

1 Microbial regime changes and indicators of eutrophication on the Mississippi  
2 River identified via a human-powered 2900 km transect

3  
4 Michael W. Henson<sup>1</sup>, Jordan Hanssen<sup>2</sup>, Greg Spooner<sup>2</sup>, Patrick Fleming<sup>2</sup>,  
5 Markus Pukonen<sup>2</sup>, Frederick Stahr<sup>3</sup>, and J. Cameron Thrash<sup>1\*</sup>

6  
7 <sup>1</sup> Department of Biological Sciences, Louisiana State University, Baton Rouge, LA  
8 70803, U.S.A.

9 <sup>2</sup> O.A.R. Northwest, Seattle, WA 98103, U.S.A.

10 <sup>3</sup> School of Oceanography, University of Washington, Seattle, WA 98195, U.S.A.

11  
12  
13  
14 \*Correspondence: thrashc@lsu.edu  
15  
16  
17  
18  
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

45 **Abstract**

46

47 Draining 31 states and roughly 3 million km<sup>2</sup>, the Mississippi River (MSR) and its  
48 tributaries constitute an essential resource to millions of people for clean drinking  
49 water, transportation, agriculture, and industry. Since the turn of the 20<sup>th</sup> century,  
50 MSR water quality has continually rated poorly due to human activity. Acting as  
51 first responders, microorganisms can mitigate, exacerbate, and/or serve as  
52 predictors for water quality, yet we know little about their community structure or  
53 ecology at the whole river scale for large rivers. We collected both biological (16S  
54 and 18S rRNA gene amplicons) and physicochemical samples from 38 MSR  
55 sites over nearly 3000 km from Minnesota to the Gulf of Mexico. These data  
56 represent the first of their kind for a top ten river in size, volume, and/or drainage  
57 and revealed distinct regime changes between upper and lower MSR microbial  
58 communities that corresponded to Strahler's river order and nutrient  
59 concentrations. Within these assemblages, we identified subgroups of OTUs  
60 from the phyla Acidobacteria, Bacteroidetes, Fungi, and Heterokonts that highly  
61 correlated with, and were predictive of, the important eutrophication nutrients  
62 nitrate and phosphate. This study offers the most comprehensive view of  
63 Mississippi River microbiota to date, establishes the groundwork for future  
64 temporal and spatial studies of river perturbations, and provides potential  
65 microbial indicators of river health related to eutrophication.

66

67

68 By connecting terrestrial, lotic, and marine systems, rivers perform vital  
69 roles in both the transport and processing of compounds in all major global  
70 biogeochemical cycles<sup>1-5</sup>. Within the carbon cycle alone, rivers collectively  
71 discharge organic carbon to the oceans at over 0.4 Pg C yr<sup>-1</sup><sup>6</sup>. Perhaps more  
72 importantly, rivers are generally net heterotrophic<sup>7</sup>, indicating that they not only  
73 transport organic matter but host active metabolic processing of it as well.  
74 Conservative estimates place heterotrophic output of the world's fluvial networks  
75 (streams, rivers, and estuaries) at 0.32 Pg C yr<sup>-1</sup><sup>5,8</sup>. Although rivers contain a  
76 small minority of global fresh water at any given moment, the considerable  
77 volumes that pass through these systems make them relevant to models  
78 attempting to quantify global elemental transformations. However, the  
79 fundamental engines of these transformations- microorganisms- have received  
80 comparatively little study in rivers relative to other aquatic systems, despite the  
81 fact that microbial functions likely play a vital role in ecosystem health for both  
82 rivers themselves and their places of discharge.

83 At 3734 km, the Mississippi River (MSR) is the fourth longest on earth,  
84 draining 31 U.S. states and two Canadian provinces- a watershed consisting of  
85 41% of the continental U.S.<sup>9,10</sup> The MSR is a major source of drinking water for  
86 many U.S. cities; a critical thoroughfare for transportation, commerce, industry,  
87 agriculture, and recreation; and conveys the vestiges of human activity to the  
88 Gulf of Mexico (GOM). In New Orleans, the average flow rate is over 600,000  
89 cubic feet s<sup>-1</sup> (cfs)<sup>11</sup>, but can exceed 3 million cfs during flood stages<sup>12</sup>, carrying  
90 over 150 x 10<sup>9</sup> kg of suspended sediment into the northern GOM annually<sup>9,13</sup>.  
91 This massive discharge includes excess nutrients (nitrogen and phosphorus),  
92 primarily from agricultural runoff<sup>14-18</sup>, and fuels one of the largest marine zones  
93 of seasonal hypoxia in the world<sup>19-22</sup>. Understanding microbial relationships to  
94 river eutrophication will inform hypotheses regarding their contributions to either  
95 mitigating or exacerbating nutrient input.

96 Far from a homogenous jumble of organisms ferried downriver, microbial  
97 community composition changes with distance from the river mouth and/or from  
98 the influence of tributaries<sup>23-25</sup>, attributable to changing nutrient concentrations<sup>26-</sup>  
99 <sup>28</sup>, dendritic length<sup>4,29</sup>, differing dissolved organic matter (DOM) sources<sup>26,30,31</sup>,  
100 and land use changes<sup>16,27,28,32</sup>. Past studies of the Thames, Danube, Yenisei,  
101 and Columbia Rivers have found that planktonic river microbiota were dominated  
102 by the phyla Actinobacteria, Proteobacteria, and Bacteroidetes, specifically, taxa  
103 such as- hgcl/cal Actinobacteria, *Polynucleobacter* spp., GKS9, and LD28  
104 *Betaproteobacteria*, CL500-29 Bacteroidetes, LD12 freshwater SAR11  
105 *Alphaproteobacteria*, and *Novosphingobium* spp.<sup>4,33,34</sup>. More recent 16S rRNA  
106 gene amplicon and metagenomic studies of the Minnesota portion of the MSR  
107 corroborated previous research in other rivers that identified an increased  
108 proportion<sup>4</sup> or richness<sup>29</sup> of freshwater taxa with river distance, and an increased  
109 abundance of "core" river taxa<sup>35</sup> with cumulative residence time<sup>44,22,27,33, 52</sup>.

110 Researchers have suggested that these patterns supported application of  
111 the River Continuum Concept (RCC)<sup>10</sup> to river microbiota. The concept postulates  
112 that as a river increases in size, the influences of riparian and other inputs will  
113 decrease as the river establishes a dominant core community<sup>36</sup>, and richness will  
114 increase from headwaters to mid-order before decreasing in higher order rivers<sup>36</sup>.  
115 Therefore, as continuous systems with increasing volumes and residence times,  
116 river microbiota should transition from experiencing strong influences of mass  
117 effects from terrestrial and tributary sources to systems where species sorting  
118 plays a more important role<sup>4,37,38</sup>. Complicating matters, particle-associated  
119 communities in rivers (frequently defined as those found on filters of > ~3  $\mu\text{m}$ )  
120 remain distinct from their free-living counterparts<sup>37-39, 53</sup>, potentially due to  
121 increased production rates from readily obtainable carbon<sup>38, 53</sup>. Typical taxa  
122 associated with particles include OTUs related to the Bacteroidetes clades  
123 *Flavobacteria* and *Cytophaga*, Planctomycetes, *Rhizobium* spp., and  
124 *Methylophilaceae* spp<sup>39,42-44</sup>. However, consistent trends in particle community  
125 composition are murky, with recent evidence suggesting organisms may switch  
126 between free-living and particle-associated lifestyles depending on substrate  
127 availability and chemical queues<sup>42,45</sup>. Thus, rivers constitute complex and highly  
128 dynamic ecosystems from a metacommunity perspective.

129 However, our knowledge of river microbial assemblages, their ecology,  
130 dispersal, and their relationship to chemical constituents remains in comparative  
131 infancy to that of lakes and oceans. Although microbes play a central role in the  
132 RCC<sup>36</sup>, microbiologically-oriented transects at the whole-river or -catchment scale  
133 have only been conducted for a handful of systems<sup>4,28,29,46,47</sup>, and until this work,  
134 none had been attempted for any of the largest rivers in the world. Furthermore,  
135 little data exists on microbial eukaryotic communities in rivers. During the fall of  
136 2014, we completed the most extensive microbiological survey of the Mississippi  
137 River to date with a continual rowed transect over two months. Rowers from the  
138 adventure education non-profit OAR Northwest collected samples from  
139 Minneapolis, MN to the Gulf of Mexico (2918 km) (Fig. 1A). They also maintained  
140 active blogging and social media content and visited 21 schools along the river  
141 that incorporated elements of the journey into their curriculum. Our findings  
142 greatly expand the current information available on microbial assemblages in  
143 major lotic ecosystems and help further delineate the relationships between  
144 microbial structure and stream order, nutrients, and volume.

145

## 146 **Results**

147 Using rowboats and simple syringe-based filtration protocol, we measured 12  
148 biological, chemical, and physical parameters (e.g. 16S and 18S rRNA gene  
149 communities,  $\text{NH}_4^+$ , river speed, etc.) from 38 sites along a 2918 km transect of  
150 the MSR (Fig. 1A). River order increases dramatically at the Missouri confluence  
151 (eighth to tenth Strahler order<sup>48</sup>), which corresponded to overall discharge (Fig.  
152 1A) and beta diversity changes discussed below, and thus we used this juncture  
153 to separate the upper MSR (0 km – 1042 km, Sites A-S) and lower MSR (1075-

154 2914 km, Sites T-AI). Within the upper MSR,  $\text{NO}_3^-$ ,  $\text{PO}_4^{3-}$ , and  $\text{NO}_2^-$  were variable  
155 but generally increased downriver until peak concentrations near the confluences  
156 of the Illinois and Missouri Rivers. This gave way to lower and more consistent  
157 concentrations along the lower MSR (Fig. 1B). Ammonium showed much greater  
158 variability along the transect. Turbidity (inversely related to secchi disk visibility)  
159 increased steadily downriver to a maximum near the Illinois and Missouri River  
160 confluences (1042 km, Site S) (Fig. 1B), then trended downwards for the rest of  
161 the transect. Planktonic ( $< 2.7 \mu\text{m}$ ) cell counts varied between 1 and  $3 \times 10^6$   
162 cells/mL in the upper MSR, and decreased to high  $10^5$  cells/mL in the lower MSR  
163 (Fig. 1B). Water temperature ranged from 19°C (133km, Site E) to 11.7°C (2552  
164 km, Site Ag), and river speed, excluding three sites sampled from shore, was  
165 between 5.5 mph at Site Y and 0.4 mph (597 km, Site L) (Table S1). Spearman  
166 rank correlations of the measured environmental parameters showed strong  
167 positive correlations between nitrate, phosphate, distance, and increased  
168 turbidity; while nitrate and phosphate both strongly correlated negatively to water  
169 temperature and river speed (Table S1).

170

#### 171 *Bacterial and archaeal communities*

172 We observed a clear distinction between the 0.2-2.7  $\mu\text{m}$  and  $> 2.7 \mu\text{m}$  16S rRNA  
173 gene communities (ANOSIM  $R = 0.65$ ,  $P = 0.001$ ) (Fig. S1A). Both size fractions  
174 had comparable species richness and evenness values that trended upwards  
175 downriver (Figs. S3A-B), although an earlier peak occurred for both at sites O-Q  
176 (761-999 km) below the Des Moines River and above the Illinois River (Figs.  
177 S3A-B). Both size fractions (stress = 0.14 for each) also showed a significant  
178 separation between sites above and below the Missouri River confluence  
179 (ANOSIM,  $> 2.7 \mu\text{m}$ :  $R = 0.44$ ,  $P = 0.001$ ; 0.2-2.7  $\mu\text{m}$ :  $R = 0.48$ ,  $P = 0.001$ ) (Figs.  
180 2A and B), that the Eukaryotic fractions mirrored (below). The environmental  
181 variables phosphate and turbidity had the highest correlation with the separation  
182 between upper and lower  $> 2.7 \mu\text{m}$  communities ( $r = 0.57$ ,  $r = 0.54$ , respectively),  
183 with water temperature and distance ( $r > 0.40$ ) also contributing (Fig. 2A). At an  
184 OTU level, taxa related to the hgcl/cal clade (Actinobacteria) and unclassified  
185 *Bacillaceae* ( $r > 0.77$ ,  $P = 0.001$ ) contributed most to the separation between the  
186 upper and lower  $> 2.7 \mu\text{m}$  communities, with OTUs related to the *Bacillales*,  
187 *Gemmatimonadaceae*, *Peptococcaceae*, and *Micromonosporaceae* clades also  
188 a factor ( $r > 0.70$ ,  $P = 0.001$ ) (Figure 2A, Table S1). For the 0.2-2.7  $\mu\text{m}$  fraction,  
189 distance and nitrate most strongly correlated with the distinction between upper  
190 and lower communities among environmental factors ( $r = 0.59$  and  $r = 0.47$ ,  
191 respectively), although phosphate, turbidity, and water temperature ( $r > 0.40$  for  
192 each) also contributed (Fig. 2B). OTUs related to *Flavobacterium* and a  
193 unclassified Bacterium (closest NCBI BLAST hit *Acidovorax* sp., KM047473),  
194 most strongly contributed to the separation between the 0.2-2.7  $\mu\text{m}$  communities  
195 ( $r > 0.52$ ,  $P = 0.001$ ). Other important OTUs belonged to the clades  
196 Bacteroidetes, *Microbacteriaceae*, *Clostridiales*, and *Holophagaceae* ( $r > 0.49$ ,  $P$   
197 = 0.001) (Figure 2B, Table S1).

198 At the phylum level, Proteobacteria, Actinobacteria, and Bacteroidetes  
199 dominated bacterial communities in both fractions (Figs. 3A and B) (Table S1).  
200 Proteobacteria in the  $> 2.7 \mu\text{m}$  fraction showed wide fluctuations in abundance  
201 (Fig. 3A), whereas their  $0.2\text{-}2.7 \mu\text{m}$  counterparts generally increased in relative  
202 abundance downriver (Fig. 3B).  $0.2\text{-}2.7 \mu\text{m}$  Bacteroidetes and Actinobacteria  
203 generally decreased in the upper river and stabilized in the lower river, but  
204 showed considerable variation in abundance in the larger fraction. Cyanobacteria  
205 in the  $> 2.7 \mu\text{m}$  fraction strongly and negatively correlated with increased turbidity  
206 (Spearman rank = 0.67). Both  $> 2.7 \mu\text{m}$  and  $0.2\text{-}2.7 \mu\text{m}$  Acidobacteria had a  
207 strong positive correlation with river distance (Wilcoxon single ranked test,  $P = <$   
208 0.01) (Fig. 3A and B). Within the  $0.2\text{-}2.7 \mu\text{m}$  fraction, the five most abundant  
209 OTUs were classified as LD12 (OTU11), two hgcl clade OTUs (OTU4, OTU7),  
210 *Limnohabitans* sp. (OTU2), and LD28 (OTU8) (Table S1). Comparatively, an  
211 unclassified *Methylophilaceae* (OTU1), *Planktothrix* sp. (OTU12), NS11-12  
212 marine group (OTU21), *Aquirestis* sp. (OTU17), and unclassified  
213 *Sphingobacteriales* (OTU25) were the most abundant OTUs in the  $> 2.7 \mu\text{m}$   
214 (Table S1). Archaea occurred at much lower relative abundances: we found only  
215 eight OTUs belonging to the Euryarchaeota and Thaumarchaeota.  $0.2\text{-}2.7 \mu\text{m}$   
216 Thaumarchaeota increased in abundance along the transect (Fig. 3B), but at less  
217 so than those in the  $> 2.7 \mu\text{m}$  fraction (Fig. 3A). In both fractions, we only  
218 detected Euryarchaeota at very low abundances.

219 We defined the core microbiome as those OTUs detectable after  
220 normalization in greater than 90% of the samples. The  $> 2.7 \mu\text{m}$  and  $0.2\text{-}2.7 \mu\text{m}$   
221 core microbiomes consisted of 95 and 106 OTUs, respectively, classified into  
222 eight different phyla- Proteobacteria, Actinobacteria, Bacteroidetes,  
223 Cyanobacteria, Verrucomicrobia, Chloroflexi, Chlorobi, Gemmatimonadetes- and  
224 some remained unclassified (Table S1). Core microbiome relative abundance in  
225 both fractions decreased along the upper river but stabilized in the lower river  
226 (Fig. 4A). We confirmed this effect by analyzing the upper and lower core  
227 microbiomes separately. Although the total OTU numbers changed (81 and 116  
228 OTUs in the upper MSR and 160 and 144 OTUs in the lower MSR for the  $> 2.7$   
229  $\mu\text{m}$  and  $0.2\text{-}2.7 \mu\text{m}$  fractions, respectively), the trends remained the same (Fig.  
230 S4).

231

### 232 *16S rRNA gene environmental ontology*

233 We successfully classified 313 of the 945 OTUs with EnvO terminology (Table  
234 S1). Of those, freshwater organisms dominated, although their relative  
235 abundance in both the  $> 2.7 \mu\text{m}$  and  $0.2\text{-}2.7 \mu\text{m}$  fractions decreased with river  
236 distance before stabilizing in the lower MSR (Fig. S5A-B). However LD12, the  
237 most abundant OTU in our dataset and a well-established freshwater organism,  
238 did not receive an EnvO classification at all, indicating the limitations of this  
239 technique with current database annotations. Terrestrial organisms from the  $>$   
240  $2.7 \mu\text{m}$  fraction decreased along the transect (Fig. S5A), while the  $0.2\text{-}2.7 \mu\text{m}$   
241 fraction remained stable (Fig. S5B). Although representing a minor fraction of

242 total OTUs, sediment-associated microorganisms remained steady along the  
243 river in both fractions.. Taxa associated with anthropogenic sources in both  
244 fractions increased along the river(Fig. S5A-B).

245

#### 246 *Microbial eukaryotic communities*

247 Eukaryotic communities, observed via the 18S rRNA gene, also showed a  
248 significant separation between  $> 2.7 \mu\text{m}$  and  $0.2\text{-}2.7 \mu\text{m}$  fractions ( $R = 0.689$ ,  $P =$   
249  $0.001$ ) (Fig. S1B). As expected due to generally larger cell sizes in microbial  
250 eukaryotes compared to prokaryotes, species richness remained higher in the  $>$   
251  $2.7 \mu\text{m}$  vs.  $0.2\text{-}2.7 \mu\text{m}$  fractions (Fig. S3C-D). Richness in the  $> 2.7 \mu\text{m}$  fraction  
252 gradually increased downriver, similarly to prokaryotic communities, but remained  
253 relatively stable among the  $0.2\text{-}2.7 \mu\text{m}$  fraction. Both the  $> 2.7 \mu\text{m}$  (stress =  $0.113$   
254 and  $0.2\text{-}2.7 \mu\text{m}$  (stress =  $0.146$ ) fractions also showed a significant separation  
255 between the lower and upper MSR (ANOSIM,  $>2.7 \mu\text{m}$ :  $R = 0.696$ ,  $P=0.001$ ;  $0.2\text{-}$   
256  $2.7 \mu\text{m}$ :  $R = 0.576$ ,  $P = 0.001$ ) (Fig. 2C, D). Distance and phosphate constituted  
257 the top two environmental factors influencing this distinction ( $r =0.75$ ,  $r =0.48$ ;  $r =$   
258  $0.70$ ,  $r =0.54$ , respectively) (Fig. 2C, D; Table S1). No other factors had  
259 correlations  $> 0.40$  (Table S1). At the OTU level, taxa related to an unclassified  
260 Ochrophyta (OTU63) and unclassified Eukaryote (OTU1) separated the MSR  
261 communities in the  $0.2\text{-}2.7 \mu\text{m}$  ( $r > 0.63$ ,  $P= 0.001$ ), while the same unclassified  
262 Eukaryote OTU (OTU1) and a second unclassified Eukaryote (OTU222)  
263 contributed most to separating the  $>2.7 \mu\text{m}$  communities ( $r > 0.80$ ,  $P = 0.001$ )  
264 (Figure 2C and D, Table S1).

265 Stramenopiles (or Heterokonts), encompassing diatoms and many other  
266 forms of algae, and OTUs that could not be classified at the phylum level,  
267 dominated both the  $> 2.7 \mu\text{m}$  and  $0.2\text{-}2.7 \mu\text{m}$  communities (Fig. 5).  
268 Stramenopiles accounted for over 25% of both communities, with higher  
269 abundances in the upper vs. lower river. We observed a similar trend of disparate  
270 abundances between the upper and lower river for  $> 2.7 \mu\text{m}$  Cryptomonadales  
271 and  $0.2\text{-}2.7 \mu\text{m}$  Nucleomycea, which include fungi (Fig. 5A; Table S1). Within the  
272  $0.2\text{-}2.7 \mu\text{m}$  fraction, we classified the five most abundant OTUs as three  
273 unclassified *Bacillariophytina* (OTU7, OTU14, OTU9), a *Pythium* sp. (OTU170),  
274 and a unclassified *Cryptomonas* (OTU11) (Table S1). Comparatively, two  
275 unclassified *Eukaryotes* (OTU2 and OTU1), a unclassified *Stramenopiles*  
276 (OTU3), a unclassified *Perkinsidae* (OTU13), and a unclassified *Chrysophyceae*  
277 (OTU6) had the highest abundance in the  $> 2.7 \mu\text{m}$  (Table S1).

278 Eighty OTUs comprised the  $> 2.7 \mu\text{m}$  core microbiome, averaging 28% of  
279 aggregate community relative abundance (Fig. 4B). We classified these as  
280 *Alveolata*, *Cryptophyceae*, *Nucleomycea*, *Stramenopiles*, or unclassified  
281 Eukaryota (Table S1). Again, consistent with larger organism sizes, and thus  
282 fewer OTUs overall, the  $0.2\text{-}2.7 \mu\text{m}$  Eukaryotic core microbiome comprised only  
283 21 OTUs that, in aggregate, averaged 20% of the community relative abundance  
284 across all samples (Fig. 4B). These OTUs consisted of *Alveolata*, *Nucleomycea*,  
285 *Stramenopiles*, or unclassified Eukaryota (Table S1). While the  $0.2\text{-}2.7 \mu\text{m}$  core

286 microbiome remained relatively stable along the river, the  $> 2.7 \mu\text{m}$  core  
287 decreased along the upper MSR before stabilizing in the lower river, similarly to  
288 that of the prokaryotes (Fig. 4B).

289

#### 290 *Network analyses identify indicator taxa associated with eutrophication*

291 Since increased nitrogen and phosphorous is of interest for stakeholders both  
292 along the river and in the GOM, we evaluated relationships between individual  
293 OTUs and these nutrients to identify taxa indicative of eutrophication. Among 0.2-  
294  $2.7 \mu\text{m}$  prokaryotes, a co-correlation network (submodule) strongly associated  
295 with nitrate ( $r = 0.6$ ,  $P = 7\text{e-}08$ ) (Fig. S6A) comprised 77 OTUs, mostly from the  
296 Proteobacteria. OTU relative abundances explained 42% of the variance in  
297 nitrate (LOOCV,  $R^2=0.65$ ;  $\text{corr} = 0.65$ ,  $P = 1.\text{e-}09$ ), and the top four OTUs by VIP  
298 score belonged to the Proteobacteria and Bacteroidetes (Fig. 6B; FigS6A; Table  
299 S1). Of these, three were highly correlated to nitrate ( $r > 0.50$ ) (Table S1):  
300 unclassified *Comamonadaceae*, *Mucilaginibacter*, and *Pseudospirillum* spp..  
301 OTUs that had more than 8 node connections, indicating taxa co-correlating with  
302 numerous others, included a *Mucilaginibacter* sp., a *Novosphingobium* sp., and a  
303 *Caulobacter* sp. (Fig. 6B, Table S1).

304 The submodule best correlated with phosphate ( $r = 0.53$ ,  $P = 2\text{e-}06$ ) (Fig.  
305 S6) comprised 151 OTUs that explained 80% of the variance (LOOCV,  $R^2=0.80$ ;  
306  $\text{corr} = 0.89$ ,  $P = < 2.3\text{e-}16$ ) (Fig. S6D). OTUs with VIP scores  $> 1$  belonged to  
307 seven different phyla (Table S1), with the top four identified as an unclassified  
308 *Holophagaceae*, an unclassified *Gemmatimonadaceae*, an unclassified  
309 *Burkholderiaceae*, and a *Pseudospirillum* sp. (Table S1). Three of these four  
310 OTUs had Pearson correlations with phosphate greater than 0.62 (Fig. 6B). The  
311 most highly interconnected OTU within the submodule was an unclassified  
312 Acidobacteria (Table S1).

313 In the  $> 2.7 \mu\text{m}$  fraction, the prokaryotic submodule with the highest  
314 correlation to nitrate ( $r = 0.56$ ,  $P = 3\text{e-}07$ ) (Fig. S7A) had 133 OTUs and  
315 explained 69% of the variation in nitrate (LOOVC,  $R^2 = 0.69$ ;  $\text{corr} = 0.83$ ,  $P = <$   
316  $2.2\text{e-}16$ ). The four highest VIP scoring OTUs, an *Anabaena* sp., a  
317 *Flavobacterium* sp., an unclassified bacterium, and a member of the  
318 *Sphingobacteriales* NS11-12 marine group, anticorrelated with nitrate (Fig. 7A,  
319 Table S1). A NCBI BLAST of the unclassified bacterium returned no significant  
320 hits  $> 90\%$  identity to named organisms. OTUs with the highest node centrality  
321 ( $> 20$ ) belonged to two clades of bacteria, *Sphingomonadales*  
322 (*Alphaproteobacteria*) and *Sphingobacteriales* (Bacteroidetes), and all correlated  
323 positively to nitrate ( $r > 0.48$ , Table S1).

324 When considering phosphate, the most significant submodule showed a  
325 modest correlation ( $r = 0.53$ ,  $P = 1\text{e-}06$ , Fig. S7A), but OTU abundances could  
326 only explain 48% of the variation of measured phosphate (LOOVC,  $R^2 = 0.48$ ;  
327  $\text{corr} = 0.77$ ,  $P = 4.88\text{e-}15$ ). Proteobacteria constituted the majority of the 80 OTUs  
328 from the submodule, and the top four scoring VIP OTUs were an unclassified  
329 *Gammaproteobacteria*, an *Arcicella* sp., an unclassified *Cytophagaceae*, and a

330 *Nitrospira* sp. (Fig. 7B, Table S1); the latter two had moderate correlations to  
331 phosphate ( $r = > 0.55$ ). The OTUs with the highest number of node connections  
332 were a *Woodsholea* sp. from the *Caulobacterales* family and the same *Nitrospira*  
333 OTU (Table S1).

334 Among 0.2-2.7  $\mu\text{m}$  Eukaryotic size fraction, a submodule strongly  
335 associated with nitrate ( $r = 0.39$ ,  $P = 6\text{e-}04$ ) (Fig. S8A) comprised 39 OTUs,  
336 mostly from the phylum Stramenopiles. OTU relative abundances explained 38%  
337 of the variance in nitrate (LOOCV,  $R^2=0.38$ ;  $\text{corr} = 0.62$ ,  $P = 2.97\text{e-}9$ ), and the top  
338 four OTUs by VIP score classified as an unclassified *Chrysophyceae*, an  
339 unclassified *Ochrophyta*, and an unclassified *Chromulinales* (Fig. 8B; FigS8A;  
340 Table S1). Of these, two highly correlated to nitrate ( $r > 0.49$ ) (Table S1): the  
341 *Chrysophyceae*, and *Chromulinales* OTUs.

342 The submodule best correlated with phosphate ( $r = 0.56$ ,  $P = 2\text{e-}07$ ) (Fig.  
343 S8A) comprised 56 OTUs that explained 80% of the variance (LOOCV,  $R^2=0.80$ ;  
344  $\text{corr} = 0.89$ ,  $P = < 2\text{e-}16$ ). OTUs with VIP scores  $> 1$  belonged to four different  
345 phyla (Table S1), with the top four identified as an unclassified  
346 *Peronosporomycetes*, an unclassified *Ochrophyta*, an unclassified Eukaryote,  
347 and an unclassified Stramenopiles. (Table S1). All four OTUs had Pearson  
348 correlations with phosphate greater than 0.60, two were negative (Fig. 8B). The  
349 most highly interconnected OTU within the submodule was an unclassified  
350 Eukaryote (Table S1).

351 Among the  $> 2.7 \mu\text{m}$  Eukaryote taxa, submodules with strongest  
352 correlations to nitrate and phosphate ( $\text{NO}_3^-$ :  $r = 0.52$ ,  $P = 2\text{e-}06$ ;  $\text{PO}_4^{3-}$ :  $r = 0.60$ ,  $P$   
353  $= 2\text{e-}08$ , Fig. S9A) had smaller membership than those for prokaryotes,  
354 comprising 59 and 39 OTUs, respectively. The OTUs in the submodule most  
355 strongly correlated with nitrate could predict 57% of observed variation in nitrate  
356 (LOOVC,  $R^2=0.572$ ;  $\text{corr} = 0.759$ ,  $P = 6.7\text{e-}15$ ). OTUs with top VIP scores were  
357 two unclassified *Chrysophyceae*, an unclassified *Ochrophyta*, and an  
358 unclassified Diatom (Fig. 8B; Table S1). When considering phosphate,  
359 submodule OTU abundances predicted 62% of measured concentrations  
360 (LOOVC,  $R^2 = 0.618$ ;  $\text{corr} = 0.799$ ,  $P = < 2\text{e-}16$ ). An unclassified Eukaryote and  
361 an unclassified *Peronosporomycetes* occupied top two phosphate-associated  
362 positions according to VIP score (Fig. 9B; Table S1).

363

## 364 **Discussion**

365 Our 2914 km transect of the MSR has supplied the longest  
366 microbiologically oriented transect of a top ten river based on volume, length, or  
367 drainage. Water samples illustrated a river continually inundated with nutrients  
368 and sediment in its upper portion that gave way to more consistent levels in the  
369 lower river (Fig. 1B). These distinct scenarios mirrored the different microbial  
370 regimes separated by the Missouri river confluence (Fig. 2), which differed from  
371 the historical distinction of the upper and lower MSR at the Ohio River confluence  
372 in Cairo, IL, but matched the separation based on changes in Strahler order from  
373 eight to ten<sup>48</sup>. In general, both prokaryotic fractions and the  $> 2.7 \mu\text{m}$  fraction of

374 eukaryotic communities increased in richness downriver while the percent of core  
375 community taxa decreased in the upper MSR before settling in the lower river  
376 (Figs. 4, S3), concomitant with increased stability of environmental factors (Fig.  
377 1B). Co-occurrence network analyses identified important potential indicator taxa  
378 for the eutrophication nutrients nitrate and phosphate that may help future efforts  
379 to detect and quantify imminent changes in river water quality.

380 In general, the most abundant OTUs throughout the MSR corresponded to  
381 typical freshwater taxa observed in other important riverine/aquatic  
382 studies<sup>4,25,29,35,49,50</sup>, such as LD12 (*Alphaproteobacteria*), hcgl-cal clade  
383 (*Actinobacteria*), *Polynucelobacter* (*Betaproteobacteria*), LD28  
384 (*Betaproteobacteria*), and *Limnohabitans* (*Betaproteobacteria*) (Table S1).  
385 Though our nutrient measurements mirrored conditions found in previous studies  
386 of the Danube, Thames, and upper Mississippi rivers<sup>4,27,34</sup>, the relative  
387 abundances of important phyla differed. Specifically, in our study Proteobacteria  
388 remained the most abundant phylum, while Bacteroidetes and Actinobacteria  
389 decreased (Fig. 3), whereas in the Minnesota portion of the MSR<sup>35</sup>, Thames<sup>29</sup>,  
390 Danube<sup>4</sup>, and Columbia Rivers<sup>40</sup>, Bacteroidetes dominated headwaters while  
391 Actinobacteria dominated further downriver. General increases in MSR species  
392 richness with distance and decreases in the percent core community along the  
393 upper river contrasted predictions by the RCC<sup>36</sup> and observations in previous  
394 studies<sup>4,27</sup> where these trends were reversed. Importantly, we did not sample the  
395 true headwaters of the MSR (Lake Itasca to above St. Cloud), and therefore at  
396 the point of first sampling, the MSR already constituted an eighth order river.  
397 Ultimately, some of these variant observations may result from different sampling  
398 methodologies, but also from biological signal related to unique environmental  
399 conditions and human impacts, changes in the level of river engineering with  
400 distance, and the magnitude of the MSR (in terms of volume and catchment  
401 complexity) and its tributaries relative to previously sampled systems. However,  
402 perhaps some differences lie in the proportion and scale at which the varied  
403 influences of dispersal and environmental filtering occur.

404 Many of our results suggest similar underlying ecological mechanisms with  
405 other river systems, even though the specific microbial community patterns differ.  
406 Multiple studies have shown evidence for the importance of mass effects in  
407 headwaters, while species sorting dominates with increased residence time as  
408 rivers grow in size<sup>44,52</sup>. Our data suggests that mass effects indeed play a role in  
409 the upper MSR, although instead of only in the headwaters, this process  
410 continues for almost a third of the length of the river. Increased turbidity  
411 correlated with decreases in freshwater bacteria (Spearman rank correlation, >  
412 2.7  $\mu\text{m}$  R = 0.55; 0.2-2.7  $\mu\text{m}$  R = 0.51) and the core microbiome (Spearman rank  
413 correlation, >2.7  $\mu\text{m}$  R = 0.53; 0.2-2.7  $\mu\text{m}$  R = 0.63) during the first ~1000 km- the  
414 upper MSR- whereas these variables and some nutrient concentrations (Fig. 1)  
415 stabilized in the lower MSR. These patterns are consistent with communities  
416 under the influence of mass effects from tributaries in the upper MSR. Once the  
417 MSR grew to a tenth order river, the large volume and size potentially buffered it

418 from allochthonous influences, allowing species sorting effects to dominate. The  
419 lower river represents a more stable environment (e.g. nutrient concentrations)  
420 with its increased size and volume, contrasting the more variant upper MSR.  
421 Though the river speed increases, the effective residence time also increases  
422 since taxa no longer experience rapidly changing environmental variables.

423 That we still observed variation between microbial communities along the  
424 lower MSR concurs with environmental filtering attributable to unmeasured  
425 bottom-up factors, such as the quality and quantity of DOM, or top-down  
426 influences such as predation or viral lysis. An overall community shift from a  
427 mixture of allochthonous members to a “native” population requires growth rates  
428 that allow taxa to overcome mass effects over a given distance<sup>53</sup>. The lower river  
429 also provides ample opportunities for microbial community differentiation based  
430 on average prokaryotic growth rates<sup>4</sup>, especially among particle-associated (>  
431 2.7  $\mu\text{m}$ ) taxa<sup>40</sup>. Thus, while the aggregate patterns in particular phyla and  
432 taxonomic richness may differ from other systems, similar ecological processes  
433 may still occur, but the relative proportion of distance whereby mass effects vs.  
434 species sorting dominate fosters unique community dynamics.

435 Within microbial communities lie taxa with unique roles as potential  
436 mitigators of human impact. MSR water quality has continually been rated low<sup>54</sup>,  
437 suffers from significant eutrophication, and causes one of the largest worldwide  
438 zones of seasonal hypoxia in the northern GOM<sup>21</sup>. Identification of  
439 microorganisms with relationships to nitrogen and phosphorous will inform efforts  
440 to model nutrient remediation and provide important targets for future study.  
441 However, such microbial indicators of biological integrity (IBIs- metrics to quantify  
442 the health of an aquatic system<sup>55</sup>), require better development<sup>55-57</sup>. Phosphate  
443 and nitrate increased along the upper MSR and generally stabilized in the lower  
444 river (Fig. 1). Using co-correlation networks and partial least squares modeling  
445 (PLS), we identified submodules with containing taxa with significant predictive  
446 power for phosphate and nitrate concentrations (Figs. 6-9). Individual OTUs  
447 within these submodules with strong correlations to nitrate or phosphate, and VIP  
448 scores > 1 within the PLS models, represent potential indicator taxa for these  
449 nutrients.

450 Taxa in a 0.2-2.7  $\mu\text{m}$  fraction submodule could predict 80% of the variance  
451 in phosphate concentrations, with much weaker correlations to nitrate. A  
452 *Holophagaceae* OTU (OTU33) was the top 0.2-2.7  $\mu\text{m}$  taxon for predicting  
453 phosphate concentrations based on VIP scores (Fig. 6B). It also occupied an  
454 important role in separating communities at a beta-diversity level (Fig. 2B), and  
455 resided in the core microbiome. A 2014 study on the tributaries of the MSR found  
456 that the Ohio River had much higher abundance of Acidobacteria relative to other  
457 tributaries<sup>39</sup>, which corroborates our finding of increasing Acidobacteria  
458 downriver. Although the *Holophagaceae* comprises diverse and abundant taxa<sup>58</sup>,  
459 many environmentally relevant taxa remain elusive to cultivation efforts<sup>59</sup>. One of  
460 the few cultivated representatives from the *Holophagaceae* family, *Holophaga*  
461 *foetida*, has the genomic capacity to accumulate phosphorous and synthesize

462 polyphosphate via a polyphosphate kinase<sup>60</sup>. While this organism has an  
463 obligately anaerobic lifestyle\* that makes it an unlikely match for the OTU, it  
464 indicates genomic potential that may span the clade. An unclassified  
465 *Gemmatimonadaceae* (OTU60) from the 0.2-2.7  $\mu\text{m}$  fraction correlated strongly  
466 with phosphate and also played a strong role in driving beta-diversity changes  
467 between the upper and lower river (Fig. 2)

468 Contrasting the 0.2-2.7  $\mu\text{m}$  taxa, a submodule in the > 2.7  $\mu\text{m}$  fraction  
469 could predict 69% of nitrate concentration variance. Within the prokaryotic > 2.7  
470  $\mu\text{m}$  taxa, an *Anabaena* sp. (OTU40) had the top VIP score among submodule  
471 taxa predicting nitrate, correlated negatively (Fig. 7A), and had membership in  
472 the core microbiome. The nitrogen-fixing *Anabaena*<sup>61</sup> typically bloom in low  
473 dissolved inorganic nitrogen (DIN) conditions, making the absence of these  
474 consistent with high DIN. Organisms in the *Sphingomonadaceae* (e.g.,  
475 *Novosphingobium* spp.) also contributed strongly to the PLS models predicting  
476 nitrate with both > 2.7  $\mu\text{m}$  (Fig. 7A) and 0.2-2.7  $\mu\text{m}$  (Fig. 6A) taxa (Table S1).  
477 *Novosphingobium* isolates have previously been associated with eutrophic  
478 environments<sup>34,64-66</sup> and some can reduce nitrate<sup>64,65</sup>, making these OTUs good  
479 candidate IBIs within the MSR basin.

480 Notably, a *Nitrospira* sp. (OTU96) significantly correlated with phosphate  
481 ( $R^2=0.60$ ,  $P=<0.001$ ) and was the second most important taxa for predicting  
482 phosphate concentrations in the > 2.7  $\mu\text{m}$  fraction (VIP = 1.8). *Nitrospira* occupy  
483 a key role as nitrite oxidizers, and have been found in various environments  
484 including the “Dead Zone” in the GoM<sup>22</sup> and in correlation with wastewater  
485 treatment effluent<sup>67</sup>. This OTU makes an intriguing IBI candidate in a watershed  
486 impacted by wastewater<sup>68</sup> and that directly influences the Dead Zone.  
487 Additionally, as an organism that creates nitrate, this might represent a taxon that  
488 exacerbates, rather than mitigates, the effects of eutrophication.

489 Within Eukaryotes, multiple algae and diatoms strongly correlated with  
490 nitrate and phosphate (Fig. 8 and 9; Table S1), and specifically *Chrysophyceae*  
491 taxa from both fractions correlated strongly with nitrate (Fig. 8A). *Chrysophyceae*  
492 (golden algae) commonly occupy river systems<sup>69</sup> including the MSR<sup>68</sup>, can be  
493 autotrophic and mixotrophic<sup>70</sup>, and may serve as predators of prokaryotes<sup>71</sup>.  
494 While we also identified many other eukaryotes as important predictors of  
495 nutrients, poor taxonomic resolution hindered our ability to discuss them further.  
496 Improved cultivation and systematics of key microbial eukaryotes will be vital to  
497 understanding river nutrient dynamics.

498 While the most geographically comprehensive to date, this study only  
499 encompasses a snapshot in time for the MSR. Seasonal changes that have been  
500 observed in the Minnesota portion of the upper MSR<sup>52</sup> and the Columbia  
501 River<sup>52,72,73</sup> undoubtedly influence this dynamic system. Future studies should  
502 incorporate microbial responses, at a full river scale, to seasonal pulse events  
503 (e.g. rain, snow melt) and how river size and volume may buffer local microbial  
504 communities from allochthonous inputs. Our current research highlights the  
505 uniqueness and complexities of the MSR ecosystem. The observed association

506 between changes in Strahler's river order, nutrient dynamics, and community  
507 composition indicates the importance of hydrology<sup>46,47</sup> on the spatial dynamics  
508 structuring microbial communities and provides baseline information for future  
509 MSR studies that incorporate greater temporal and spatial resolution. With water  
510 quality and river health of growing local and global importance<sup>54,74</sup>, the  
511 determination of candidate microbial IBIs also provides impetus for targeted  
512 research on their functions and further investigation of these organisms as  
513 sentinels of river health.

514

515

## 516 **Materials and Methods**

517

### 518 *Sampling and Cell Counts*

519 We used rowboats and a simple filtration protocol (Supplementary Information) to  
520 collect water from 39 sites along a continually rowed transect of the MSR,  
521 starting in Lake Itasca and ending in the GOM, over 70 days from September 18<sup>th</sup>  
522 to November 26<sup>th</sup>, 2014. Sites were chosen to be near major cities and above  
523 and below large tributaries. After some samples were removed due to insufficient  
524 sequence data, contamination, or incomplete metadata (see below), the final  
525 usable set of samples included 38 sites starting at Minneapolis (Fig. 1A, Table  
526 S1). Most sampling occurred within the body of the river, although due to safety  
527 issues, three samples were collected from shore (Table S1). We collected  
528 duplicate samples at each site, but because separate rowboat teams frequently  
529 collected these sometimes several dozen meters apart, they cannot be  
530 considered true biological replicates and we have treated them as independent  
531 samples. At each site, we filtered 120 mL of water sequentially through a 2.7  $\mu\text{m}$   
532 GF/D filter (Whatman GE, New Jersey, USA) housed in a 25mm polycarbonate  
533 holder (TISCH, Ohio, USA) followed by a 0.2  $\mu\text{m}$  Sterivex filter (EMD Millipore,  
534 Darmstadt, Germany) with a sterile 60 mL syringe (BD, New Jersey, USA). We  
535 refer to fractions collected on the 2.7  $\mu\text{m}$  and 0.22  $\mu\text{m}$  filters as  $> 2.7 \mu\text{m}$  and 0.2-  
536 2.7  $\mu\text{m}$ , respectively. Flow-through water from the first 60 mL was collected in  
537 autoclaved acid-washed 60 mL polycarbonate bottles. Both filters were wrapped  
538 in parafilm, and together with the filtrate, placed on ice in Yeti Roadie 20 coolers  
539 (Yeti, Austin, TX) until shipment to LSU. Further, 9 mL of whole water for cell  
540 counts was added to sterile 15 mL Falcon tubes containing 1 mL of formaldehyde  
541 and placed into the cooler. We monitored cooler temperature with HOBO loggers  
542 (Onset, Bourne, MA) to ensure samples stayed at  $\leq 4^\circ\text{C}$ . The final cooler  
543 containing samples from sites P-AI had substantial ice-melt. Though our filters  
544 were wrapped in parafilm, we processed melted cooler water alongside our other  
545 samples to control for potential contamination in these filters. Given that some of  
546 our samples were expected to contain low biomass, we also included duplicate  
547 process controls for kit contamination<sup>75,76</sup> with unused sterile filters. Flow-through  
548 0.2  $\mu\text{m}$  filtered water from each collection was analyzed for  $\text{SiO}_4$ ,  $\text{PO}_4^{3-}$ ,  $\text{NH}_4^+$ ,  
549  $\text{NO}_3^-$ , and  $\text{NO}_2^-$  ( $\mu\text{g/L}$ ) at the University of Washington Marine Chemistry

550 Laboratory  
551 (<http://www.ocean.washington.edu/story/Marine+Chemistry+Laboratory>). Aboard-  
552 rowboat measurements were taken for temperature and turbidity. We determined  
553 turbidity by deploying a secchi disk (Wildco, Yulee, FL), while drifting with the  
554 current so the line hung vertically. It was lowered until no longer visible, then  
555 raised until just visible, and measured for its distance below the waterline. We  
556 then calculated secchi depth from the average of two measurements.  
557 Temperature was measured with probes from US Water Systems (Indianapolis,  
558 IN), rinsed with distilled water between samples. Samples for cell counts were  
559 filtered through a 2.7  $\mu\text{m}$  GF/D filter, stained with 1x Sybr Green (Lonza), and  
560 enumerated using the Guava EasyCyte (Millipore) flow cytometer as previously  
561 described<sup>77</sup>.

### 562 563 *DNA extraction and Sequencing*

564 DNA was extracted from both filter fractions and controls using a MoBio  
565 PowerWater DNA kit (MoBio Laboratories, Carlsbad, CA) following the  
566 manufacturer's protocol with one minor modification: in a biosafety cabinet (The  
567 Baker Company, Stanford, ME), sterivex filter housings were cracked open using  
568 sterilized pliers and filters were then removed by cutting along the edge of the  
569 plastic holder with a sterile razor blade before being placed into bead-beating  
570 tubes. DNA was eluted with sterile MilliQ water, quantified using the Qubit2.0  
571 Fluorometer (Life Technologies, Carlsbad, CA), and stored at -20° C. Bacterial  
572 and archaeal sequences were amplified at the V4 region of the 16S rRNA gene  
573 using the 515f and 806r primer set<sup>78</sup>, and eukaryotic sequences from the V9  
574 region of the 18S rRNA gene using the 1391r and EukBR primer set<sup>79</sup>. Amplicons  
575 were sequenced on an Illumina MiSeq as paired-end 250 bp reads at Argonne  
576 National Laboratory. Sequencing of the 16S and 18S rRNA gene amplicons  
577 resulted in 13253140 and 13240531 sequences, respectively.

### 578 579 *Sequence Analysis*

580 We analyzed amplicon data with Mothur v.1.33.3<sup>80</sup> using the Silva v.119  
581 database<sup>81,82</sup>. Briefly, 16S and 18S rRNA gene sequences were assembled into  
582 contigs and discarded if the contig had any ambiguous base pairs, possessed  
583 repeats greater than 8 bp, or were greater than 253 or 184 bp in length,  
584 respectively. Contigs were aligned using the Silva rRNA v.119 database,  
585 checked for chimeras using UCHIME<sup>83</sup>, and classified also using the Silva rRNA  
586 v.119 database. Contigs classified as chloroplast, eukaryotes, mitochondria, or  
587 "unknown;" or as chloroplast, bacteria, archaea, mitochondria, or "unknown;"  
588 were removed from 16S or 18S rRNA gene data, respectively. The remaining  
589 contigs were clustered into operational taxonomic units (OTUs) using a 0.03  
590 dissimilarity threshold (OTU<sub>0.03</sub>). After these steps, 146725 and 131352 OTUs  
591 remained for the 16S and 18S rRNA gene communities, respectively.

### 592 593 *Sample quality control*

594 To evaluate the potential for contamination from extraction kits, cooler water in  
595 the last set of samples, or leaking/bursting of pre-filters, all samples were  
596 evaluated with hierarchical clustering and NMDS analysis. Hierarchical clustering  
597 was performed in R using the *hclust* function with methods set to “average”, from  
598 the *vegan* package<sup>84</sup>. Samples were removed from our analysis if they were  
599 observed to be outliers in both the NMDS and hierarchical clustering such that  
600 they grouped with our process controls. The process and cooler water controls  
601 were extreme outliers in both, as was sample L2 (Fig. S1, S2). Sterivex and  
602 prefilter samples generally showed strong separation with the exception of three  
603 16S rRNA gene samples- STER X2, W2, S2 (Fig. S1, S2). The only other  
604 samples that were removed were due to missing chemical data (Lake Itasca1-2,  
605 A1-2) or failed sequencing (16S STER Af1; 16S PRE S2, X2; 18S PRE O1). Not  
606 including process or cooler water controls, 152 samples were sequenced each  
607 for prokaryotic and eukaryotic communities. After these QC measures, 144 and  
608 149 samples remained in the analyses from the 16S and 18S rRNA gene  
609 amplicons, respectively. Further, to control for potential contaminants, any OTU  
610 with greater 20 reads in the process or cooler controls was removed from the  
611 data set. 146725 and 131327 OTUs remained after these steps for 16S and 18S  
612 rRNA gene communities, respectively.

613

#### 614 *Alpha and Beta Diversity*

615 OTU<sub>0.03</sub> analyses were completed with the R statistical environment v.3.2.1<sup>85</sup>.  
616 Using the package *PhyloSeq*<sup>86</sup>, alpha-diversity was first calculated on the  
617 unfiltered OTUs using the “estimate richness” command within *PhyloSeq*, which  
618 calculates Chao1<sup>86</sup>. After estimating chao1, potentially erroneous rare OTUs,  
619 defined here as those without at least two sequences in 20% of the data, were  
620 discarded. After this filter, the dataset contained 950 and 724 16S and 18S rRNA  
621 gene OTUs, respectively. For site-specific community comparisons, OTU counts  
622 were normalized using the package *DESeq2*<sup>87</sup> with a variance stabilizing  
623 transformation<sup>88</sup>. Beta-diversity between samples was examined using Bray-  
624 Curtis distances via ordination with non-metric multidimensional scaling (NMDS).  
625 Analysis of similarity (ANOSIM) was used to test for significant differences  
626 between groups of samples of the NMDS analyses using the *anosim* function in  
627 the *vegan* package<sup>84</sup>. The influence of environmental parameters on beta-  
628 diversity was calculated in R with the *envfit* function.

629

#### 630 *Network analyses and modeling*

631 To identify specific OTUs with strong relationships to environmental parameters  
632 (e.g. turbidity, NO<sub>3</sub><sup>-</sup>), we employed weighted gene co-expression network  
633 analysis (WGCNA)<sup>89</sup> as previously described<sup>90</sup> for OTU relative abundances.  
634 First, a similarity matrix of nodes (OTUs) was created based on pairwise Pearson  
635 correlations across samples. This was transformed into an adjacency matrix by  
636 raising the similarity matrix to a soft threshold power ( $p$ ;  $p = 6$  for 16S and 18S >  
637 2.7  $\mu\text{m}$ ,  $p = 4$  for 16S 0.2-2.7  $\mu\text{m}$ ) that ensured scale-free topology. Submodules

638 of highly co-correlating OTUs were defined with a topological overlap matrix and  
639 hierarchical clustering. Each submodule, represented by an eigenvalue, was  
640 pairwise Pearson correlated to individual environmental parameters (Figs. S6-8  
641 A). To explore the relationship of submodule structure to these parameters,  
642 submodule OTUs were plotted using their individual correlation to said parameter  
643 (here nitrate or phosphate) and their submodule membership, defined as the  
644 number of connections within the module (Figs. S6-8 B, D). Strong correlations  
645 between submodule structure and an environmental parameter facilitate  
646 identification of OTUs that are highly correlated to that parameter. To evaluate  
647 the predictive relationship between a submodule and a parameter, we employed  
648 partial least square regression (PLS) regression analysis. PLS maximizes the  
649 covariance between two parameters (e.g. OTU abundance and nitrate  
650 concentration) to define the degree to which the measured value (OTU  
651 abundance) can predict the response variable (nutrient concentration). The PLS  
652 model was permuted a 1000 times and Pearson correlations were calculated  
653 between the response variable and leave-one-out cross-validation (LOOCV)  
654 predicted values. Modeled values were then compared with measured values to  
655 determine the explanatory power of the relationships (Figs. S6-8 C, E). Relative  
656 contributions of individual OTUs to the PLS regression were calculated using  
657 value of importance in the projection (VIP)<sup>91</sup> determination. PLS was run using  
658 the R package *pls*<sup>92</sup>, while VIP was run using an additional code found here:  
659 <http://mevik.net/work/software/VIP.R>.

660

#### 661 *Environmental Ontology*

662 Environmental ontology of individual 16S rRNA gene OTUs was determined  
663 using the SEQenv (v1.2.4) pipeline (<https://github.com/xapple/seqenv>) as  
664 previously described<sup>4</sup>. Briefly, representative sequences of our OTUs were  
665 searched against the NCBI nt database (updated on 07/01/2016) using BLAST  
666 and filtered for hits with a minimum of 99% identity. From each hit, a text query of  
667 the metadata was performed to identify terms representing the sequence's  
668 environmental source. The text was mined for EnvO terms  
669 (<http://environmentontology.org/>) and the frequency in which the terms appeared  
670 for each OTU was recorded. Using the seq\_to\_names output provided by  
671 SEQenv, EnvO terms were formed into eight groups: Freshwater, Aquatic  
672 Undetermined, Salt Water, Anthropogenic, Other, Terrestrial, Sediment, and  
673 Unclassified (Table S1). To be assigned a group, an OTU had to have the  
674 majority (> 50%) of its hits classified to that term, while equal distribution between  
675 two or more groups were classified as Unclassified. OTUs that returned no  
676 significant hits to an EnvO term were assigned to a ninth category, NA. OTUs  
677 and their corresponding relative abundances were merged based on the  
678 assigned group and plotted.

679

#### 680 *Accession numbers*

681 Community 16S and 18S rRNA gene community sequence fastq files are  
682 available at the NCBI Sequence Read Archive under the accession numbers:  
683 SRR3485674- SRR3485971 and SRR3488881- SRR3489315.

684

#### 685 *Code Availability*

686 All code used for Mothur, SeqENV, PhyloSeq, WGCNA, and PLS regression  
687 analyses can be found on the Thrash lab website  
688 (<http://thethrashlab.com/publications>) with the reference to this manuscript linked  
689 to “Supplementary Data”.

690

#### 691 **Acknowledgements**

692 This work was supported by the Department of Biological Sciences, College of  
693 Science, and the Office of Research and Economic Development at Louisiana  
694 State University, and the College of the Environment at University of Washington.  
695 We would like to thank Dr. Gary King and Dr. Caroline Fortunato for their friendly  
696 reviews. The authors also thank the countless volunteers, schools, and  
697 organizations that facilitated the research. We specifically thank Pete Weess,  
698 Jessica Zimmerman, Katy Welch, Brian Moffitt, and David Cheney for helping  
699 organize the shipment of coolers between sites, and the OAR Northwest  
700 sponsors- Seattle Yacht Club Foundation, Yetti Coolers, and the National  
701 Mississippi River Museum and Aquarium. A full list of sponsors and volunteers  
702 can be found on the OAR Northwest website ([oarnorthwest.org](http://oarnorthwest.org)). We would also  
703 like to thank Dr. Matthew Sullivan and Dr. Simon Roux for their support and help  
704 scripting the code for the WGCNA and sPLS analyses. Lastly, we would like to  
705 thank Mrs. Ginger Thrash, who connected the Thrash lab to OAR Northwest.

706

#### 707 **Author Contributions**

708 M.W.H., J.H., F.S., and J.C.T. designed the study, J.H., G.S., P.F., and M.P.  
709 collected the data, M.W.H. processed the samples, M.W.H. and J.C.T. analyzed  
710 the data, M.W.H. and J.C.T. wrote the manuscript, and all authors contributed to  
711 the editing of the manuscript.

#### 712 **Conflict of Interest**

713 The authors declare no competing financial interests.

#### 714 **Literature Cited**

715

- 716 1. Richey, J. E., Melack, J. M., Aufdenkampe, A. K., Ballester, V. M. & Hess,  
717 L. L. Outgassing from Amazonian rivers and wetlands as a large tropical  
718 source of atmospheric CO<sub>2</sub>. *Nature* **416**, 617–620 (2002).
- 719 2. Ensign, S. H. & Doyle, M. W. Nutrient spiraling in streams and river  
720 networks. *J. Geophys. Res. Biogeosciences* **111**, (2006).
- 721 3. Withers, P. J. A. & Jarvie, H. P. Delivery and cycling of phosphorus in  
722 rivers: A review. *Sci. Total Environ.* **400**, 379–395 (2008).

- 723 4. Savio, D. *et al.* Bacterial diversity along a 2600 km river continuum.  
724 *Environ. Microbiol.* **17**, n/a–n/a (2015).
- 725 5. Battin, T. J. *et al.* Biophysical controls on organic carbon fluxes in fluvial  
726 networks. *Nat. Geosci.* **2**, 595–595 (2009).
- 727 6. Cauwet, G., Hansell, D. A. & Carlson, C. A. in *Biogeochemistry of marine*  
728 *dissolved organic matter*. 579–609 (Academic Press, San Diego, CA,  
729 2002).
- 730 7. Cole, J. J. *et al.* Plumbing the global carbon cycle: Integrating inland waters  
731 into the terrestrial carbon budget. *Ecosystems* **10**, 171–184 (2007).
- 732 8. Cole, J. J. & Caraco, N. F. Carbon in catchments: Connecting terrestrial  
733 carbon losses with aquatic metabolism. in *Marine and Freshwater*  
734 *Research* **52**, 101–110 (2001).
- 735 9. Dagg, M., Benner, R., Lohrenz, S. & Lawrence, D. Transformation of  
736 dissolved and particulate materials on continental shelves influenced by  
737 large rivers: Plume processes. *Cont. Shelf Res.* **24**, 833–858 (2004).
- 738 10. Turner, R. E. & Rabalais, N. N. Linking landscape and water quality in the  
739 Mississippi river basin for 200 years. *Bioscience* **53**, 563–572 (2003).
- 740 11. Rabalais, N. N. *et al.* Nutrient Changes in the Mississippi River and System  
741 Responses on the Adjacent Continental Shelf. *Estuaries* **19**, 386 (1996).
- 742 12. Singh, V. *Application of Frequency and Risk in Water Resources:*  
743 *Proceedings of the International Symposium on Flood Frequency and Risk*  
744 *Analyses, 14–17 May 1986, Louisiana State University, Baton Rouge,*  
745 *USA.* (Springer Science & Business Media, 2012).
- 746 13. Dagg, M. J. *et al.* Biogeochemical characteristics of the lower Mississippi  
747 River, USA, during June 2003. *Estuaries* **28**, 664–674 (2005).
- 748 14. Turner, R. E. & Rabalais, N. N. Suspended sediment, C, N, P, and Si yields  
749 from the Mississippi River Basin. *Hydrobiologia* **511**, 79–89 (2004).
- 750 15. Schilling, K. E., Chan, K.-S., Liu, H. & Zhang, Y.-K. Quantifying the effect of  
751 land use land cover change on increasing discharge in the Upper  
752 Mississippi River. *J. Hydrol.* **387**, 343–345 (2010).
- 753 16. Staley, C. *et al.* Bacterial community structure is indicative of chemical  
754 inputs in the Upper Mississippi River. *Front Microbiol* **5**, 524 (2014).
- 755 17. McIsaac, G. F., David, M. B., Gertner, G. Z. & Goolsby, D. a. Nitrate flux in  
756 the Mississippi River. *Nature* **414**, 166–167 (2001).
- 757 18. Duan, S., Powell, R. T. & Bianchi, T. S. High frequency measurement of  
758 nitrate concentration in the Lower Mississippi River, USA. *J. Hydrol.* **519**,  
759 376–386 (2014).
- 760 19. Rabalais, N. N. *et al.* Hypoxia in the northern Gulf of Mexico: Does the  
761 science support the Plan to Reduce, Mitigate, and Control Hypoxia?  
762 *Estuaries and Coasts* **30**, 753–772 (2007).
- 763 20. Bianchi, T. S. *et al.* The science of hypoxia in the northern Gulf of Mexico:  
764 A review. *Science of the Total Environment* **408**, 1471–1484 (2010).
- 765 21. Rabalais, N. N., Turner, R. E. & Wiseman, W. J. Gulf of Mexico Hypoxia,  
766 a.K.a. ‘the Dead Zone’. *Annu. Rev. Ecol. Syst.* **33**, 235–263 (2002).

- 767 22. Bristow, L. A. *et al.* Biogeochemical and metagenomic analysis of nitrite  
768 accumulation in the Gulf of Mexico hypoxic zone. *Limnol. Oceanogr.* **60**,  
769 1733–1750 (2015).
- 770 23. Lemke, M. J. *et al.* Description of freshwater bacterial assemblages from  
771 the upper Paran river floodpulse system, Brazil. *Microb. Ecol.* **57**, 94–103  
772 (2009).
- 773 24. Christian Winter Gerhard Kavka, Robert L. Mach, and, T. H. & Farnleitner,  
774 A. H. Longitudinal Changes in the Bacterial Community Composition of the  
775 Danube rRiver: a Whole-River Approach. *AEM* (2007).
- 776 25. Kolmakova, O. V., Gladyshev, M. I., Rozanov, A. S., Peltek, S. E. &  
777 Trusova, M. Y. Spatial biodiversity of bacteria along the largest Arctic river  
778 determined by next-generation sequencing. *FEMS Microbiol. Ecol.* **89**,  
779 442–450 (2014).
- 780 26. Staley, C. *et al.* Core functional traits of bacterial communities in the Upper  
781 Mississippi River show limited variation in response to land cover. *Front.*  
782 *Microbiol.* **5**, 414 (2014).
- 783 27. Van Rossum, T. *et al.* Year-long metagenomic study of river microbiomes  
784 across land use and water quality . *Front. Microbiol.* **6** , 1–15  
785 (2015).
- 786 28. Meziti, A., Tsementzi, D., Ar. Kormas, K., Karayanni, H. & Konstantinidis,  
787 K. T. Anthropogenic effects on bacterial diversity and function along a river-  
788 to-estuary gradient in Northwest Greece revealed by metagenomics.  
789 *Environmental Microbiology* (2016). doi:10.1111/1462-2920.13303
- 790 29. Read, D. S. *et al.* Catchment-scale biogeography of riverine  
791 bacterioplankton. *ISME J* **9**, 516–526 (2015).
- 792 30. Blanchet, M. *et al.* When riverine dissolved organic matter (DOM) meets  
793 labile DOM in coastal waters: changes in bacterial community activity and  
794 composition. *Aquat. Sci.* (2016). doi:10.1007/s00027-016-0477-0
- 795 31. Ruiz-González, C. *et al.* Differences in organic matter and bacterioplankton  
796 between sections of the largest Arctic river: Mosaic or continuum?. *Front.*  
797 *Microbiol.* **6**, 196–206 (2015).
- 798 32. Zeglin, L. H. Stream microbial diversity in response to environmental  
799 changes: review and synthesis of existing research. *Front. Microbiol.* **6**,  
800 454 (2015).
- 801 33. Newton, R. J. & McLellan, S. L. A unique assemblage of cosmopolitan  
802 freshwater bacteria and higher community diversity differentiate an  
803 urbanized estuary from oligotrophic Lake Michigan. *Front. Microbiol.* **6**, 1–  
804 13 (2015).
- 805 34. Zwart, G., Crump, B. C., Kamst-van Agterveld, M. P., Hagen, F. & Han, S.  
806 K. Typical freshwater bacteria: An analysis of available 16S rRNA gene  
807 sequences from plankton of lakes and rivers. *Aquat. Microb. Ecol.* **28**, 141–  
808 155 (2002).
- 809 35. Staley, C. *et al.* Application of Illumina next-generation sequencing to  
810 characterize the bacterial community of the Upper Mississippi River. *J.*

- 811 *Appl. Microbiol.* **115**, 1147–1158 (2013).
- 812 36. RL, V., GW, M., KW, C., JR, S. & CE, C. River continuum concept. *Can J*  
813 *Fish Aquat Sci* **37**, 130–137 (1980).
- 814 37. Fortunato, C. S., Herfort, L., Zuber, P., Baptista, A. M. & Crump, B. C.  
815 Spatial variability overwhelms seasonal patterns in bacterioplankton  
816 communities across a river to ocean gradient. *ISME J* **6**, 554–563 (2012).
- 817 38. Besemer, K. *et al.* Headwaters are critical reservoirs of microbial diversity  
818 for fluvial networks. *Proc. Biol. Sci.* **280**, 20131760 (2013).
- 819 39. Jackson, C. R., Millar, J. J., Payne, J. T. & Ochs, C. a. Free-living and  
820 particle-associated bacterioplankton in large rivers of the Mississippi River  
821 Basin demonstrate biogeographic patterns. *Appl. Environ. Microbiol.* **80**,  
822 7186–7195 (2014).
- 823 40. Crump, B. C., Armbrust, E. V. & Baross, J. A. Phylogenetic Analysis of  
824 Particle-Attached and Free-Living Bacterial Communities in the Columbia  
825 River, Its Estuary, and the Adjacent Coastal Ocean. *Appl. Environ.*  
826 *Microbiol.* **65**, 3192–3204 (1999).
- 827 41. Riemann, L. & Winding, A. Community dynamics of free-living and particle-  
828 associated bacterial assemblages during a freshwater phytoplankton  
829 bloom. *Microb. Ecol.* **42**, 274–285 (2001).
- 830 42. D’Ambrosio, L., Ziervogel, K., Macgregor, B., Teske, A. & Arnosti, C.  
831 Composition and enzymatic function of particle-associated and free-living  
832 bacteria: a coastal/offshore comparison. *ISME J.* 1–13 (2014).  
833 doi:10.1038/ismej.2014.67
- 834 43. Allgaier, M. & Grossart, H.-P. Seasonal dynamics and phylogenetic  
835 diversity of free-living and particle-associated bacterial communities in four  
836 lakes in northeastern Germany. *Aquat. Microb. Ecol.* **45**, 115–128 (2006).
- 837 44. Crump, B. C. & Baross, J. A. Particle-associated bacteria and heterotrophic  
838 plankton associated with the Columbia River estuarine turbidity maxima.  
839 *Part. Bact. heterotrophic Plankt. Assoc. with Columbia River Estuar.*  
840 *Turbid. maxima* (1996).
- 841 45. Grossart, H. P. Ecological consequences of bacterioplankton lifestyles:  
842 Changes in concepts are needed. *Environmental Microbiology Reports* **2**,  
843 706–714 (2010).
- 844 46. Freimann, R., Bürgmann, H., Findlay, S. E. G. & Robinson, C. T.  
845 Hydrologic linkages drive spatial structuring of bacterial assemblages and  
846 functioning in alpine floodplains TL - 6. *Front. Microbiol.* **6 VN - re**, (2015).
- 847 47. Niño-García, J. P., Ruiz-González, C. & Giorgio, P. A. del. Interactions  
848 between hydrology and water chemistry shape bacterioplankton  
849 biogeography across boreal freshwater networks. *ISME J.* (2016).  
850 doi:10.1038/ismej.2015.226
- 851 48. Pierson, S. M., Rosenbaum, B. J., McKay, L. D. & Dewald, T. G. *Strahler*  
852 *Stream Order and Strahler Calculator Values in NHDPlus*. (2008).
- 853 49. Gladyshev, M. I. *et al.* Differences in organic matter and bacterioplankton  
854 between sections of the largest Arctic river: Mosaic or continuum? *Limnol.*

- 855 *Oceanogr.* **60**, 1314–1331 (2015).
- 856 50. Ghai, R. *et al.* Metagenomics of the water column in the pristine upper  
857 course of the Amazon river. *PLoS One* **6**, (2011).
- 858 51. Staley, C. *et al.* Bacterial community structure is indicative of chemical  
859 inputs in the Upper Mississippi River. *Front Microbiol* **5**, 524 (2014).
- 860 52. Staley, C. *et al.* Species sorting and seasonal dynamics primarily shape  
861 bacterial communities in the Upper Mississippi River. *Sci. Total Environ.*  
862 **505**, 435–45 (2015).
- 863 53. Crump, B. C., Hopkinson, C. S., Sogin, M. L. & Hobbie, J. E. Microbial  
864 Biogeography along an Estuarine Salinity Gradient: Combined Influences  
865 of Bacterial Growth and Residence Time. *Appl. Environ. Microbiol.* **70**,  
866 1494–1505 (2004).
- 867 54. Russell, T. A. & Weller, L. *State of the River Report: Water Quality and*  
868 *River Health in the Metro Mississippi River.* (Friends of the Mississippi  
869 River, 2013).
- 870 55. Karr, J. R. Biological integrity: a long-neglected aspect of water resource  
871 management. *Ecological Applications* **1**, 66–84 (1991).
- 872 56. Sims, A., Zhang, Y., Gajaraj, S., Brown, P. B. & Hu, Z. Toward the  
873 development of microbial indicators for wetland assessment. *Water*  
874 *Research* **47**, 1711–1725 (2013).
- 875 57. Karr, J. R. Assessment of biotic integrity using fish communities. *Fisheries*  
876 **6**, 21–27 (1981).
- 877 58. Thrash, J. C. & Coates, J. D. in *Bergey's Manual® of Systematic*  
878 *Bacteriology* 725–735 (Springer, 2010).
- 879 59. Kielak, A. M., Barreto, C. C., Kowalchuk, G. A., van Veen, J. A. &  
880 Kuramae, E. E. The Ecology of Acidobacteria: Moving beyond Genes and  
881 Genomes. *Front. Microbiol.* **7**, 744 (2016).
- 882 60. Anderson, I. *et al.* Genome sequence of the homoacetogenic bacterium  
883 *Holophaga foetida* type strain (TMBS4 T). *Stand. Genomic Sci.* **6**, 174  
884 (2012).
- 885 61. Allen, M. B. & Arnon, D. I. Studies on nitrogen-fixing blue-green algae. I.  
886 Growth and nitrogen fixation by *Anabaena cylindrica* Lemm. *Plant Physiol.*  
887 **30**, 366 (1955).
- 888 62. Van Geel, B., Mur, L. R., Ralska-Jasiewiczowa, M. & Goslar, T. Fossil  
889 akinetes of *Aphanizomenon* and *Anabaena* as indicators for medieval  
890 phosphate-eutrophication of Lake Gosciadz (Central Poland). *Rev.*  
891 *Palaeobot. Palynol.* **83**, 97–105 (1994).
- 892 63. Wood, S. A., Prentice, M. J., Smith, K. & Hamilton, D. P. Low dissolved  
893 inorganic nitrogen and increased heterocyte frequency: precursors to  
894 *Anabaena planktonica* blooms in a temperate, eutrophic reservoir. *J.*  
895 *Plankton Res.* **32**, 1315–1325 (2010).
- 896 64. Addison, S. L., Foote, S. M., Reid, N. M. & Lloyd-Jones, G.  
897 *Novosphingobium nitrogenifigens* sp. nov., a polyhydroxyalkanoate-  
898 accumulating diazotroph isolated from a New Zealand pulp and paper

- 899 wastewater. *Int. J. Syst. Evol. Microbiol.* **57**, 2467–2471 (2007).
- 900 65. Li, H.-F. *et al.* *Novosphingobium sediminis* sp. nov., isolated from the  
901 sediment of a eutrophic lake. *J. Gen. Appl. Microbiol.* **58**, 357–362 (2012).
- 902 66. Trusova, M. Y. & Gladyshev, M. I. Phylogenetic diversity of winter  
903 bacterioplankton of eutrophic siberian reservoirs as revealed by 16S rRNA  
904 gene sequence. *Microb. Ecol.* **44**, 252–259 (2002).
- 905 67. Cebron, A. & Garnier, J. Nitrobacter and Nitrospira genera as  
906 representatives of nitrite-oxidizing bacteria: detection, quantification and  
907 growth along the lower Seine River (France). *Water Res.* **39**, 4979–4992  
908 (2005).
- 909 68. Korajkic, A. *et al.* Changes in bacterial and eukaryotic communities during  
910 sewage decomposition in Mississippi river water TL - 69. *Water Res.* **69**  
911 **VN - r**, 30–39 (2015).
- 912 69. Necchi Jr, O. in *River Algae* 153–158 (Springer, 2016).
- 913 70. Jansson, M., Blomqvist, P., Jonsson, A. & Bergström, A. Nutrient limitation  
914 of bacterioplankton, autotrophic and mixotrophic phytoplankton, and  
915 heterotrophic nanoflagellates in Lake Öträsket. *Limnol. Oceanogr.* **41**,  
916 1552–1559 (1996).
- 917 71. Caron, D. A., Porter, K. G. & Sanders, R. W. Carbon, nitrogen, and  
918 phosphorus budgets for the mixotrophic phytoflagellate *Poterioochromonas*  
919 *malhamensis* (Chrysophyceae) during bacterial ingestion. *Limnol.*  
920 *Oceanogr.* **35**, 433–443 (1990).
- 921 72. Smith, M. W. *et al.* Seasonal Changes in Bacterial and Archaeal Gene  
922 Expression Patterns across Salinity Gradients in the Columbia River  
923 Coastal Margin. *PLoS One* **5**, e13312 (2010).
- 924 73. Fortunato, C. S. *et al.* Determining indicator taxa across spatial and  
925 seasonal gradients in the Columbia River coastal margin. *ISME J* **7**, 1899–  
926 1911 (2013).
- 927 74. Vorosmarty, C. J. *et al.* Global threats to human water security and river  
928 biodiversity. *Nature* **467**, 555–561 (2010).
- 929 75. Weiss, S. *et al.* Tracking down the sources of experimental contamination  
930 in microbiome studies. *Genome Biol.* **15**, 564 (2014).
- 931 76. Salter, S. J. *et al.* Reagent and laboratory contamination can critically  
932 impact sequence-based microbiome analyses. *BMC Biol.* **12**, 87 (2014).
- 933 77. Thrash, J. C., Weckhorst, J. L. & Pitre, D. M. in *Protocols for Metagenomic*  
934 *Library Generation and Analysis in Petroleum Hydrocarbon Microbe*  
935 *Systems* 1–22 (Humana Press, 2015). doi:10.1007/8623\_2015\_67
- 936 78. Caporaso, J. G. *et al.* Ultra-high-throughput microbial community analysis  
937 on the Illumina HiSeq and MiSeq platforms. *ISME J.* **6**, 1621–1624 (2012).
- 938 79. Amaral-Zettler, L. A., McCliment, E. A., Ducklow, H. W. & Huse, S. M. A  
939 Method for Studying Protistan Diversity Using Massively Parallel  
940 Sequencing of V9 Hypervariable Regions of Small-Subunit Ribosomal RNA  
941 Genes. *PLoS One* **4**, e6372 (2009).
- 942 80. Schloss, P. D. *et al.* Introducing mothur: Open-source, platform-

- 943 independent, community-supported software for describing and comparing  
944 microbial communities. *Appl. Environ. Microbiol.* **75**, 7537–7541 (2009).
- 945 81. Quast, C. *et al.* The SILVA ribosomal RNA gene database project:  
946 Improved data processing and web-based tools. *Nucleic Acids Res.* **41**,  
947 590–596 (2013).
- 948 82. Pruesse, E. *et al.* SILVA: a comprehensive online resource for quality  
949 checked and aligned ribosomal RNA sequence data compatible with ARB.  
950 *Nucleic Acids Res.* **35**, 7188–7196 (2007).
- 951 83. Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C. & Knight, R.  
952 UCHIME improves sensitivity and speed of chimera detection.  
953 *Bioinformatics* **27**, 2194–2200 (2011).
- 954 84. Oksanen, J. *et al.* vegan: Community Ecology Package. R package version  
955 2.2-1. *R package version 1*, R package version 2.2–1. (2015).
- 956 85. R, C. T. *R: A language and environment for statistical computing* (2013).
- 957 86. McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible  
958 Interactive Analysis and Graphics of Microbiome Census Data. *PLoS One*  
959 **8**, e61217 (2013).
- 960 87. Love, M. I., Anders, S. & Huber, W. Differential analysis of count data - the  
961 DESeq2 package. *Genome Biol.* **15**, 550 (2014).
- 962 88. Learman, D. R. *et al.* Biogeochemical and microbial variation across 5500  
963 km of Antarctic surface sediment implicates organic matter as a driver of  
964 benthic community structure. *Front. Microbiol.* **7**, (2016).
- 965 89. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted gene co-  
966 expression network analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 967 90. Guidi, L. *et al.* Plankton networks driving carbon export in the oligotrophic  
968 ocean. *Nature* in review (2015). doi:10.1038/nature16942
- 969 91. Chong, I. G. & Jun, C. H. Performance of some variable selection methods  
970 when multicollinearity is present. *Chemom. Intell. Lab. Syst.* **78**, 103–112  
971 (2005).
- 972 92. Mevik, B.-H. & Wehrens, R. The pls Package: Principle Component and  
973 Partial Least Squares Regression in R. *J. Stat. Softw.* **18**, 1–24 (2007).

974

975

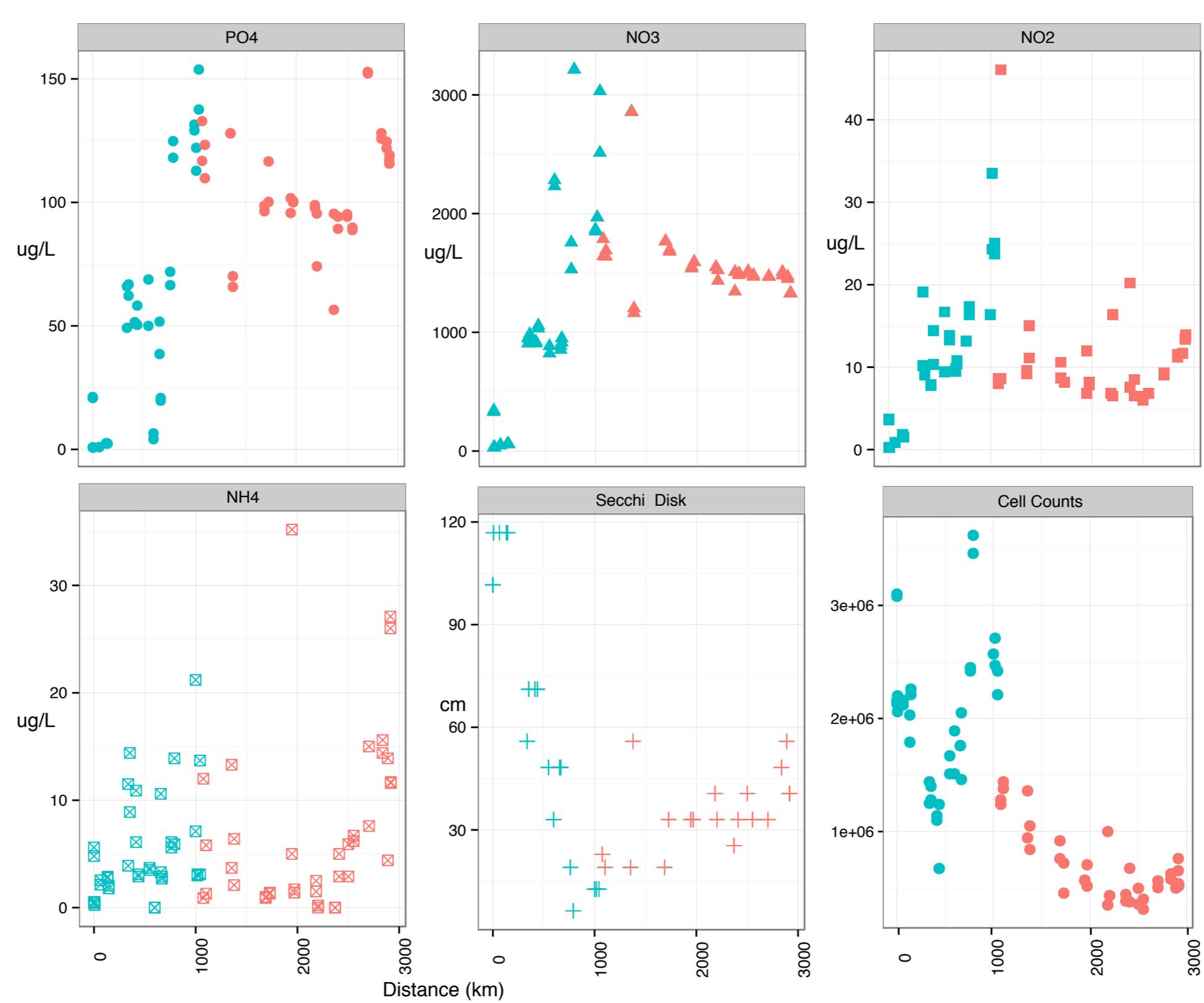
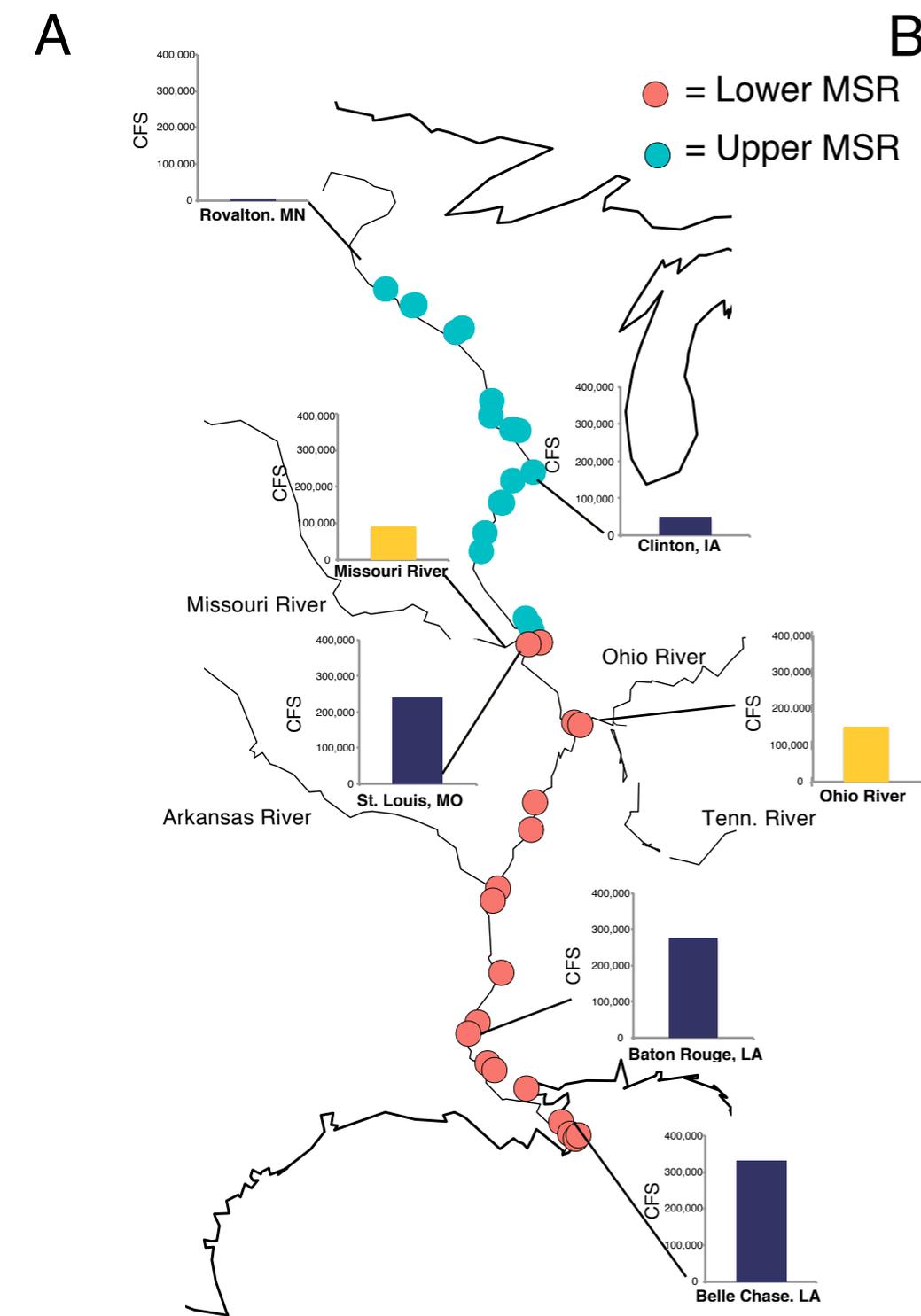
## 976 **Figure Legends**

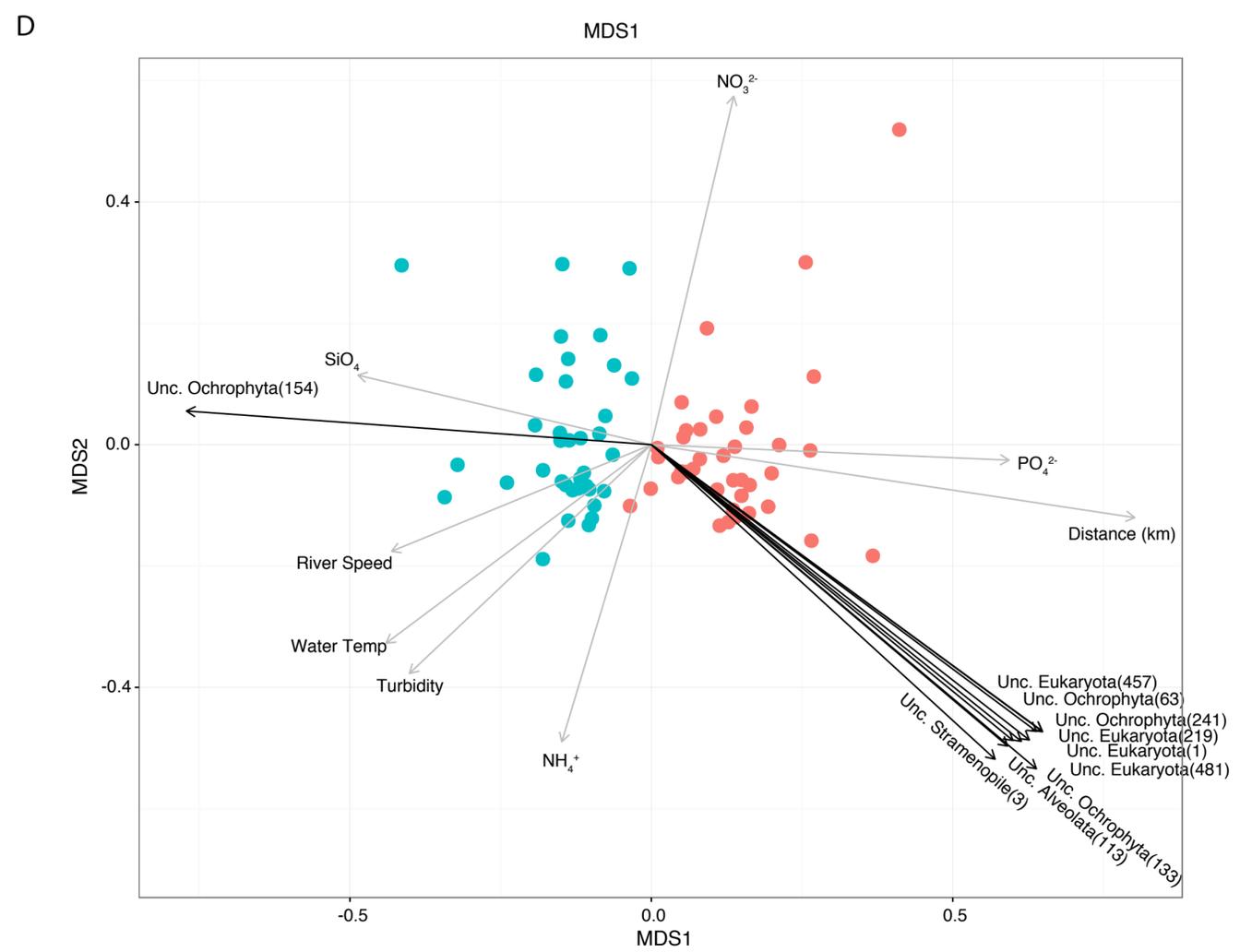
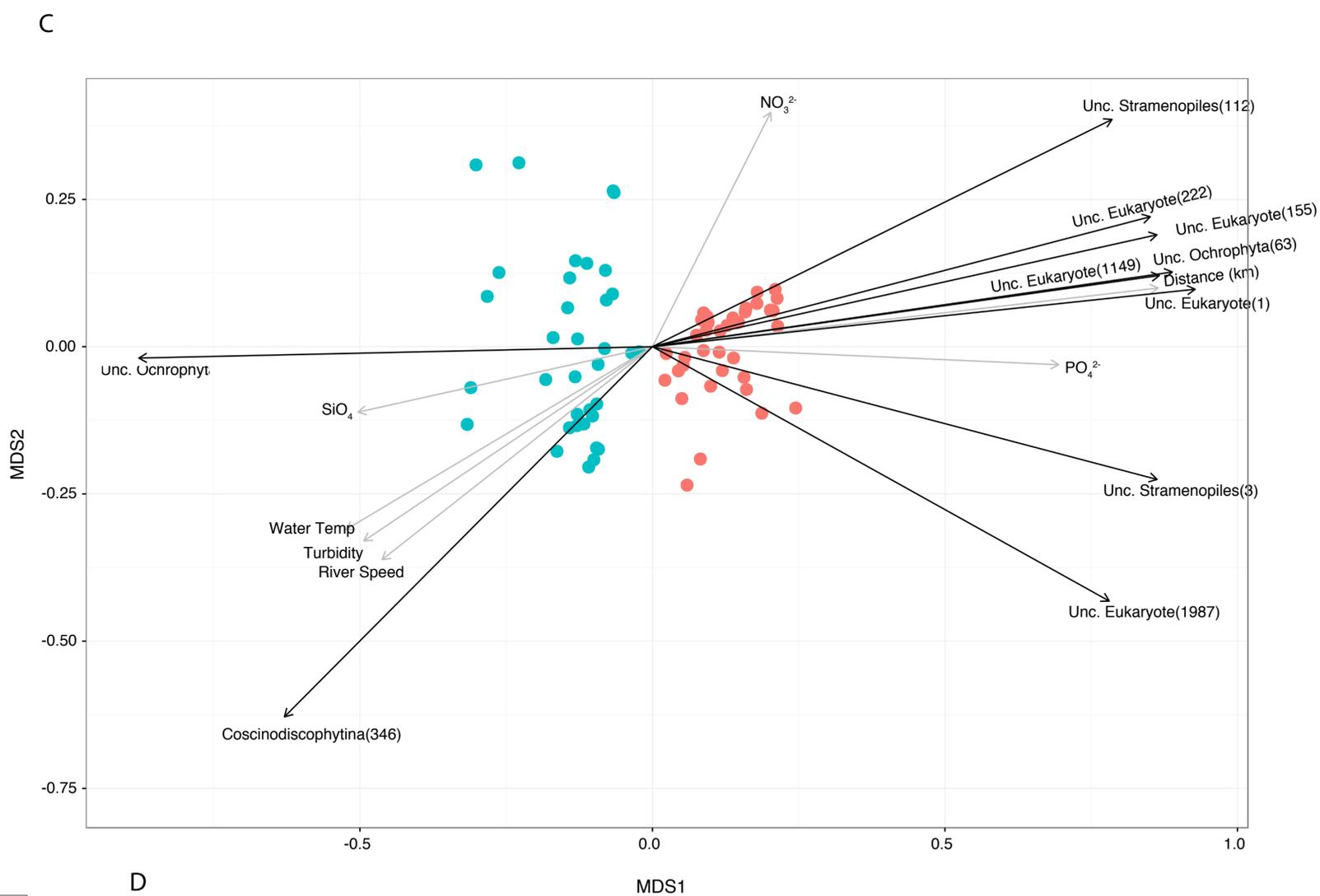
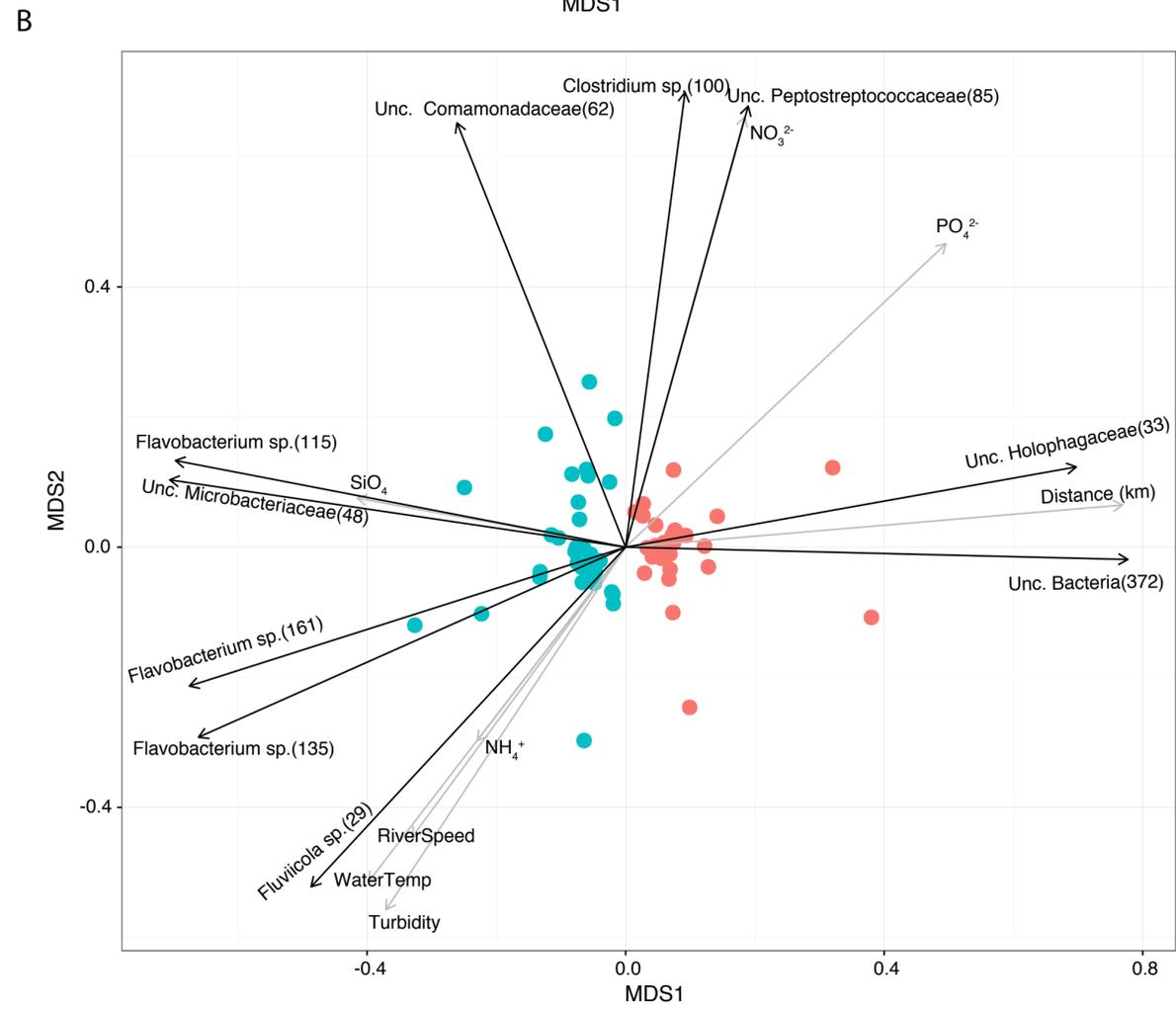
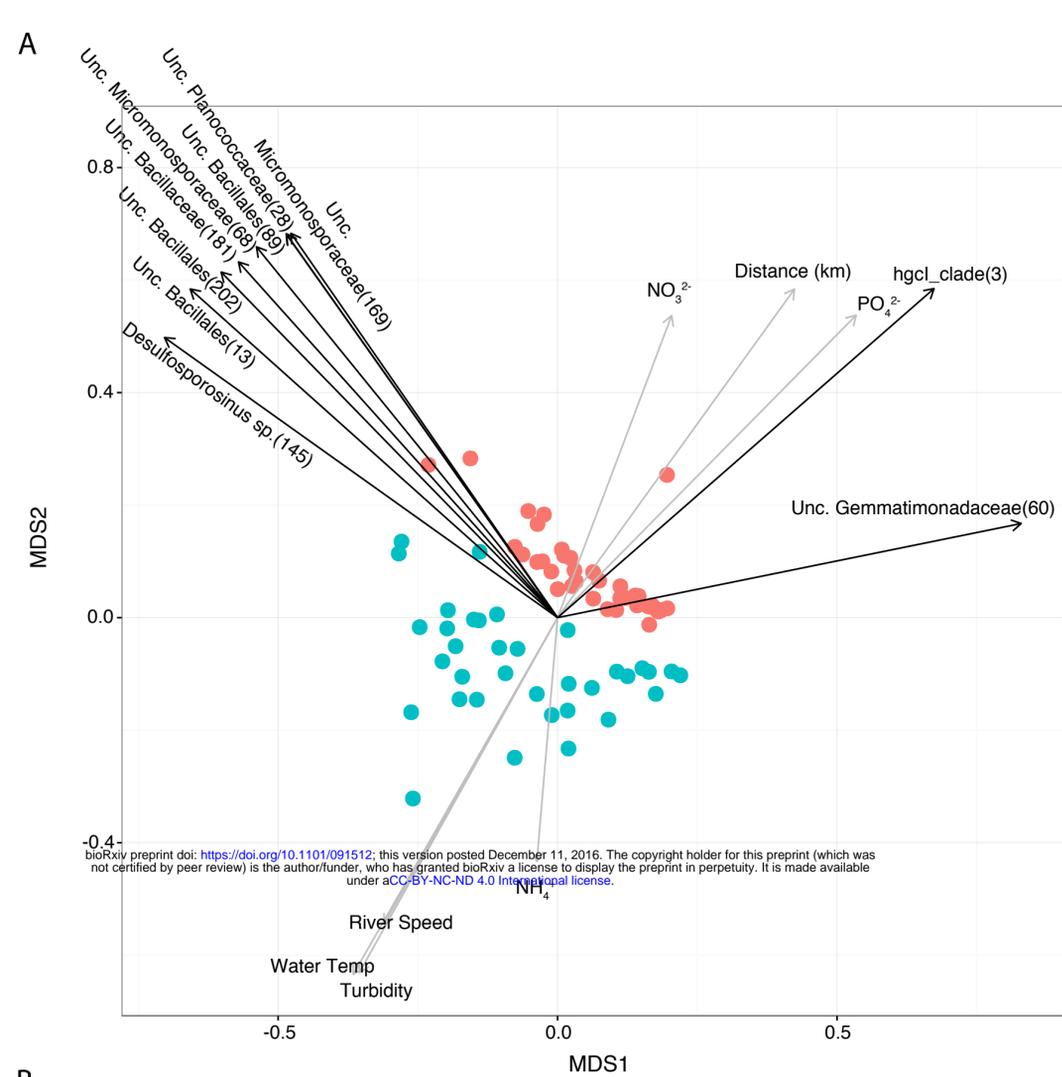
977

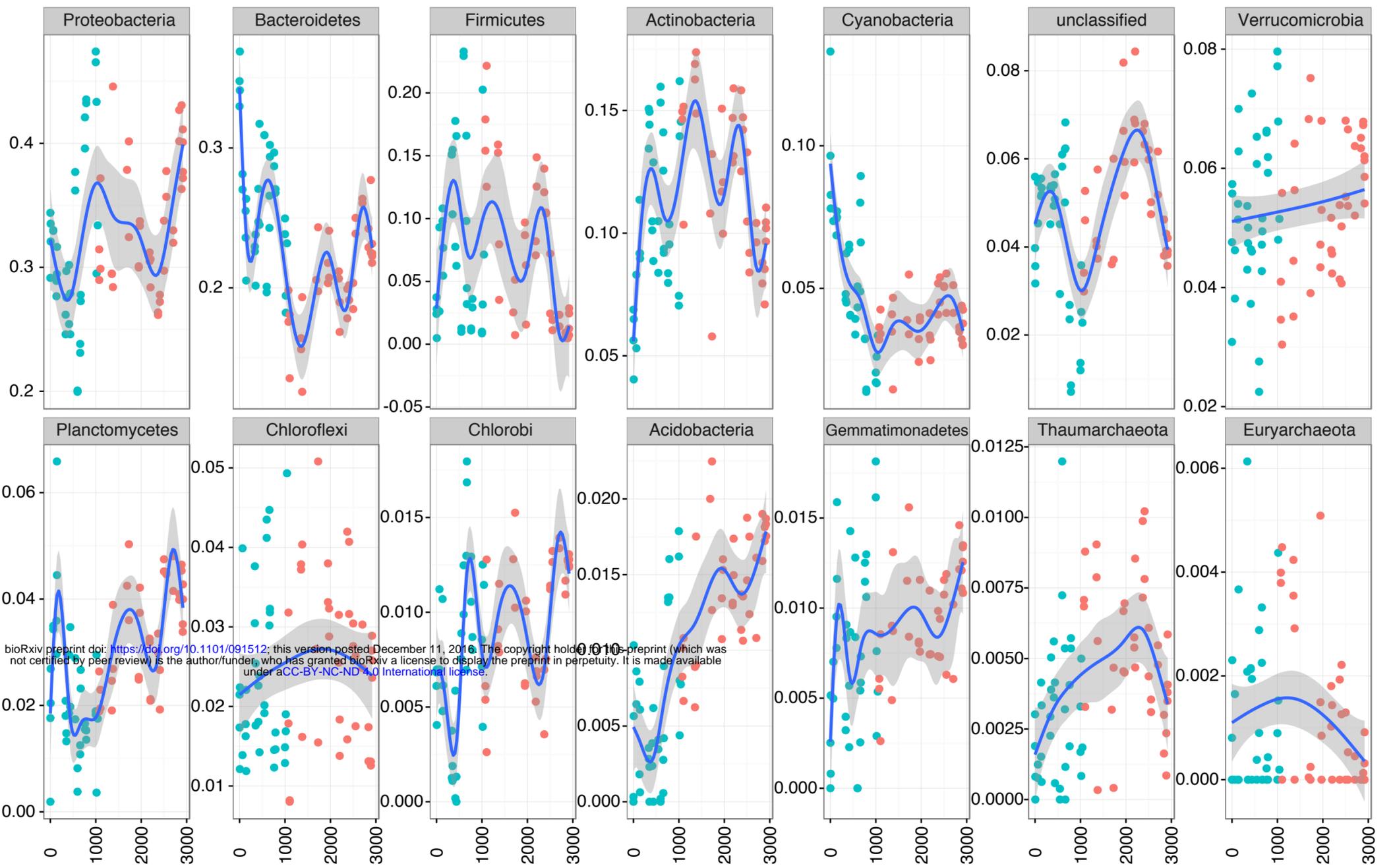
978 **Figure 1.** Sampling map (A) with graph inserts that represent the measured  
979 discharge rate (cubic feet second<sup>-1</sup> [CFS]) as recorded on the USGS  
980 gauges, and six environmental parameters measured along the transect,  
981 according to concentration or visible depth of secchi disk by distance (B).  
982 Throughout the figure and text, blue and red dots represent sampling  
983 locations above and below the Missouri River confluence, respectively, and  
984 are designated throughout as “Upper” and “Lower.” Cell Counts only  
985 represent the < 2.7  $\mu\text{m}$  fraction. C) An OAR Northwest rowboat dwarfed by  
986 a paddleboat.

- 987 **Figure 2.** Non-metric multidimensional scaling (NMDS) results for whole  
988 community correlations with environmental parameters (phosphate, nitrate,  
989 nitrite, ammonia, distance (km), water temperature, turbidity (cm), river  
990 speed (mph)) and the top ten OTUs based on significance (P) and strength  
991 of correlation (r). The four plots represent > 2.7  $\mu\text{m}$  (A, C) and 0.2-2.7  $\mu\text{m}$   
992 (B, D) fractions for the 16S (A, B) and 18S (C, D) rRNA gene communities.  
993 Vector length is proportional to the strength of the correlation.
- 994 **Figure 3.** Relative abundance, by phylum, according to transect distance, for  
995 phyla accounting for > 0.1% of the total reads for the 16S rRNA gene > 2.7  
996  $\mu\text{m}$  (A) and 0.2-2.7  $\mu\text{m}$  (B) communities. Non-linear regressions with 95%  
997 CI (gray shading) are provided for reference.
- 998 **Figure 4.** Core microbiome aggregate abundance for the 16S (A) and 18S (B)  
999 rRNA gene. In each, triangles and circles points represent 0.2-2.7  $\mu\text{m}$  and  
1000 > 2.7  $\mu\text{m}$  fractions, respectively. Non-linear regressions with 95%  
1001 confidence intervals (CI) (gray shading) are provided for reference.
- 1002 **Figure 5.** Relative abundance, by phylum, according to transect distance, for  
1003 phyla accounting for > 0.1% of the total reads for the 18S rRNA gene > 2.7  
1004  $\mu\text{m}$  (A) and 0.2-2.7  $\mu\text{m}$  (B) communities. Non-linear regressions with 95%  
1005 CI (gray shading) are provided for reference.
- 1006 **Figure 6.** PLS results for the 0.2-2.7  $\mu\text{m}$  16S rRNA gene community for selected  
1007 submodules with nitrate and phosphate and a VIP score > 1. Correlation of  
1008 submodule OTUs to nitrate (A) and phosphate (B), according to the  
1009 number of co-correlations (node centrality). Circle size is proportional to  
1010 VIP scores, with top 10 VIP scoring and top node centrality OTUs labeled  
1011 with their highest-resolution taxonomic classification and OTU number.  
1012 Colors represent the taxonomic classification the phylum level.
- 1013 **Figure 7.** PLS results for the > 2.7  $\mu\text{m}$  16S rRNA gene community for selected  
1014 submodules with nitrate and phosphate and a VIP score of > 1. Correlation  
1015 of submodule OTUs to nitrate (A) and phosphate (B), according to the  
1016 number of co-correlations (node centrality). Circle size is proportional to  
1017 VIP scores, with top 10 VIP scoring and top node centrality OTUs labeled  
1018 with their highest-resolution taxonomic classification and OTU number.  
1019 Colors represent the taxonomic classification the phylum level.
- 1020 **Figure 8.** PLS results for the 0.2-2.7  $\mu\text{m}$  16S rRNA gene community for selected  
1021 submodules with nitrate and phosphate and a VIP score of > 1. Correlation  
1022 of submodule OTUs to nitrate (A) and phosphate (B), according to the  
1023 number of co-correlations (node centrality). Circle size is proportional to  
1024 VIP scores, with top 10 VIP scoring and top node centrality OTUs labeled  
1025 with their highest-resolution taxonomic classification and OTU number.  
1026 Colors represent the taxonomic classification the phylum level.
- 1027 **Figure 9.** PLS results for the > 2.7  $\mu\text{m}$  18S rRNA gene community for selected  
1028 submodules with nitrate and phosphate and a VIP score of > 1. Correlation  
1029 of submodule OTUs to nitrate (A) and phosphate (B), according to the  
1030 number of co-correlations (node centrality). Circle size is proportional to

1031 VIP scores, with top 10 VIP scoring and top node centrality OTUs labeled  
1032 with their highest-resolution taxonomic classification and OTU number.  
1033 Colors represent the taxonomic classification the phylum level.

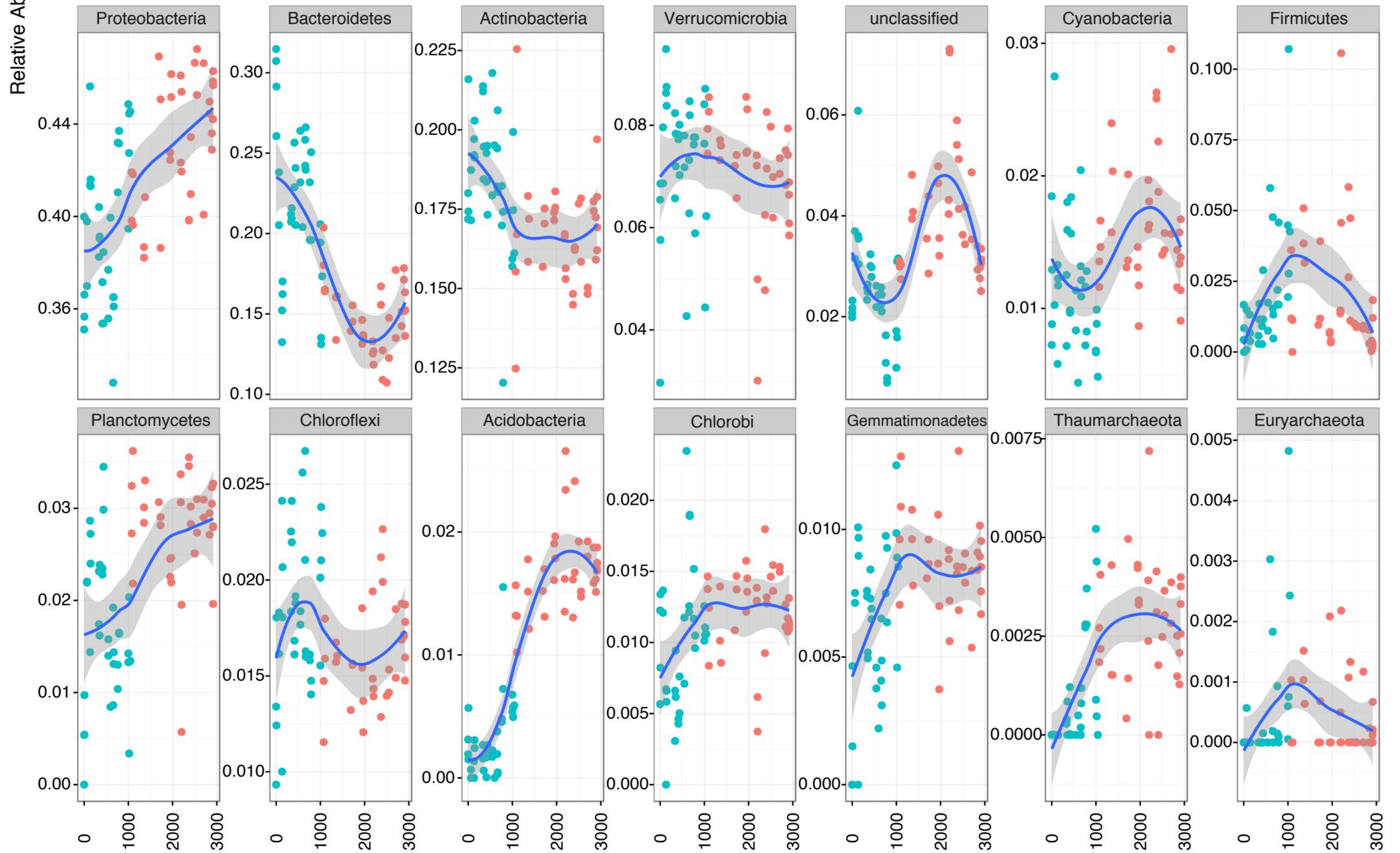


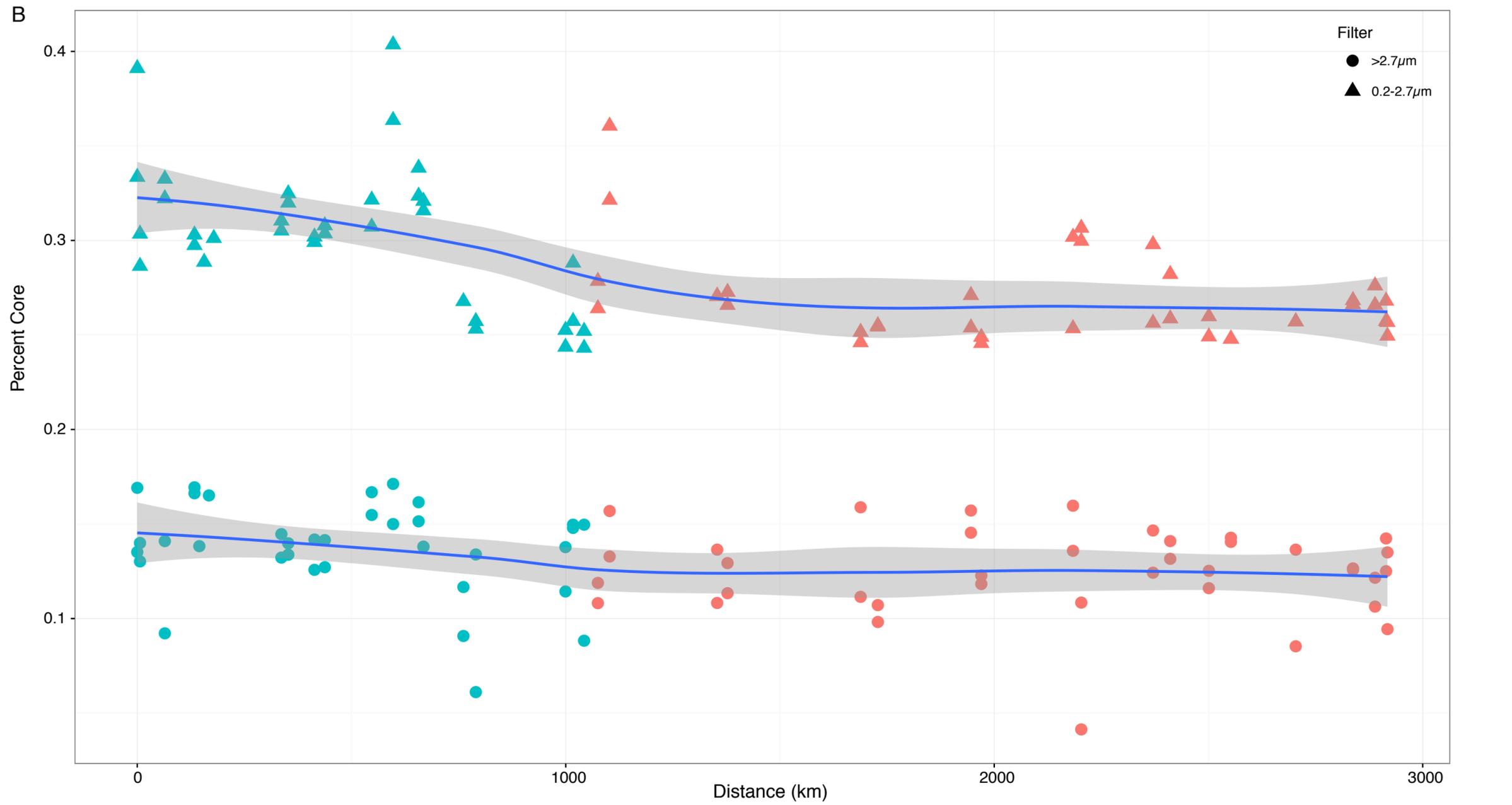
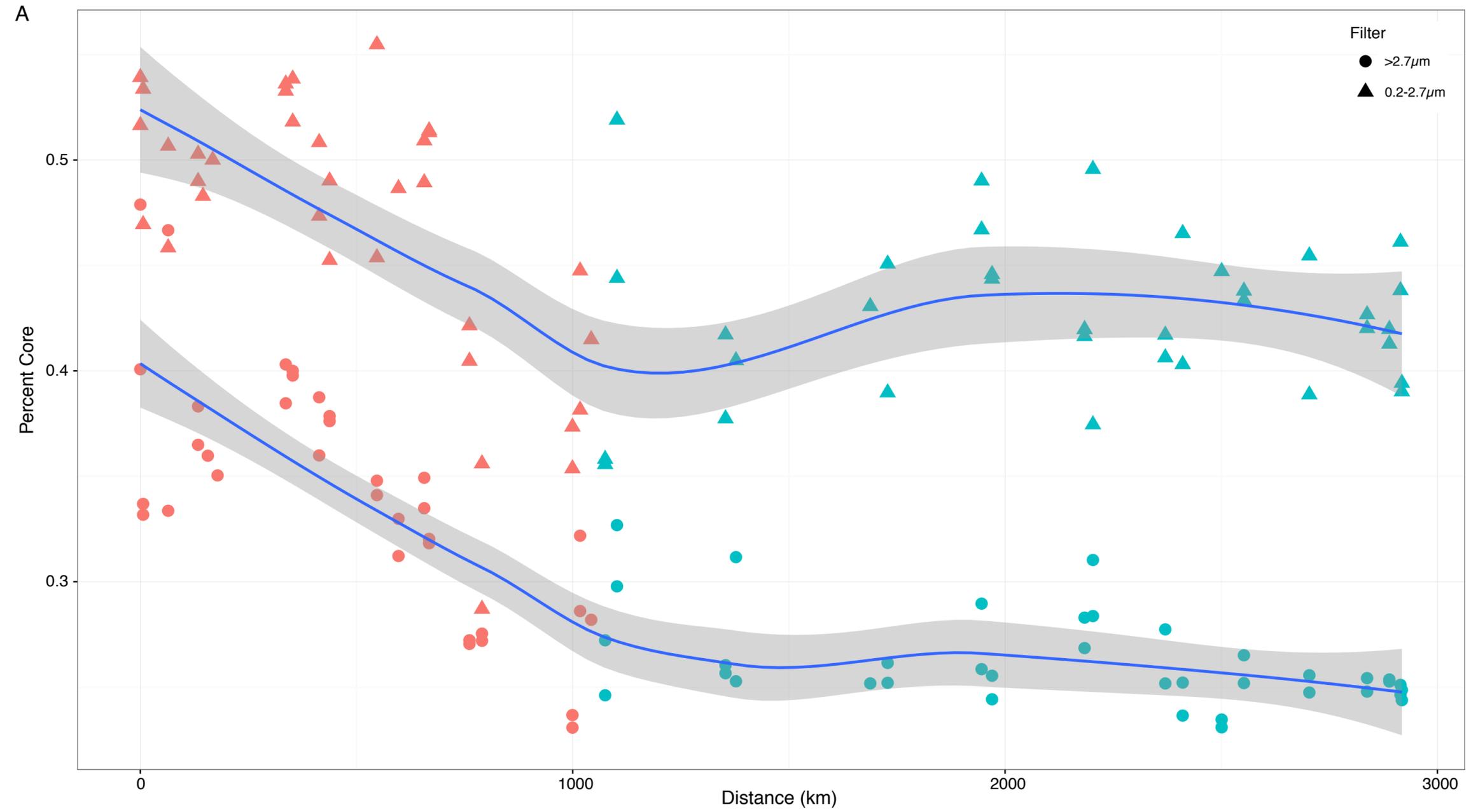




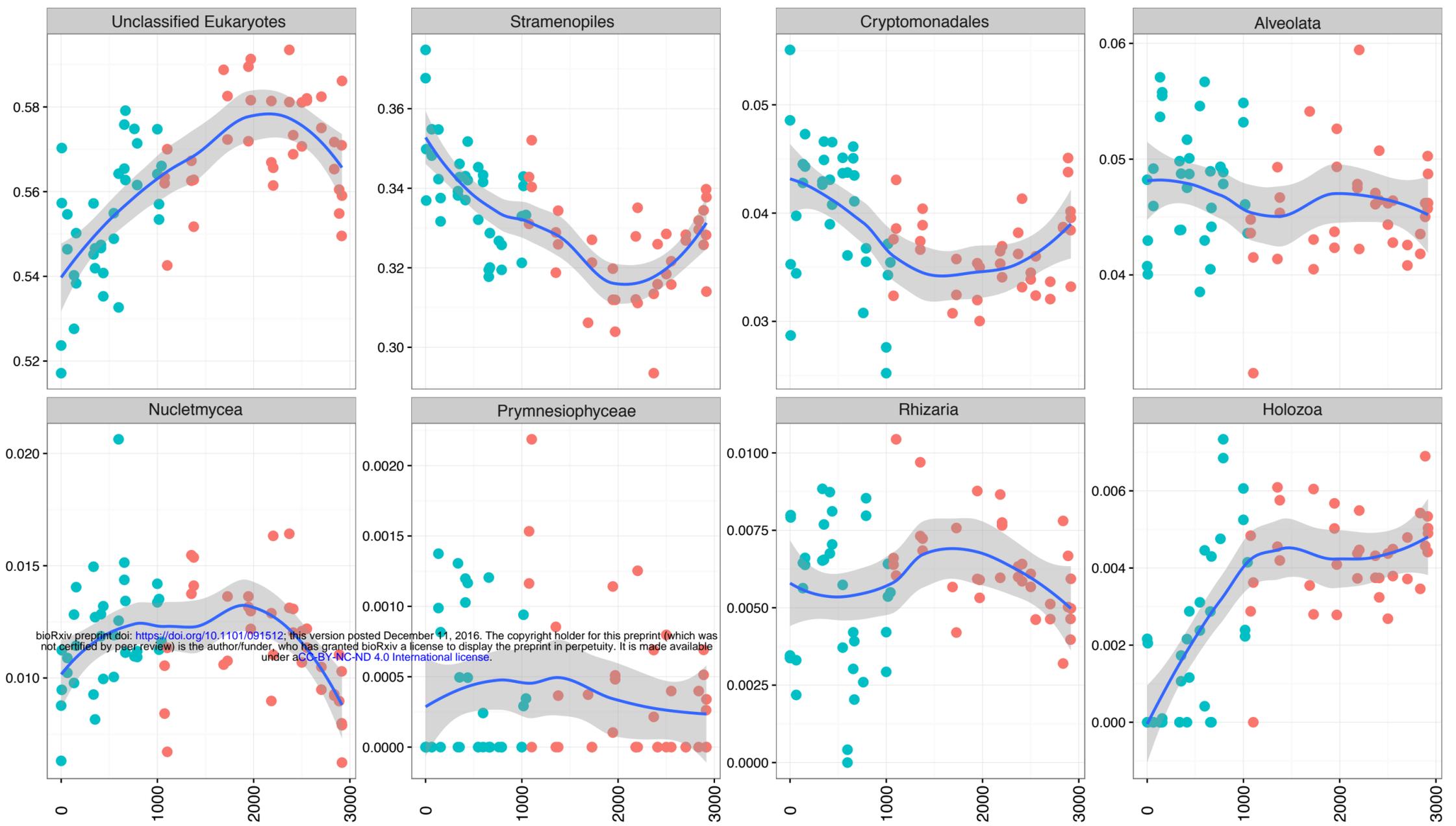
bioRxiv preprint doi: <https://doi.org/10.1101/091512>; this version posted December 11, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

B

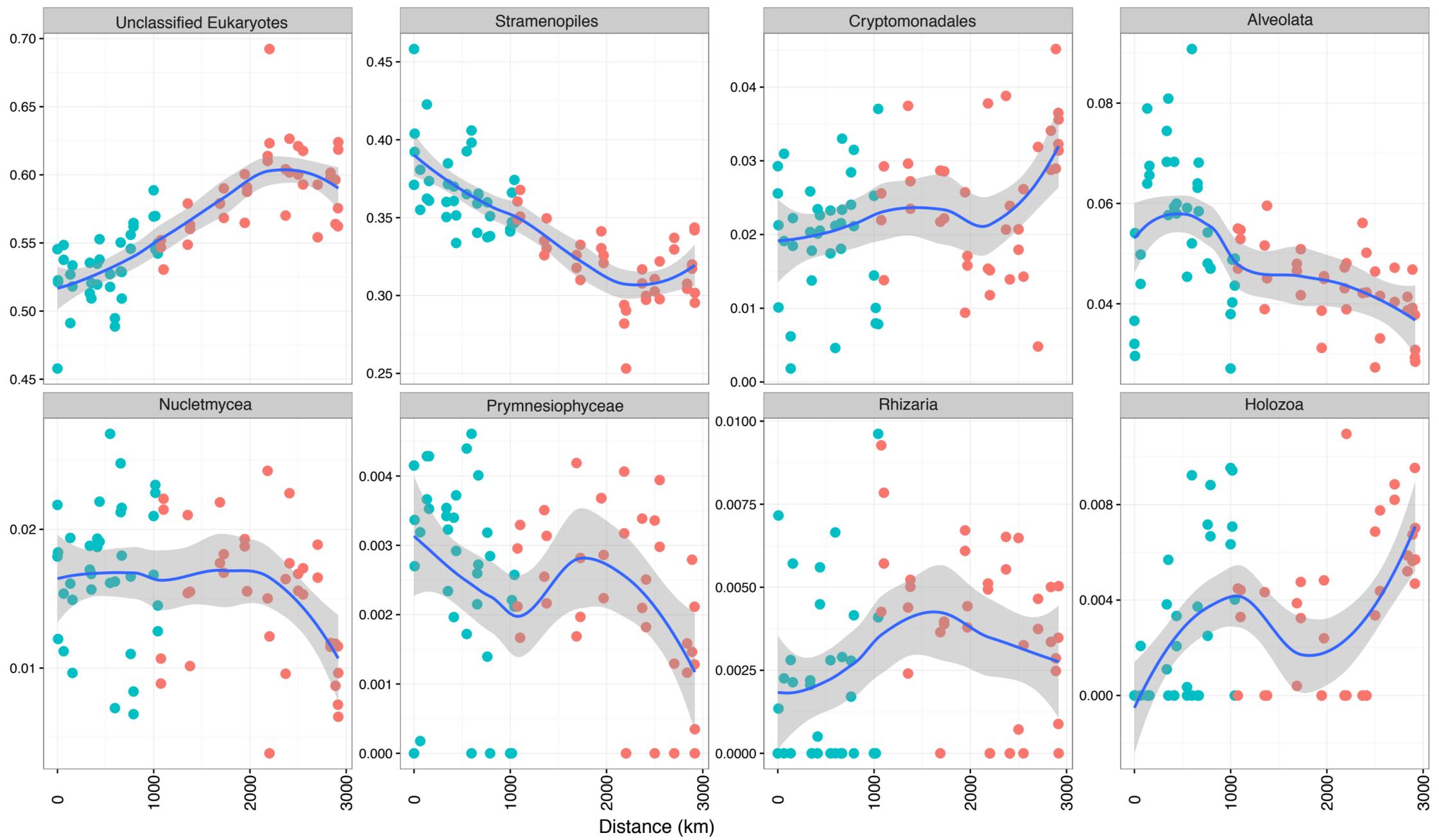


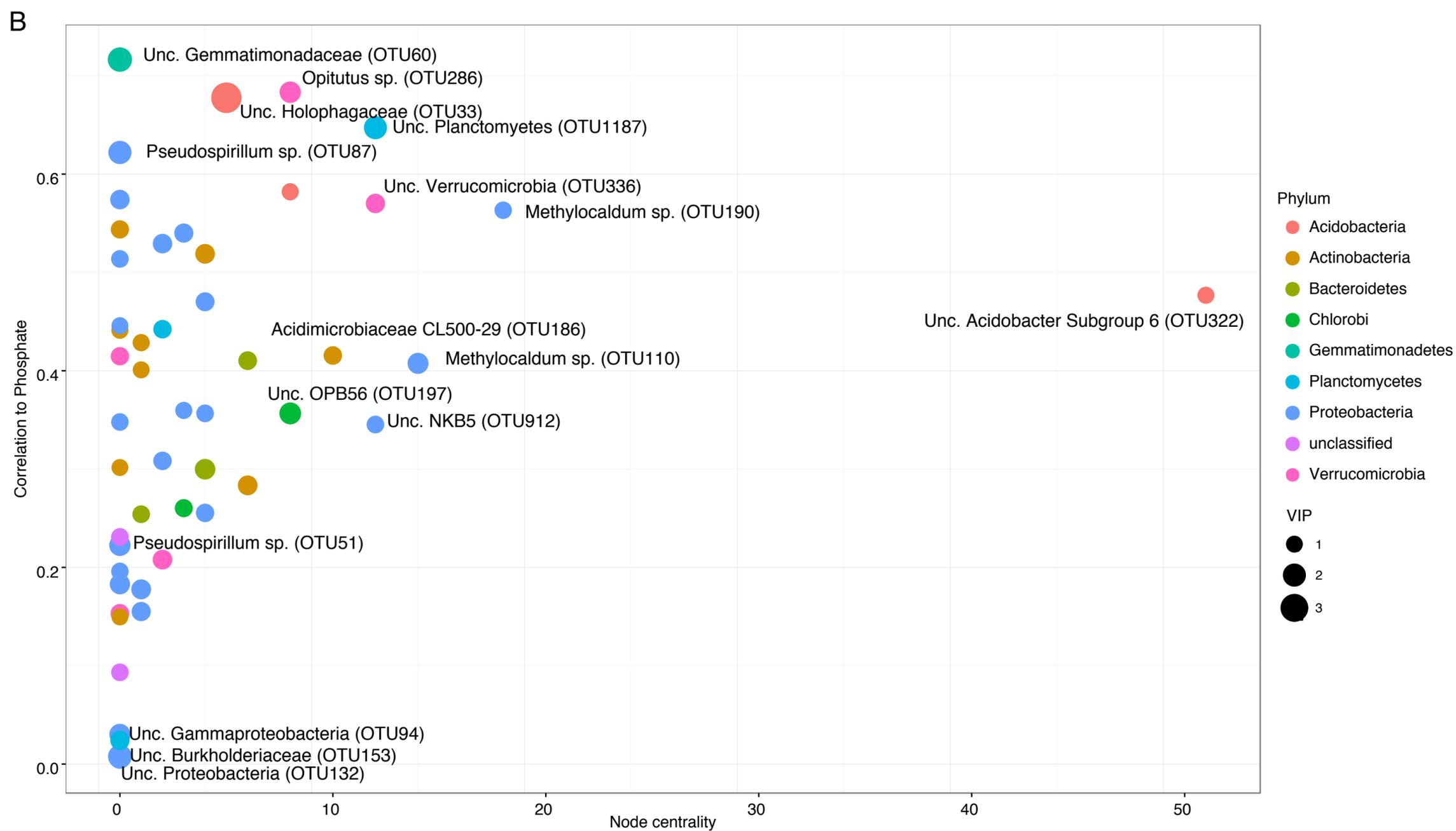
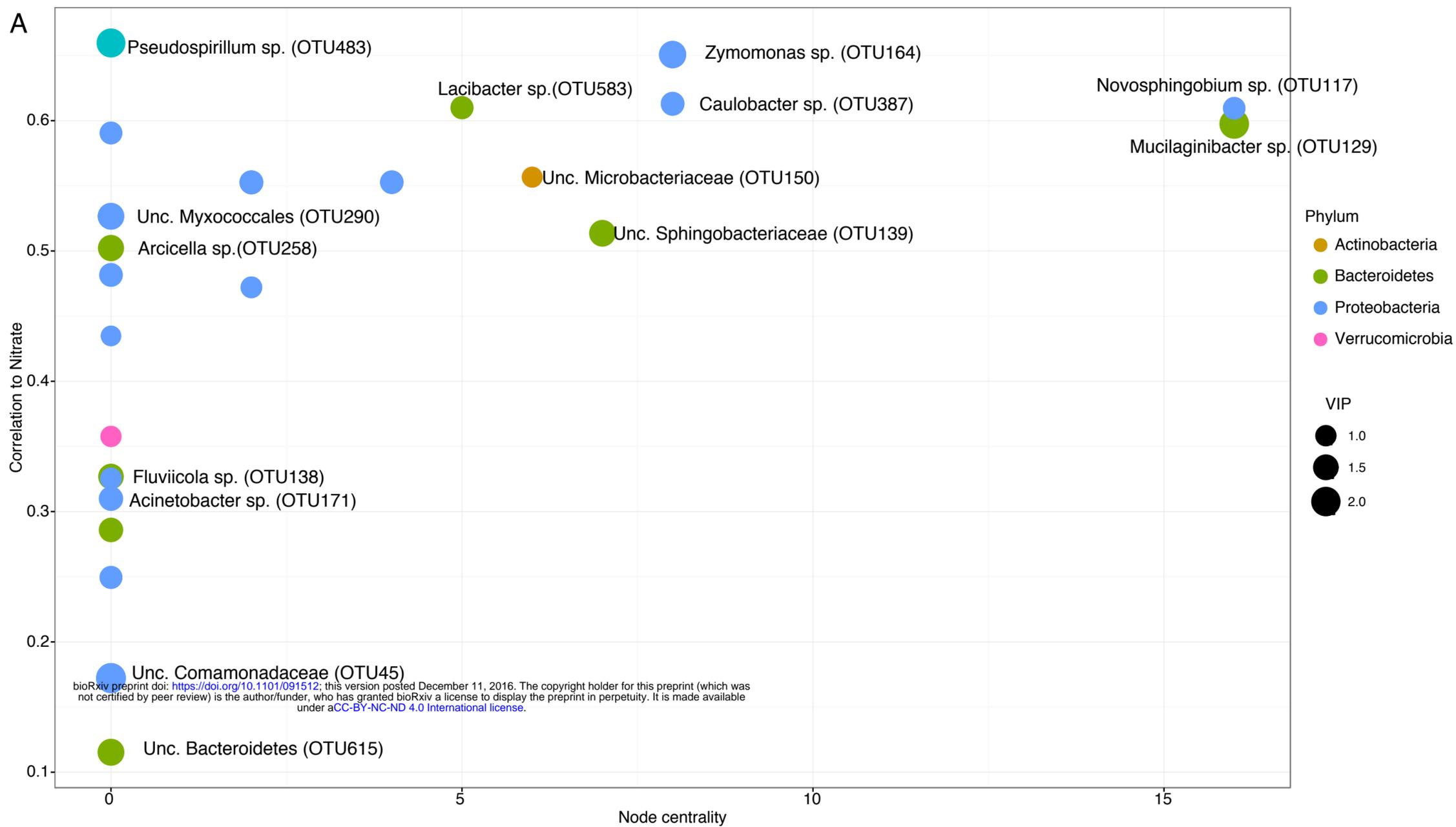


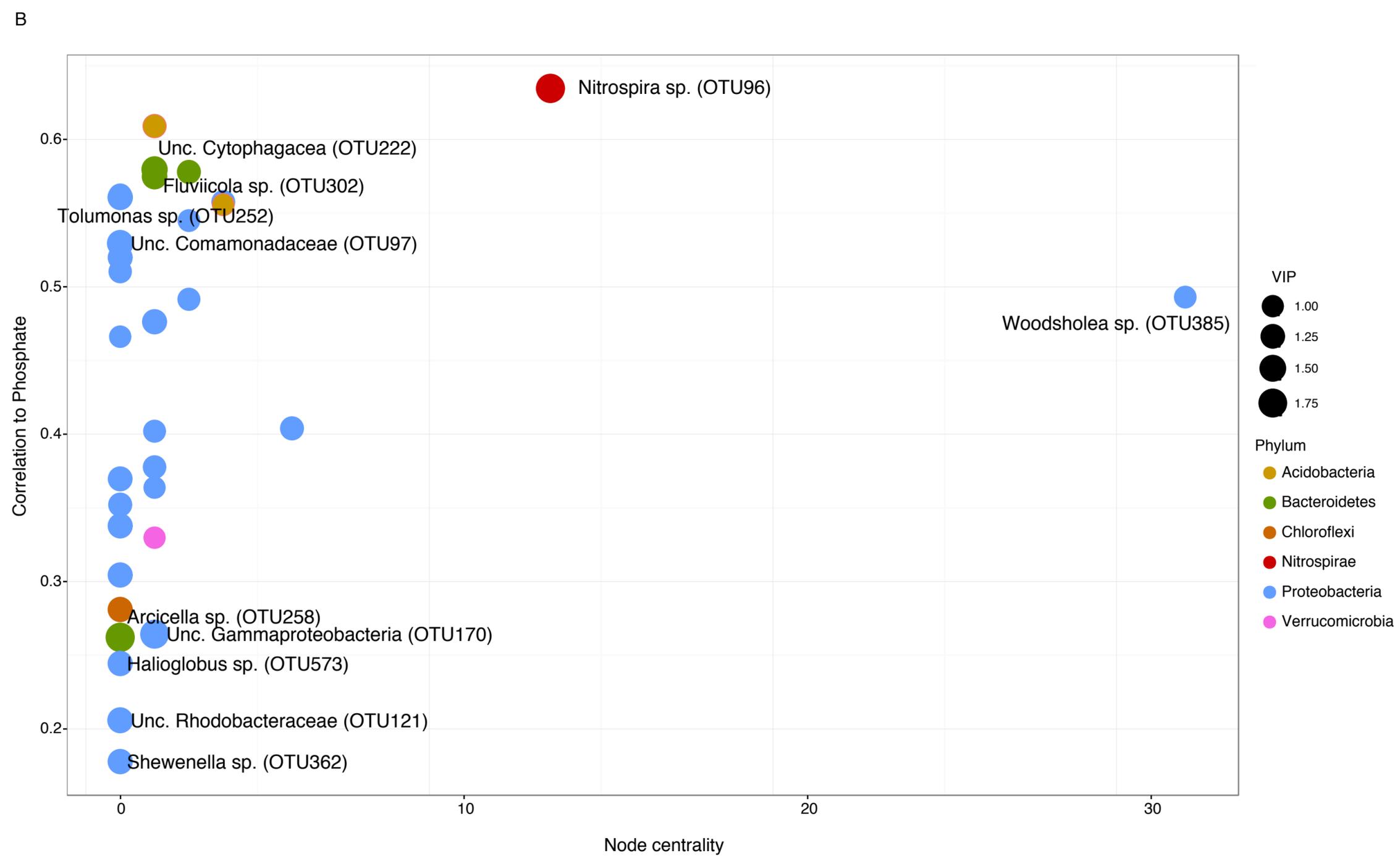
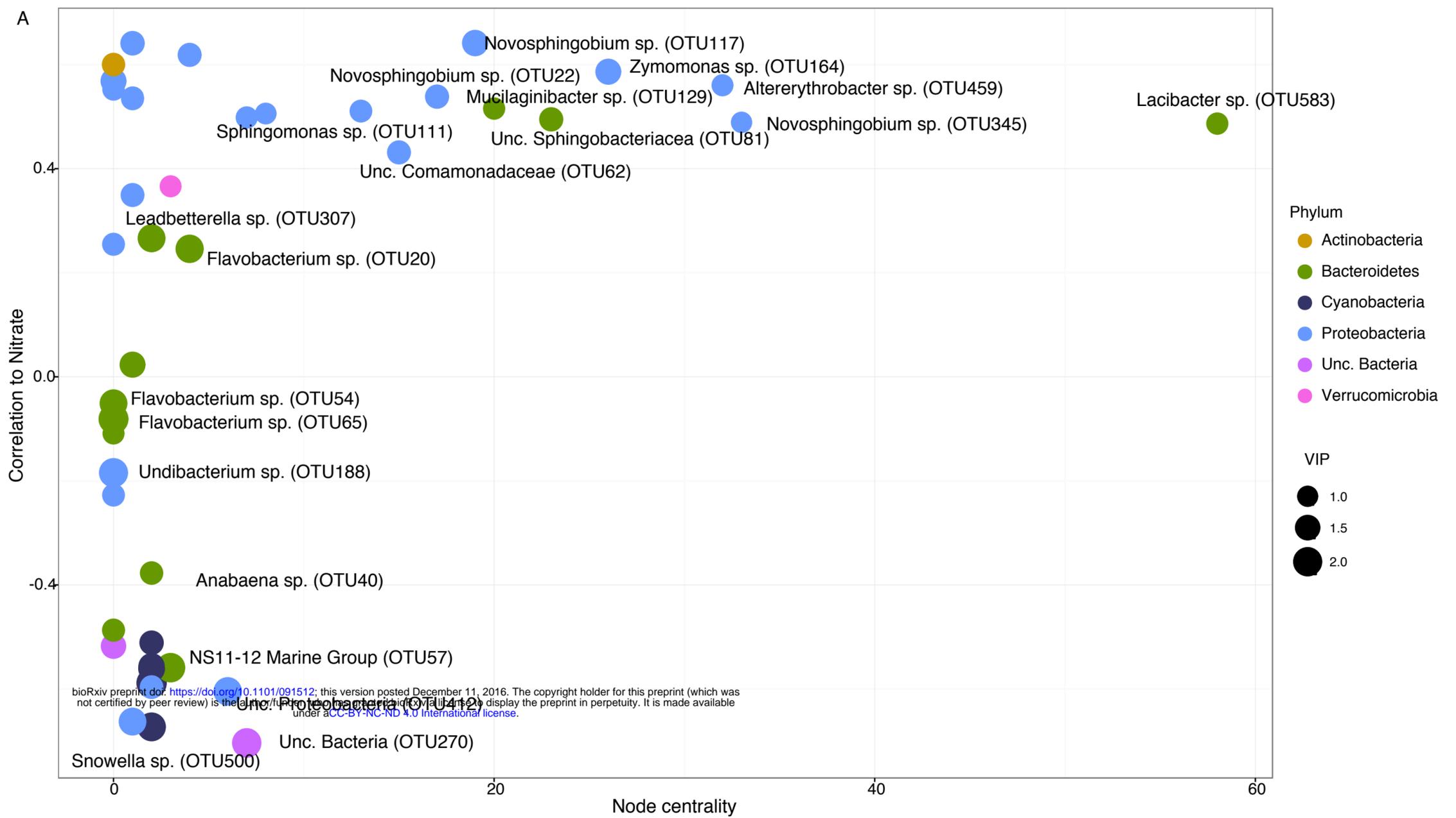
A

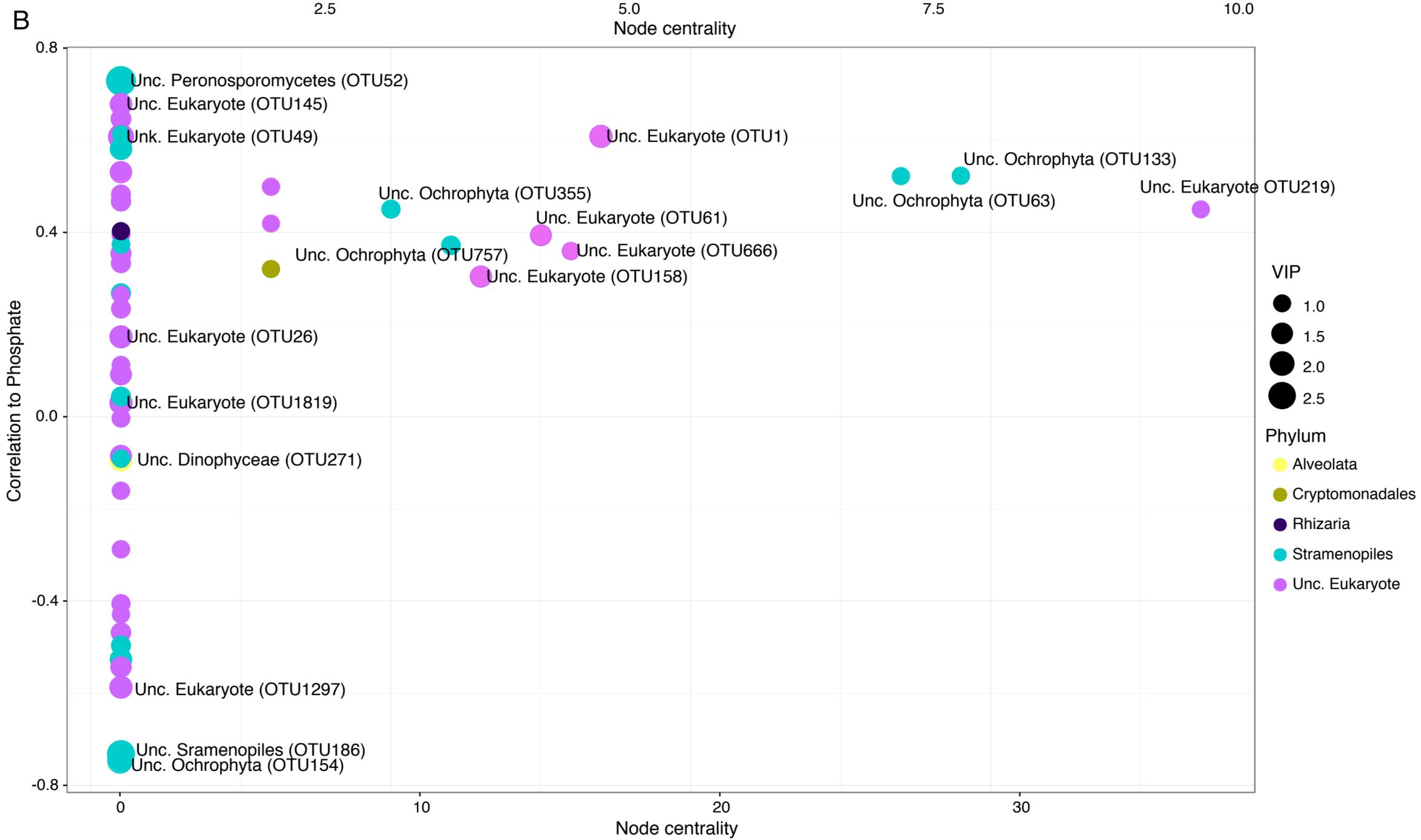
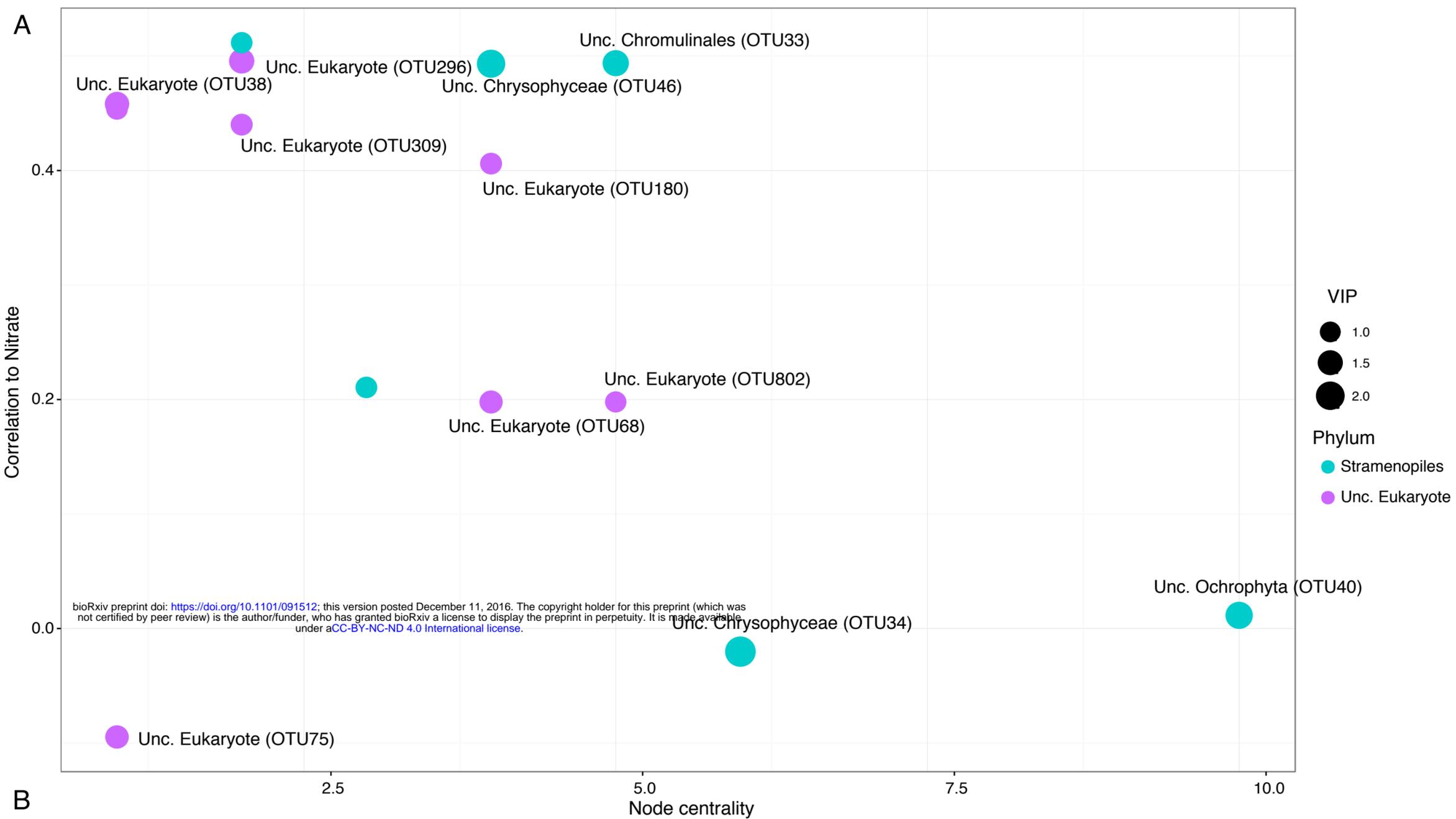


B









A



B

