  
1005 **flagellar motility based on the presence/absence of 21 flagellar genes.** Flagellar  
1006 motility information was obtained from the bacterial phenotypic trait data compiled in  
1007 [26].  $N_{\text{Training set}} = 858$  genomes,  $N_{\text{Full set}} = 1225$ .



1008

1009

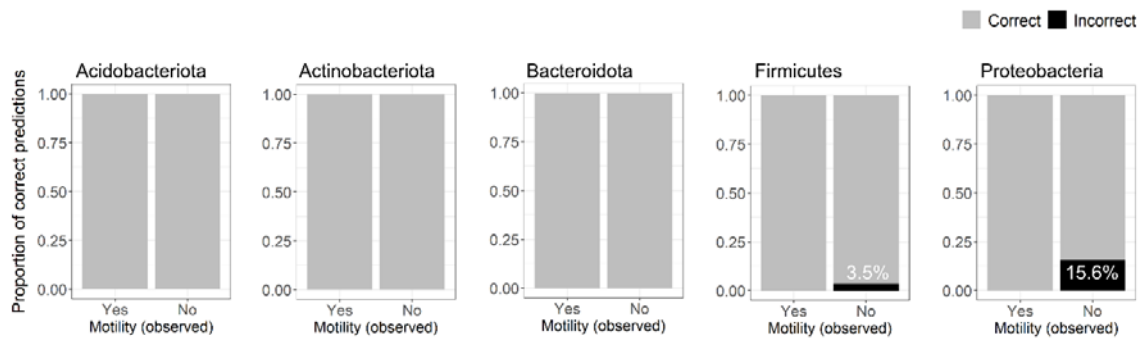
1010

1011



1012 **Supplementary Figure 3. Predictive accuracy across phyla of a boosted**  
1013 **regression machine learning model for the prediction of the capacity for**  
1014 **flagellar motility in any given bacterial genome based on the**  
1015 **presence/absence of 21 flagellar genes.** Accuracy was tested on 30% of the  
1016 original genome set (116 genomes from flagellated taxa and 251 genomes from non-  
1017 flagellated taxa). Genomes were obtained from the Genome Taxonomy Database  
1018 (GTDB r207; [27]).  $N_{\text{Acidobacteriota}} = 8$ ,  $N_{\text{Actinobacteriota}} = 87$ ,  $N_{\text{Bacteroidota}} = 52$ ,  $N_{\text{Firmicutes}} =$   
1019  $66$ ,  $N_{\text{Proteobacteria}} = 123$ .

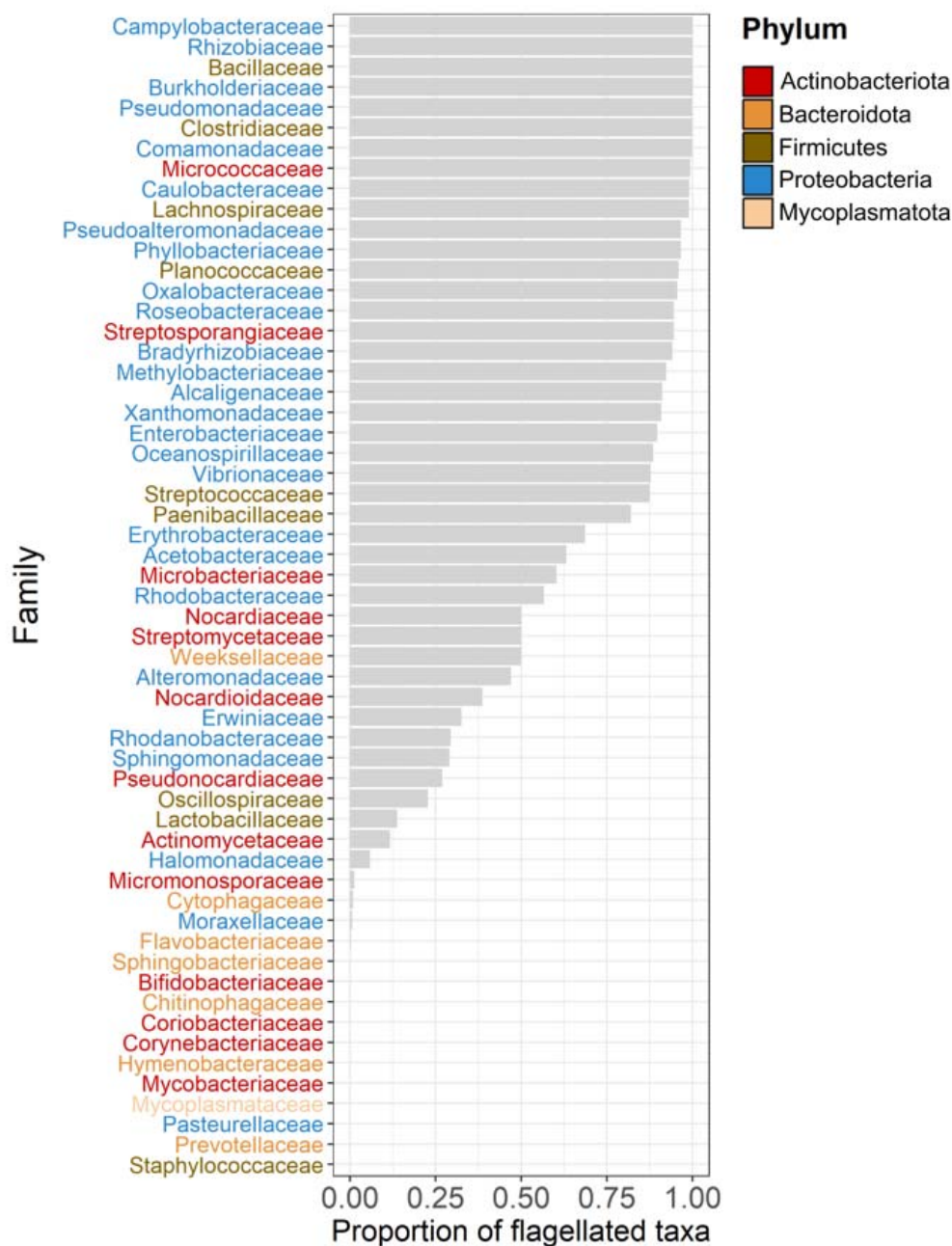
1020



1021

1022

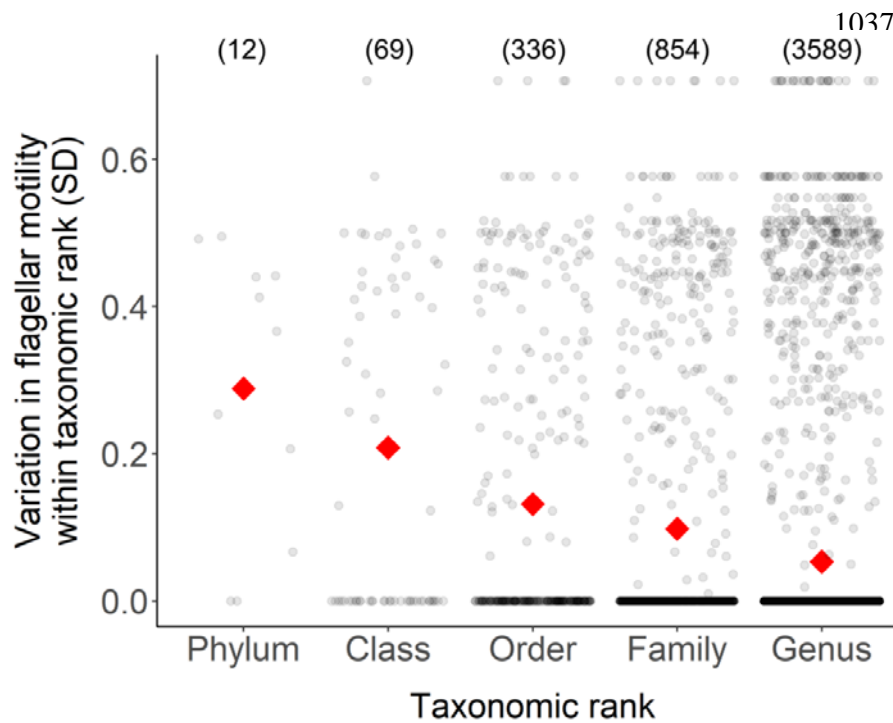
1023 **Supplementary Figure 4. Prevalence of flagellar motility across bacterial**  
1024 **families containing more than 100 high-quality genomes in the Genome**  
1025 **Taxonomy Database (GTDB r207; [27]).** We only included genomes that were  
1026 >95% complete, contained an assembled 16S rRNA gene, and showed no signs of  
1027 chimerism. N = 23,256 genomes.



1028

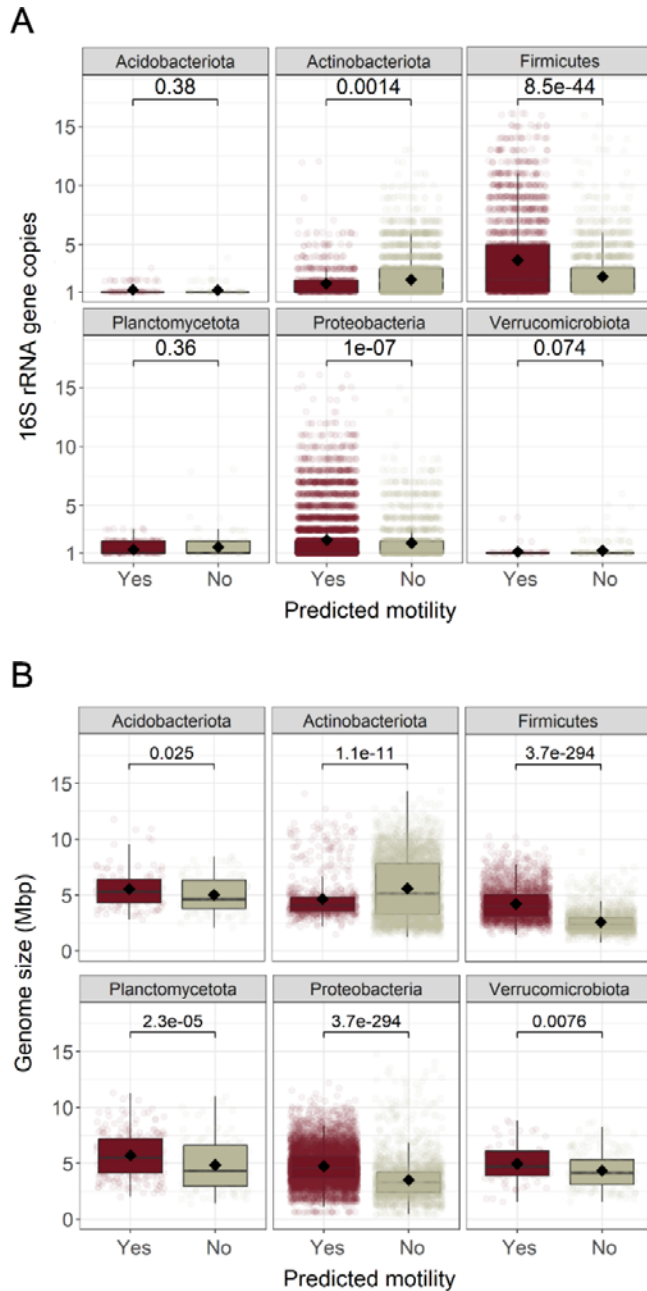
1029 **Supplementary Figure 5. Variation in flagellar motility status across taxonomic**  
1030 **ranks.** A measure of variation in the flagellar motility status of taxa belonging to  
1031 different taxonomic ranks was obtained from the standard deviation (SD) of their  
1032 flagellar motility status (1, flagellated; 0, non-flagellated). Numbers in brackets  
1033 indicate the total number of unique taxa within each of the taxonomic ranks. Red  
1034 diamonds indicate the mean of the standard deviation of the flagellar motility status  
1035 within each taxonomic rank.

1036



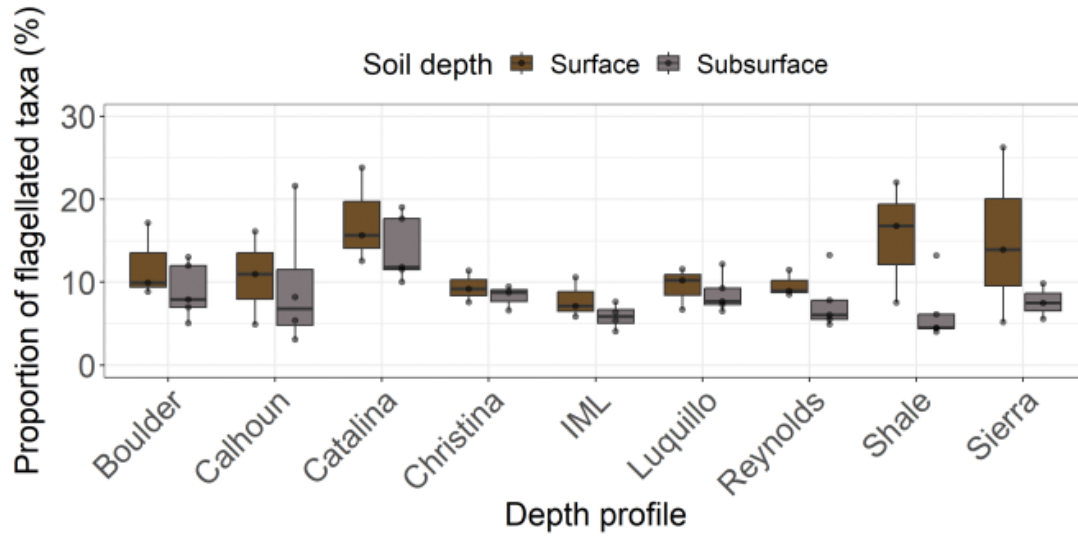
1038 **Supplementary Figure 6. Distribution of the total number of 16S rRNA gene**  
1039 **copies per genome and genome size across the 6 phyla with even proportions**  
1040 **of taxa predicted to be flagellated and non-flagellated.** A. Number of 16S rRNA  
1041 gene copies in genomes of taxa predicted to be flagellated and non-flagellated. B.  
1042 Genome size of taxa predicted to be flagellated and non-flagellated. Statistical  
1043 significance was obtained from Mann-Whitney U tests ( $P < 0.05$ ),  $N = 21,551$ .

1044

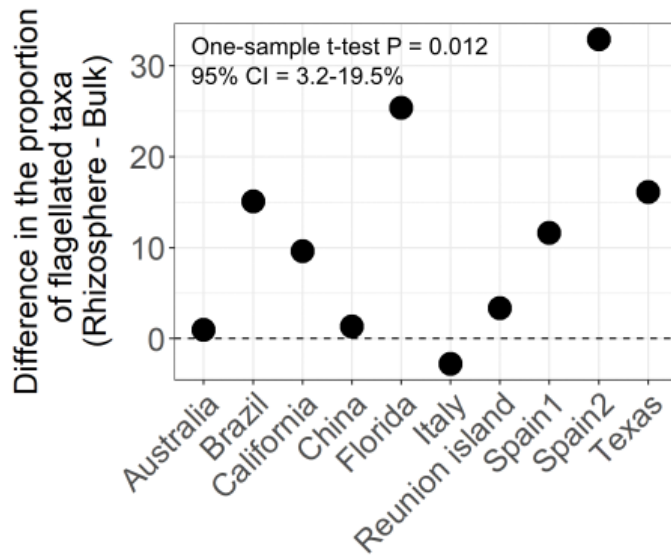


1045 **Supplementary Figure 7. Estimated prevalence of flagellar motility in bacterial**  
1046 **communities from 9 soil depth profiles collected across the USA (Surface, 0-**  
1047 **20cm; Subsurface, 20-90cm, N = 66; [51]).**

1048



1049 **Supplementary Figure 8. Difference in the prevalence of flagellar motility in**  
1050 **rhizosphere and bulk soil bacterial communities collected from citrus species**  
1051 **across the globe (N = 10; [55]).**



1052

1053

1054

1055 **Supplementary Figure 9. Taxonomic composition of the Amplicon Sequence**  
1056 **Variants (ASVs) that responded to glucose amendment in soil.** Bacterial  
1057 communities from a 117 day soil incubation experiment with daily glucose  
1058 amendment were characterized using amplicon sequencing of the 16S rRNA gene  
1059 [57]. ASVs were considered responsive to glucose amendment based on a  
1060 differential abundance analysis comparing bacterial communities from soils  
1061 amended with glucose versus communities from soils that did not receive any  
1062 external carbon inputs (N = 28 responsive ASVs).

1063

