

1 **Using machine learning to associate bacterial taxa with functional groups through flow**
2 **cytometry, 16S rRNA gene sequencing, and productivity data**

3

4 Peter Rubbens^{1*}, Marian L. Schmidt^{2*#}, Ruben Props³, Bopaiah A. Biddanda⁴, Nico Boon³,
5 Willem Waegeman¹, Vincent J. Denef²

6

7 *Peter Rubbens and Marian L. Schmidt contributed equally to this work.

8 #Corresponding author: marschmi@umich.edu

9

10 ¹KERMIT, Department of Data Analysis and Mathematical Modelling, Ghent University,
11 Coupure Links 653, B-9000, Gent, Belgium; ²Department of Ecology and Evolutionary Biology,
12 University of Michigan, 830 North University Ave., Ann Arbor, MI 48109, USA; ³CMET,
13 Center for Microbial Ecology and Technology, Ghent University, Coupure Links 653, B-9000,
14 Gent, Belgium; ⁴Annis Water Resources Institute, Grand Valley State University, 740 West
15 Shoreline Drive, Muskegon, MI 49441, USA

16 **Abstract**

17 High- (HNA) and low-nucleic acid (LNA) bacteria are two separated flow cytometry (FCM)
18 groups that are ubiquitous across aquatic systems. HNA cell density often correlates strongly
19 with heterotrophic production. However, the taxonomic composition of bacterial taxa within
20 HNA and LNA groups remains mostly unresolved. Here, we associated freshwater bacterial taxa
21 with HNA and LNA groups by integrating FCM and 16S rRNA gene sequencing using a
22 machine learning-based variable selection approach. There was a strong association between
23 bacterial heterotrophic production and HNA cell abundances ($R^2 = 0.65$), but not with more
24 abundant LNA cells, suggesting that the smaller pool of HNA bacteria may play a
25 disproportionately large role in the freshwater carbon flux. Variables selected by the models
26 were able to predict HNA and LNA cell abundances at all taxonomic levels, with highest
27 accuracy at the OTU level. There was high system specificity as the selected OTUs were mostly
28 unique to each lake ecosystem and some OTUs were selected for both groups or were rare. Our
29 approach allows for the association of OTUs with FCM functional groups and thus the
30 identification of putative indicators of heterotrophic activity in aquatic systems, an approach that
31 can be generalized to other ecosystems and functioning of interest.

32 **Introduction**

33 A key goal in the field of microbial ecology is to understand the relationship between microbial
34 diversity and ecosystem functioning. However, it is challenging to associate bacterial taxa to
35 specific ecosystem processes. Marker gene surveys have shown that natural bacterial
36 communities are extremely diverse, however, the presence of a taxon does not imply their
37 activity. Taxa present in these surveys may have low metabolic potential, be dormant, or have
38 recently died [1, 2]. Therefore, new methodologies which integrate different data types are
39 needed to associate bacterial taxa with ecosystem functions in order to ultimately model and
40 predict them [3].

41
42 One such advance is the use of flow cytometry (FCM), which has been used extensively to study
43 aquatic microbial communities [4–6]. This single-cell technology partitions individual microbial
44 cells into phenotypic groups based on their observable optical characteristics. Most commonly,
45 cells are stained with a nucleic acid stain (*e.g.* SYBR Green I) and upon analysis assigned to
46 either a low nucleic acid (LNA) or a high nucleic acid (HNA) group [7–10]. HNA cells differ
47 from LNA cells in both a considerable increase in fluorescence due to cellular nucleic acid
48 content and scatter intensity due to cell morphology. The HNA group is thought to correspond to
49 the ‘active’ fraction, whereas the LNA population has been considered as the ‘dormant’ or
50 ‘inactive’ group of a microbial community [4, 11–13]. This is based on positive linear
51 relationships between HNA abundance and (a) bacterial heterotrophic production (BP) [8, 12–
52 15], (b) bacterial activity measured using the dye 5-cyano-2,3-ditolyl tetrazolium chloride [16,
53 17], and (c) phytoplankton abundance [18]. Additionally, growth rates are higher for HNA than

54 LNA cells [11, 14, 19] and HNA cells accrue cell damage significantly faster than the LNA cells
55 under temperature [20] and chemical oxidant stress [21].

56

57 One main research question that still remains is whether HNA and LNA groups are composed of
58 unique taxa or if they are different physiological states of the same taxa. Bouvier et al. [9]
59 proposed four possible scenarios: (1) bacteria start their life cycle in the HNA group and move to
60 the LNA group upon death or inactivity; (2) cells in the HNA group originate from LNA cells
61 undergoing cell division; (3) HNA and LNA consist of different non-overlapping taxa; (4)
62 bacteria switch between groups from time to time in addition to having part of the community
63 that is unique to each fraction. The view that HNA cells are more active is in line with scenario 1
64 and 2. On the other hand, several studies have found distinct groups with little taxonomic overlap
65 and proposed scenario 3 [22, 23] or 3 and 4 [24]. In this case, HNA and LNA groups have been
66 associated with different life strategies in bacterioplankton communities, such as large cell size
67 (HNA) versus small cell size (LNA) [13, 23], genome size [15] and ploidy [22]. By combining
68 FCM with taxonomic identification of bacterial communities, one can associate individual taxa
69 with population dynamics and functioning.

70

71 In this study, we developed a novel approach to associate the dynamics of individual taxa with
72 those of the LNA and HNA groups in freshwater lakes by using a machine learning variable
73 selection strategy. We applied two variable selection methods, the Randomized Lasso [25] and
74 the Boruta algorithm [26] to associate individual taxa with HNA and LNA cell abundances. This
75 approach allowed us to associate specific taxa to FCM functional groups, and via the observed
76 HNA-productivity relationship, to functioning. In addition, this approach enabled us to test the

77 influence of rare taxa on these two groups as recent research has found that rare taxa may have a
78 strong impact on community structure and functioning [27, 28]. To validate the RL-based
79 association with the HNA and/or LNA group, we correlated taxon abundances with specific
80 regions in the FCM fingerprint without prior knowledge of the HNA/LNA group. Furthermore,
81 we tested for phylogenetic conservation of HNA and LNA functional groups and for the
82 association between the selected taxa and productivity. The combination of FCM and 16S rRNA
83 gene sequencing allows for the inference and assessment of the taxonomic structure of HNA and
84 LNA groups, therefore advancing our ability to link bacterial taxa to their functionality in nature.
85 This knowledge will help identify the taxa that drive carbon fluxes in freshwater ecosystems,
86 which are disproportionately large relative to the global freshwater surface area [29].

87 **Results**

88 In this study, we developed a machine learning variable selection strategy to integrate FCM and
89 16S rRNA gene sequencing with the aim of inferring the bacterial drivers of functional groups in
90 freshwater lake systems. We studied a set of oligo- to eutrophic small inland lakes, a short
91 residence time mesotrophic freshwater estuary lake (Muskegon Lake), and a large oligotrophic
92 Great Lake (Lake Michigan), all located in Michigan, USA. We showed that abundance variation
93 of these FCM functional groups is predicted by a small subset of all taxa that are present in the
94 environment. Selected taxa were mostly FCM groups and lake system specific, and across
95 systems, association with HNA or LNA was not phylogenetically conserved. The relationship
96 between selected taxa and productivity measurements was assessed for one of the lake systems
97 (Muskegon Lake), thereby showing that HNA cells (and their putative bacterial taxa) likely turn
98 over faster and disproportionately contribute to the freshwater carbon flux.

99

100 *Study lakes are dominated by LNA cells*

101 The inland lakes (6.3×10^6 cells/mL) and Muskegon Lake (6.0×10^6 cell/mL) had significantly
102 higher total cell abundances than Lake Michigan (1.7×10^6 cell/mL; $p = 2.7 \times 10^{-14}$). Across all
103 lakes, the mean proportion of HNA cell counts (HNAcc) to total cell counts was much lower
104 (29-33%) compared to the mean proportion of LNA cell counts (LNAcc; 67-71%). Through
105 ordinary least squares regression, there was a strong correlation between HNAcc and LNAcc
106 across all data ($R^2 = 0.45$, $P = 2 \times 10^{-24}$; **Figure 1A**), however, only Lake Michigan ($R^2 = 0.59$, P
107 $= 5 \times 10^{-11}$) and Muskegon Lake ($R^2 = 0.44$, $P = 2 \times 10^{-9}$) had significant correlations when the
108 three ecosystems were considered separately.

109

110 *HNA cell counts and heterotrophic bacterial production are strongly correlated*

111 For mesotrophic Muskegon Lake, there was a strong correlation between total bacterial
112 heterotrophic production and HNAcc ($R^2 = 0.65$, $p = 1e-05$; **Figure 1B**), no correlation between
113 BP and LNAcc ($R^2 = 0.005$, $p = 0.31$; **Figure 1C**), and a weak correlation between heterotrophic
114 production and total cell counts ($R^2 = 0.18$, $p = 0.03$; **Figure 1D**). There was a positive (HNA)
115 and negative (LNA) correlation between the fraction of HNA or LNA to total cells and
116 productivity, however, the relationship was weak and not significant ($R^2 = 0.14$, $p = 0.057$).

117

118 *Association of OTUs to functional groups by Randomized Lasso regression*

119 The relevance of specific OTUs for predicting freshwater FCM functional group abundance was
120 assessed using the Randomized Lasso (RL) approach, which assigns a score between 0
121 (unimportant) to 1 (highly important) to each taxon in function of the target variable: HNAcc or
122 LNAcc. This score can be interpreted as the probability that an OTU will be included in the

123 Lasso model to predict HNA or LNA cell abundances. Variations of HNAcc and LNAcc were
124 modelled in function of relative changes of OTUs. To address the negative correlation bias
125 intrinsic to compositional data, compositions were first transformed using a centered log-ratio
126 (CLR) transformation.

127

128 The RL score was used to implement a recursive variable elimination scheme. Specifically, we
129 iteratively removed the lowest-ranked OTUs based on the RL score (*i.e.* OTUs were ranked
130 according to the score from high to low) and the Lasso was fitted to the data to predict HNAcc
131 and LNAcc based on the corresponding subset of OTUs. The performance was expressed in
132 terms of the R_{CV}^2 , the R^2 between predicted and true values of HNAcc and LNAcc of samples
133 that were held-out using a leave-one-group-out cross-validation scheme, in which samples were
134 grouped according to year and location of measurement. If R_{CV}^2 equals 1, predictions were equal
135 to the true values, a value of 0 is equivalent to random guessing.

136

137 There was taxonomic dependency for both HNAcc and LNAcc across lake systems (**Figure 2**).
138 R_{CV}^2 increased when lower-ranked OTUs were removed (moving from right to left on **Figure 2**),
139 which was gradual for the inland lakes (**Figure 2A**) and Muskegon Lake (**Figure 2C**) but was
140 abrupt for Lake Michigan (**Figure 2B**). The number of taxa that resulted in the highest R_{CV}^2
141 contained less than a quarter of the total amount of taxa that were present (*see solid (HNA) and*
142 *dotted (LNA) lines in Figure 2*), being 10.2% HNA and 15.3% LNA for the inland lakes, 4.0%
143 HNA and 3.0% LNA for Lake Michigan, and 25.0% for both HNA and LNA in Muskegon Lake.
144 This behavior was consistent for each lake system and FCM population. The Lake Michigan
145 results differed the most from other lake systems, having the lowest R_{CV}^2 , a sharp increase in R_{CV}^2

146 instead of gradual, and a considerably lower minimal amount of OTUs (13 for HNacc, 10 for
147 LNacc). No relationship could be established between rankings of variable selection methods
148 and the relative abundance of individual OTUs (**Figure S1**). Multiple taxa with low average
149 abundance were included in the minimal set of predictive variables, whereas few highly
150 abundant OTUs were included. HNacc and LNacc could be predicted with equivalent
151 performance to relative HNA and LNA proportions, yet the increase between initial and optimal
152 performance was bigger (**Figure S2**). The final predictive performance was lower when
153 compositional data was not transformed using the CLR-transformation (**Figure S3**).

154

155 *Identification on different taxonomic levels: OTUs outperform all other taxonomic levels*

156 To assess whether HNA and LNA groups were taxonomically conserved, compositional data
157 was analyzed on all possible taxonomic levels for Muskegon Lake (**Figure 3**), using the same
158 strategy as outlined in previous paragraph. The resulting R_{CV}^2 values were considerably higher
159 than zero on all taxonomic levels, meaning that at all levels individual taxonomic changes can be
160 related to changes in HNacc and LNacc. Even though the OTU level resulted in the best
161 prediction of HNacc and LNacc (**Figure 3**), each individual OTU has a lower RL score
162 compared to other taxonomic levels, which on average became lower as the taxonomic level
163 decreased (**Figure S4**). The fraction of variables (taxa) that could be removed to reach the
164 maximum R_{CV}^2 decreased as the taxonomic level became less resolved.

165

166 *Validation of OTU selection results with the Boruta algorithm*

167 The OTU results were validated with an additional variable selection strategy, called the Boruta
168 algorithm. This approach allowed the further generalization of the findings presented above. In

169 addition, it connects with Random Forest results from other studies, which have been described
170 recently in microbiome studies of other systems (*see [30] and [31]*). The Boruta algorithm
171 selects relevant variables based on statistical hypothesis testing between the variable importance
172 of an original variable and the importance of the most important permuted variable (*see*
173 *materials and methods*), as retrieved from multiple Random Forest models. Selected variables
174 are ranked as ‘1’, tentative variables as ‘2’, and all other variables get lower ranks, depending on
175 the stage in which they were eliminated. The Boruta algorithm was applied for all three lake
176 systems at the OTU-level, selected OTUs are visualized in **Figure S5**. The fraction of selected
177 OTUs was always smaller than 1% across lake systems and functional groups (**Figure S6**). The
178 top scored OTU according to the RL was also selected according to the Boruta algorithm for
179 HNacc for all lake systems; for LNacc both methods only agreed for Lake Michigan (**Table 1**).
180 OTU060 (Proteobacteria;Sphingomonadales;alfIV_unclassified) was the only OTU selected in
181 function of LNacc across all lake systems, whereas no OTUs were selected across lake systems
182 for HNacc. As Random Forest regressions are the base method of the Boruta algorithm, we
183 compared the predictive power of Boruta selected OTUs to those of all OTUs using Random
184 Forest regression. For all lake systems and functional group performance increased when only
185 selected OTUs were included in the model (**Table S1**). Lasso predictions, in which OTUs were
186 selected according to the RL, were better as opposed to Random Forest predictions in which
187 OTUs were selected according to the Boruta algorithm (**Figure S7**). The fraction of selected
188 OTUs according to the Boruta algorithm was lower than the optimal amount of OTUs according
189 to the RL.
190

191 In this way, a number of findings could be generalized independent of a specific method: 1)
192 Selected OTUs were mostly lake systems specific, 2) a small fraction of OTUs was needed to
193 predict changes in community composition, 3) selected OTUs are often rare and do not show a
194 relationship with abundance and 4) top RL-ranked HNA OTUs were also selected according to
195 the Boruta algorithm, suggesting to inspect more closely the phylogeny of these taxa.

196

197 *HNA- and LNA-associated OTUs differed across lake systems*

198 Selected OTUs were mostly assigned to either the HNA or LNA groups and there was limited
199 correspondence across lake systems between the selected OTUs (**Figure 4**). In Muskegon Lake,
200 OTU173 (Bacteroidetes;Flavobacteriales;bacII-A) was selected as the major HNA-associated
201 taxon while OTU29 (Bacteroidetes;Cytophagales;bacIII-B) had the highest RL score for LNA
202 OTUs. In Lake Michigan, OTU25 (Bacteroidetes;Cytophagales;bacIII-A), was selected as the
203 major HNA-associated taxon while OTU168 (Alphaproteobacteria:Rhizobiales:alfVII) was
204 selected as a major LNA-associated taxon. For the inland lakes, OTU369
205 (Alphaproteobacterial;Rhodospirillales;alfVIII) was the major HNA-associated OTU while the
206 OTU555 (Deltaproteobacteria;Bdellovibrionaceae;OM27) was the major LNA-associated taxon.
207 Many more OTUs were selected in Muskegon Lake (197 OTUs; compared to 134 OTUs from
208 the Inland Lakes and 21 OTUs from Lake Michigan) and these OTUs were often associated
209 with both HNA and LNA groups.

210

211 RL scores were correlated for HNAcc and LNA within each lake system (Inland $r = 0.25$, $P <$
212 0.001 ; Michigan $r = 0.59$, $P < 0.001$, Muskegon $r = 0.59$, $P < 0.001$). Only OTUs that were
213 present in all three freshwater environments were considered to calculate correlations between

214 lake systems (190 in total, **Figure S8**). RL scores were lake ecosystem specific, with only a
215 significant similarity between the Inland lakes and Muskegon lake using the RL for HNacc ($r =$
216 0.21 , $P = 0.0042$). Note that the correlation within a lake system therefore differs from
217 previously reported values (as not all OTUs were considered), yet differences were small and
218 results were comparable. The Boruta algorithm selected mostly OTUs which were unique both
219 for the lake system and functional population (**Figure S5**).

220

221 *Selected HNA and LNA OTUs do not have a phylogenetic signal*

222 While many of the 258 OTUs selected by the RL were one of a few members of their phylum
223 (*e.g.* Firmicutes; Epsilonproteobacteria; OTU717 in Lentisphaerae; OTU267 in Omnitrophica;
224 etc), the Bacteroidetes (60 OTUs), Betaproteobacteria (36 OTUs), Alphaproteobacteria (22
225 OTUs), and Verrucomicrobia (21 OTUs) were a total of 54% of the selected OTUs (**Figure 5**).
226 Of these top four phyla, the majority of their membership were within the LNA group (41-52%
227 of selected OTUs), with the minority of OTUs within the HNA group (14-30% of selected
228 OTUs), and a quarter to a third of the OTUs were selected as members of both the LNA and
229 HNA groups (23-36% of selected OTUs).

230

231 To evaluate how much phylogenetic history explains whether a selected taxon was associated
232 with the HNA and/or LNA group(s), we calculated the phylogenetic signal, which is a measure
233 of the dependence among species' trait values on their phylogenetic history [32]. If the
234 phylogenetic signal is very strong, taxa belonging to similar phylogenetic groups (*e.g.* a Phylum)
235 will share the same trait (*i.e.* association with HNacc or LNacc). Alternatively, if the
236 phylogenetic signal is weak, taxa within a similar phylogenetic group will have different traits.

237 For the most part, Pagel's lambda was used [33] to test for phylogenetic signal where lambda
238 varies between 0 and 1. A lambda value of 1 indicates complete phylogenetic patterning whereas
239 a lambda of 0 indicates no phylogenetic patterning and leads to a tree collapsing into a single
240 polytomy. There was no phylogenetic signal with FCM functional group used as a discrete
241 character (*i.e.* HNA, LNA, or Both). As a continuous character using the RL scores for HNA
242 (**Figure S9**), there was also no phylogenetic signal ($\lambda = 0.16$; $P = 1$). There was a
243 significant LNA signal ($p = 0.003$), however, the lambda value was 0.66, suggesting weak
244 phylogenetic structuring in the LNA group. However, this significant result in the LNA was not
245 replicated with other measures of phylogenetic signal (Blomberg's K (HNA: $p = 0.63$; LNA: $p =$
246 0.54), and Moran's I (HNA: $p = 0.88$; LNA: $p = 0.12$)) indicating that there is likely no
247 phylogenetic signal in the taxa that drive the dynamics in either the HNA or the LNA group.

248

249 ***Flow cytometry fingerprints confirm associated taxa and reveal complex relationships between***
250 ***taxonomy and flow cytometric fingerprints***

251 To confirm the association of the final selected OTUs with the HNA and LNA groups, we
252 calculated the correlation between the density of individual regions (*i.e.* "bins") in the flow
253 cytometry data with the relative abundances of the OTUs. The Kendall rank correlation
254 coefficient between OTU abundances and counts in the flow cytometry fingerprint was
255 calculated for each of the top HNA OTUs selected by the RL within each of the three systems.
256 The correlation coefficient was visualized for each bin in the flow cytometry fingerprint (**Figure**
257 **6**). As these values denote correlations, they do not indicate actual presence. OTU25 correlated
258 with almost the entire HNA region, whereas OTU173 was limited to the lower part of the HNA
259 region. In contrast, OTU369 was positively correlated to both the LNA and HNA regions of the

260 cytometric fingerprint, highlighting results from **Figure 4** where OTU369 was selected in
261 function of both HNA and an LNA. The threshold that was used to define HNacc and LNacc
262 lies very close to the actual corresponding regions.

263

264 **Proteobacteria and rare taxa correlate with productivity measurements**

265 The Kendall rank correlation coefficient was calculated between CLR-transformed abundances
266 of individual OTUs and productivity measurements. OTU481 was significantly correlated after
267 correction for multiple hypothesis testing using the Benjamini-Hochberg procedure ($P < 0.001$,
268 $P_{\text{adj}} = 0.016$). This OTU had however a low RL score (0.022) and was not selected according
269 to the Boruta algorithm. Of the top 10 OTUs according to the RL, three still had significant P-
270 values (OTU614: $P = 0.0064$; OTU412, $P = 0.044$; OTU487, $P = 0.014$). Some OTUs that had a
271 high RL score also had a positive response to productivity measurements (**Figure S10**). At the
272 phylum level, only Proteobacteria were significantly correlated to productivity measurements
273 after Benjamini-Hochberg correction ($P < 0.001$, $P_{\text{adj}} = 0.010$).

274 **Discussion**

275 Our study introduces a novel computational workflow to investigate relationships between
276 microbial diversity and ecosystem functioning. Specifically, we aimed to study the ecology of
277 flow cytometric functional groups (i.e. HNA and LNA) by associating their dynamics with those
278 of bacterial taxa (i.e. OTUs). We simultaneously collected flow cytometry and 16S rRNA gene
279 sequencing data from three types of freshwater lake systems in the Great Lakes region, and
280 bacterial heterotrophic productivity from one lake ecosystem, and used a machine learning based
281 variable selection strategy, known as the Randomized Lasso, to associate one with another. Our
282 results showed that (1) there was a strong correlation between bacterial heterotrophic

283 productivity and HNA cell abundances, (2) HNA and LNA cell abundances were best predicted
284 by a small subset of OTUs that were unique to each lake type, (3) some OTUs were included in
285 the best model for both HNA and LNA abundance, (4) there was no phylogenetic conservation
286 of HNA and LNA group association and (5) freshwater FCM fingerprints display more complex
287 patterns related to OTUs and productivity than compared to the traditional dichotomy of HNA
288 and LNA. While HNA and LNA groups are universal across aquatic ecosystems, our data
289 suggest that some bacterial taxa contribute to both HNA and LNA groups and that the taxa
290 driving HNA and LNA abundance are unique to each lake system.

291
292 Although high-nucleic acid cell counts (HNAcc) and low-nucleic acid cell counts (LNAcc) were
293 correlated with each other, only the association between bacterial heterotrophic production (BP)
294 and HNAcc was strong and significant. This correlation between BP and HNA is higher than
295 previously reported values, though previous reports have focused on the proportion of HNA
296 rather than absolute cell abundances with the majority of data collected from marine systems.
297 For example, Bouvier et al. [9] found a correlation between the fraction of HNA cells and BP
298 within a large dataset of 640 samples across various freshwater to marine samples ($r = 0.49$),
299 whereas a study off the coast of the Antarctic Peninsula found a moderate correlation ($R^2 = 0.36$;
300 [15]). Another study in the Bay of Biscay also found this association ($R^2 = 0.16$; [13]), however,
301 the authors attributed this difference to be related to cell size and not due to the activity of HNA.
302 Notably, these studies were predominantly testing the association of marine HNA and the reason
303 for the stronger correlation in our study may be due to the nature of the freshwater samples. As
304 such, future studies in freshwater environments should test this hypothesis, which is especially
305 important for understanding the broader influence that HNA bacteria may have in the context of

306 the disproportionately large role that freshwater systems play as hotspots in the global carbon
307 cycle [29]. Finally, as our correlations with proportional HNA abundance also indicated less
308 strong correlations than with absolute HNAcc, we suggest absolute HNAcc should be used to
309 best predict heterotrophic bacterial production with FCM data.

310

311 The use of machine learning methods, such as the Lasso and Random Forest, are becoming more
312 common in microbiome literature as these approaches are able to deal with multi-dimensional
313 data and test the predictive power of a combined set of variables ([34–36]. Although the Lasso
314 already uses an intrinsic variable selection strategy, it has been noted that the Lasso method is
315 not suited for compositional data because the regression coefficients have an unclear
316 interpretation, and single variables may be selected when correlated to other variables [37].

317 When performing variable selection with Random Forests, traditional variable importance
318 measures such as the mean decrease in accuracy can be biased towards correlated variables [38].

319 Our approach included algorithms which extended on these traditional machine learning
320 algorithms, i.e. the Randomized Lasso or Boruta algorithm [25, 26]. These methods make use of
321 resampling and randomization which allow to either assign a probability of selection (RL) or
322 statistically decide which OTU to select (Boruta). Both the RL and Boruta algorithm have been
323 applied to microbiome studies before. Examples for RL include the selection of genera in the gut
324 microbiome relation to BMI [34] or the selection of OTUs from the oral microbiome in function
325 of salivary pH and lysozyme activity [39]. The Boruta algorithm has been applied to select
326 relevant genera, for example in the gut microbiome in relation to multiple sclerosis [31] or in
327 function of different diets during pregnancy of primates [30]. Moreover, the Boruta algorithm
328 has been recently proposed as one of the top-performing variable selection methods that make

329 use of Random Forests [40]. The ability of our approach to identify unique sets of OTUs
330 predictive of HNacc and LNacc despite the correlation between HNacc and LNacc (**Figure**
331 **1A**) illustrates the power of the machine learning based-variable selection methods. However,
332 there is still room for improvement when attempting to integrate these different types of data. For
333 example, 16S rRNA gene sequencing still faces the hurdles of DNA extraction [41] and 16S
334 copy number bias [42]. Moreover, detection limits are different for FCM (expressed in the
335 number of cells) and 16S rRNA gene sequencing (expressed in the number of gene counts or
336 relative abundance), which create data that may be different in resolution. Future work may
337 focus on developing ways around these shortcomings to further improve the integration of FCM
338 with 16S rRNA gene sequencing.

339

340 In our study, only a minority of OTUs was needed to predict specific flow cytometric group
341 abundances. While each OTU individually had low predictive power, the selected group of
342 OTUs was generally a strong predictor of HNacc and LNacc. In addition, the selected OTUs
343 were often rare and thus no relationship could be established between the RL score and the
344 abundance of an OTU (**Figure S3**). These results are in line with recent findings of Herren &
345 McMahon [28], who reported that a minority of low abundance taxa explained temporal
346 compositional changes of microbial communities. The selection of different sets of HNA and
347 LNA OTUs across the three freshwater systems indicates that different taxa underlie the
348 universally observed HNA and LNA functional groups across aquatic systems. This is in line
349 with strong species sorting in lake systems [43, 44], shaping community composition through
350 diverging environmental conditions between the lake systems presented here [45]. This high
351 system specificity also explains the low RL scores for individual OTUs, as the spatial dynamics

352 of an OTU diverged strongly across systems. (For example, an OTU that has an RL score of 0.5
353 implies that on average it will only be chosen one out of two times in a Lasso model).

354

355 Based on the high correlation of BP with HNAcc and low correlation with BP and LNAcc, the
356 high proportion of LNA cells across all lake systems might indicate that the majority of cells in
357 the bacterial community are dormant or have very low activity. This agrees with previous
358 research showing that up to 40% [46] or even 64-95% [47] of cells in freshwater systems to be
359 inactive or dormant. In fact, up to 60-80% of the OTUs in freshwater lakes have been reported to
360 be dormant [48]. Based on variable environmental conditions sampled across our dataset, some
361 of the taxa that are predominantly dormant in one sample may contribute to activity in another
362 sample. If this differing contribution to activity also covaries with a taxon's abundance, these
363 taxa may be considered to be 'conditionally rare taxa' [49] and previously 1-2% of freshwater
364 lake OTUs have been reported to be conditionally rare [27]. It has also been shown that marine
365 heterotrophic bacteria can survive for at least 8 months (maximum tested length) in a starved
366 state [50]. These factors may explain why some OTUs were included in both the HNAcc and
367 LNAcc models and is in line with scenario 1 from Bouvier et al [9] (*i.e.* the transitioning of cells
368 from active growth to death or inactivity). Alternatively, the same OTU may occur in both HNA
369 and LNA groups due to phenotypic plasticity. Phenotypic plasticity has been shown for bacterial
370 morphology and size, for example during predation and carbon starvation [51]. The fact that
371 HNA and LNA groups have been suggested to correspond to cells of differing size, with HNA
372 harboring larger cell sizes [10, 23], is in line with this hypothesis. Finally, the OTU level
373 grouping of bacterial taxa can disguise genomic and phenotypic heterogeneity [52–55], which
374 may be an explanation for inconsistent associations between OTUs and FCM functional groups.

375

376 While all taxonomic levels resulted in a model with predictive power, the best model was at the
377 most resolved taxonomy (*i.e.* OTU) indicating that it is unlikely that OTUs within the HNA and
378 LNA groups are phylogenetically conserved. Indeed, when analyzing the data at an OTU level,
379 very little phylogenetic conservation was found between selected OTUs for HNA and LNA
380 groups. This is in contrast to a recent study that found a clear signal at the phylum level [23].
381 Proctor et al. [23] showed separate bacterial clusters between HNA and LNA groups across
382 different aquatic systems. However, this was not the case for lake water samples. It is notable
383 that Proctor et al. [23] separated HNA and LNA cells based on cell size (where HNA cells were
384 >0.4 μm and LNA cells were 0.2-0.4 μm , based on 50-90% removal of HNA cells after filtering),
385 while our study separated these FCM functional groups on the basis of fluorescence intensity
386 alone. Moreover, our study assessed associations between OTUs and population dynamics, while
387 Proctor et al. [23] assessed actual presence.

388

389 The Boruta algorithm and RL scores agreed on the top-ranked HNA OTU for all lake systems,
390 which motivates further investigation of the ecology of these OTUs. While little information on
391 the identities of HNA and LNA freshwater lake bacterial taxa exists, several studies identified
392 Bacteroidetes among the most prominent HNA taxa and is in line with our findings. Vila-Costa
393 et al. [24] found that the HNA group was dominated by Bacteroidetes in summer samples from
394 the Mediterranean Sea, Read et al. [17] showed that HNA abundances correlated with
395 Bacteroidetes, and Schattenhofer et al. [22] reported that the Bacteroidetes accounted for the
396 majority of HNA cells in the North Atlantic Ocean. In Muskegon Lake, OTU173 was the
397 dominant HNA taxon and is a member of the Order *Flavobacteriales* (bacII-A). The bacII group

398 is a very abundant freshwater bacterial group and has been associated with senescence and
399 decline of an intense algal bloom [56]. BacII-A has also made up ~10% of the total microbial
400 community during cyanobacterial blooms, reaching its maximum density immediately following
401 the bloom [57]. In Lake Michigan, OTU25, a member of the Bacteroidetes Order *Cytophagales*
402 known as bacIII-A, was the top HNA OTU. However, much less is known about this specific
403 group of Bacteroidetes. Though, the bacII-A/bacIII-A group has been strongly associated with
404 more heterotrophically productive headwater sites (compared to higher order streams) from the
405 River Thames, showing a negative correlation in rivers with dendritic distance from the
406 headwaters, indicating that these taxa may contribute more to productivity [17]. In the inland
407 lakes, OTU369 was the major HNA taxon and is associated with the Alphaproteobacteria Order
408 Rhodospirillales (alfVIII), which to our knowledge is a group with very little information
409 available in the literature. In contrast to our findings of Bacteroidetes and Alphaproteobacterial
410 HNA selected OTUs, Tada & Suzuki [58] found that the major HNA taxon from an oceanic algal
411 culture was from the Betaproteobacteria whereas LNA OTUs were within the Actinobacteria
412 phylum.

413

414 **Conclusions**

415 Our results indicate that there are taxonomic differences between HNA and LNA groups in
416 freshwater lake systems, though these are lake system specific. This result may be due to taxa
417 switching between these groups, potentially due to genomic or phenotypic plasticity. The
418 difference between selected taxa is larger between lake systems as opposed to differences
419 between HNA and LNA groups, which were not conserved phylogenetically. Thus, our results
420 correspond most with research presented by Vila-Costa et al. [24], in which a taxonomic division

421 was found between HNA and LNA groups, yet this was not rigid and followed seasonal trends.
422 Overall, our results motivate scenario 4 proposed by Bouvier et al. [9], where HNA and LNA
423 exhibit a different taxonomy, but this taxonomy changes over time and space and may overlap.
424 With this study, we show that different types of microbial ecological data can be integrated with
425 machine learning to learn about the composition and functioning of bacterial populations in
426 aquatic systems. Future studies on HNA and LNA bacterial groups should use genome-resolved
427 metagenomics, metatranscriptomics, or single-cell genomics to decipher whether the traits that
428 underpin the association of a taxon with a FCM group are related to genomic or phenotypic
429 plasticity.

430

431 **Materials and Methods**

432 *Data collection and DNA extraction, sequencing and processing*

433 In this study, we used a total of 173 samples collected from three types of lake systems described
434 previously [45], including: (1) 49 samples from Lake Michigan (2013 & 2015), (2) 62 samples
435 from Muskegon Lake (2013-2015; one of Lake Michigan's estuaries), and (3) 62 samples from
436 twelve inland lakes in Southeastern Michigan (2014-2015). For more details on sampling, please
437 see **Figure 1** and the *Field Sampling, DNA extraction, and DNA sequencing and processing*
438 sections within Chiang et al. [45]. In all cases, water for microbial biomass samples were
439 collected and poured through a 210 μm and 20 μm bleach sterilized nitex mesh and sequential in-
440 line filtration was performed using 47 mm polycarbonate in-line filter holders (Pall Corporation,
441 Ann Arbor, MI, USA) and an E/S portable peristaltic pump with an easy-load L/S pump head
442 (Masterflex®, Cole Parmer Instrument Company, Vernon Hills, IL, USA) to filter first through a
443 3 μm isopore polycarbonate (TSTP, 47 mm diameter, Millipore, Billerica, MA, USA) and

444 second through a 0.22 μm Express Plus polyethersulfone membrane filters (47 mm diameter,
445 Millipore, MA, USA). The current study only utilized the 3 - 0.22 μm fraction for analyses.
446
447 DNA extractions and sequencing were performed as described in Chiang et al. [45]. Fastq files
448 were submitted to NCBI sequence read archive under BioProject accession number
449 PRJNA412984 and PRJNA414423. We analyzed the sequence data using MOTHUR V.1.38.0
450 (seed = 777; [59] based on the MiSeq standard operating procedure and put together at the
451 following link: https://github.com/rprops/Mothur_oligo_batch. A combination of the Silva
452 Database (release 123; [60]) and the freshwater TaxAss 16S rRNA database and pipeline [61]
453 was used for classification of operational taxonomic units (OTUs).

454
455 For the taxonomic analysis, each of the three lake datasets were analyzed separately and treated
456 with an OTU abundance threshold cutoff of at least 5 sequences in 10% of the samples in the
457 dataset (similar strategy to [62]). For comparison of taxonomic abundances across samples, each
458 of the three datasets were then rarefied to an even sequencing depth, which was 4,491 sequences
459 for Muskegon Lake samples, 5,724 sequences for the Lake Michigan samples, and 9,037
460 sequences for the inland lake samples. Next, the relative abundance at the OTU level was
461 calculated using the *transform_sample_counts()* function in the phyloseq R package [63] by
462 taking the count value and dividing it by the sequencing depth of the sample. For all other
463 taxonomic levels, the taxonomy was merged at certain taxonomic ranks using the *tax_glom()*
464 function in phyloseq [63] and the relative abundance was re-calculated.

465

466 ***Heterotrophic bacterial production measurements***

467 Muskegon Lake samples from 2014 and 2015 were processed for heterotrophic bacterial
468 production using the [³H] leucine incorporation into bacterial protein in the dark method [64, 65].
469 At the end of the incubation with [³H]-leucine, cold trichloroacetic acid-extracted samples were
470 filtered onto 0.2 µm filters that represented the leucine incorporation by the bacterial community.
471 Measured leucine incorporation during the incubation was converted to bacterial carbon
472 production rate using a standard theoretical conversion factor of 2.3 kg C per mole of leucine
473 [65].

474

475 *Flow cytometry, measuring HNA and LNA*

476 In the field, a total of 1 mL of 20 µm filtered lake water were fixed with 5 µL of glutaraldehyde
477 (20% vol/vol stock), incubated for 10 minutes on the bench (covered with aluminum foil to
478 protect from light degradation), and then flash frozen in liquid nitrogen to later be stored in -
479 80°C freezer until later processing with a flow cytometer. Flow cytometry procedures followed
480 the protocol laid out in Props et al. [66], which also uses the samples presented in the current
481 study. Samples were stained with SYBR Green I and measured in triplicate. The lowest number
482 of cells collected after denoising was 2342. HNA and LNA groups were selected using the fixed
483 gates introduced in Prest et al. [67] and plotted in **Figure S11**. Cell counts were determined per
484 HNA and LNA group and averaged over the three replicates (giving rise to HNAcc and LNAcc).
485

486 **Data analysis**

487 Processed data and analysis code for the following analyses can be found on the GitHub page for
488 this project at https://deneflab.github.io/HNA_LNA_productivity/.

489

490 ***HNA-LNA and HNA-Productivity Statistics and Regressions***

491 We tested the difference in absolute number of cells within HNA and LNA functional groups
492 across running analysis of variance with a post-hoc Tukey HSD test (*aov()* and *TukeyHSD()*;
493 *stats* R package; [68]). In addition, we tested the association of HNA and LNA to each other and
494 with productivity by running ordinary least squares regression with the *lm()* (*stats* R package;
495 [68]).

496

497 ***Ranking correlation***

498 Ranking correlation between variables was calculated using the Kendall rank correlation
499 coefficient, using the *kendalltau()* function in Scipy (v1.0.0) or *cor()* in R (v3.2). The ‘tau-b’
500 implementation was used, which is able to deal with ties. Values range from -1 (strong
501 disagreement) to 1 (strong agreement). The same statistic was used to assess the similarity
502 between rankings of variable selection methods.

503

504 ***Centered-log ratio transform***

505 First, following guidelines from Paliy & Shanker, Gloor et al. and Quinn et al.[69–71], relative
506 abundances of OTUs were transformed using a centered log-ratio (CLR) transformation before
507 variable selection was applied. This means that the relative abundance x_i of a taxa was
508 transformed according to the geometric mean of that sample, in which there are p taxa present:

509
$$x'_i = \log(x_i / (\prod_{j=1}^p x_j)^{1/p})$$

510 Zero values were replaced by $\delta = 1/p^2$. This was done using the scikit-bio package
511 (www.scikit-bio.org, v0.4.1).

512

513 *Lasso & stability selection*

514 Scores were assigned to taxa based on an extension of the Lasso estimator, which is called
515 *stability selection* [25]. In the case of n samples, the Lasso estimator fits the following regression
516 model:

$$517 \hat{\beta}^\lambda = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|y - X\beta\|_2^2 + \lambda \sum_{j=1}^p |\beta_j| ,$$

518 in which X denotes the abundance table, y the target to predict, which either is HNA cell
519 abundances (HNAcc) or LNA cell abundances (LNAcc), and λ is a regularization parameter
520 which controls the complexity of the model and prevents overfitting. The Lasso performs an
521 intrinsic form of variable selection, as the weights of certain variables will be put to zero.

522

523 Stability selection, when applied to the Lasso, is in essence an extension of the Lasso regression.
524 It implements two types of randomizations to assign a score to the variables, and is therefore also
525 called the *Randomized Lasso* (RL). The resulting RL score can be seen as the probability that a
526 certain variable will be included in a Lasso regression model (*i.e.*, its weight will be non-zero
527 when fitted). When performing stability selection, the Lasso is fitted to B different subsamples of
528 the data of fraction $n/2$, denoted as X' and corresponding y' . A second randomization is added by
529 introducing a weakness parameter α . In each model, the penalty λ changes to a randomly chosen
530 value in the set $[\lambda, \lambda/\alpha]$, which means that a higher penalty will be assigned to a random subset
531 of the total amount of variables. The Randomized Lasso therefore becomes:

$$532 \hat{\beta}^\lambda = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|y' - X'\beta\|_2^2 + \lambda \sum_{j=1}^p \frac{|\beta_j|}{w_j} ,$$

533 where w_j is a random variable which is either α or 1. Next, the Randomized Lasso score (RL
534 score) is determined by counting the number of times the weight of a variable was non-zero for
535 each of the B models and divided by B . Meinshausen and Bühlmann show that, under stringent
536 conditions, the number of falsely selected variables is controlled for the Randomized Lasso when
537 the RL score is higher than 0.5. If λ is varied, one can determine the stability path, which is the
538 relationship between π and λ for every variable. For our implementation, $B = 500$, $\alpha = 0.5$ and
539 the highest score was selected in the stability path for which λ ranged from 10^{-3} until 10^3 ,
540 logarithmically divided in 100 intervals. The *RandomizedLasso()* function from the scikit-learn
541 machine learning library was used [72], v0.19.1).

542

543 ***Random Forests & Boruta***

544 The Boruta algorithm is a *wrapper* algorithm that makes use of Random Forests as a base
545 classification or regression method in order to select all relevant variables in function of a
546 response variable [26]. Similar to stability selection, the method uses an additional form of
547 randomness in order to perform variable selection. Random Forests are fitted to the data multiple
548 times. To remove the correlation to the response variable, each variable gets per iteration a so-
549 called *shadow variable*, which is a permuted copy of the original variable. Next, the Random
550 Forest algorithm is run with the extended set of variables, after which variable importances are
551 calculated for both original and shadow variables. The shadow variable that has the highest
552 importance score is used as reference, and every variable with significantly lower importance, as
553 determined by a Bonferroni corrected t-test, is removed. Likewise, variables containing an
554 importance score that is significantly higher are included in the final list of selected variables.
555 This procedure can be repeated until all original variables are either discarded or included in the

556 final set; variables that remain get the label ‘tentative’ (i.e., after all repetitions it is still not
557 possible to either select or discard a certain variable). We used the `boruta_py` package to
558 implement the Boruta algorithm (https://github.com/scikit-learn-contrib/boruta_py). Random
559 Forests were implemented using `RandomForestRegressor()` function from scikit-learn [72],
560 v0.19.1). Random Forests were run with 200 trees, the number of variables considered at every
561 split of a decision tree was $p/3$ and the minimal number of samples per leaf was set to five. The
562 latter were based on default values for Random Forests in a regression setting [73]. The Boruta
563 algorithm was run for 300 iterations, variables were selected or discarded at $P < .05$ after
564 performing Bonferroni correction.

565

566 ***Recursive variable elimination***

567 Scores of the Randomized Lasso were evaluated using a recursive variable elimination strategy
568 [74]. Variables were ranked according to the RL score. Next, the lowest-ranked variables were
569 eliminated from the dataset, after which the Lasso was applied to predict HNAcc and LNAcc
570 respectively. This process was repeated until only the highest-scored taxa remained. In this way,
571 performance of the Randomized Lasso was assessed from a minimal-optimal evaluation
572 perspective [75]. In other words, the lowest amount of variables that resulted in the highest
573 predictive performance was determined.

574

575 ***Performance evaluation***

576 In order to account for the spatiotemporal structure of the data, a blocked cross-validation
577 scheme was implemented [76]. Samples were grouped according the site and year that they were
578 collected. This results in 5, 10 and 16 distinctive groups for the Michigan, Muskegon and Inland

579 lake systems respectively. Predictive models were optimized in function of the R^2 between
580 predicted and true values of held-out groups using a leave-one-group-out cross-validation
581 scheme with the *LeaveOneGroupOut()* function. This results in a cross-validated R_{CV}^2 value. For
582 the Lasso, λ was determined using the *lassoCV()* function, with setting $\text{eps}=10^{-4}$ and
583 $\text{n_alphas}=400$. The Random Forest object was optimized using a grid search where max_features
584 was chosen in the interval $[1, \sqrt{p}, 2\sqrt{p}, \dots, p]$ (all variables) or $[1, \dots, p]$ (Boruta selected variables)
585 and min_samples_leaf in the interval $[1, \dots, 5]$, using the *GridSearchCV()* function. The number
586 of decision trees (n_trees) was set to 200. All functions are part of scikit-learn ([72]; v0.19.1)

587

588 *Stability of the Randomized Lasso*

589 Similarity of RL scores between lake systems and functional groups was quantified using the
590 Pearson correlation. This was done using the *pearsonr()* function in Scipy (v1.0.0).

591

592 *Patterns of HNA and LNA OTUs across ecosystems and phylogeny*

593 To visualize patterns of selected HNA and LNA OTUs across the three ecosystems, a heatmap
594 was created with the RL scores of each OTU from the Randomized Lasso regression that were
595 higher than specified threshold values. The heatmap was created with the *heatmap.2()* function
596 (*gplots* R package) using the euclidean distances of the RL scores and a complete linkage
597 hierarchical clustering algorithm (**Figure 4**).

598

599 *Correlations between taxa and productivity measurements*

600 Kendall tau ranking correlations between productivity measurements and individual abundances
601 were calculated on the phylum and OTU level using the *kendalltau()* function from Scipy

602 (v1.0.0). P-values were corrected using Benjamini-Hochberg correction, reported as P_{adj}. This
603 was done using the *multitest()* function from the Python module Statsmodels ([77]; v0.5.0).

604

605 ***Phylogenetic tree construction and signal calculation***

606 We calculated the best performing maximum likelihood tree using the GTR-CAT model (-gtr -
607 fastest) model of nucleotide substitution with fasttree (version 2.1.9 No SSE3; [78]).

608 Phylogenetic signal with both discrete (*i.e.* HNA, LNA, or both) and continuous traits (*i.e.* the
609 RL score) using the newick tree from FastTree was then used to model phylogenetic signal using
610 Pagel's lambda (discrete trait: fitDiscrete() from the geiger R package [79]; continuous trait:
611 phylsig() from the phytools R [80]), Blomberg's K (phylsig() function from the phytools R
612 package [80]), and Moran's I (abouheif.moran() function from the adephylo R package [81]).

613 **Acknowledgements**

614 PR was supported by Ghent University (BOFSTA2015000501) and MLS was supported by the
615 National Science Foundation Graduate Research Fellowship Program (Grant No. DGE
616 1256260). Part of the computational resources (Stevin Supercomputer Infrastructure) and
617 services used in this work were provided by the VSC (Flemish Supercomputer Center), funded
618 by Ghent University, the Hercules Foundation and the Flemish Government department EWI.
619 Flow cytometry analysis was supported through a Geconcerteerde Onderzoeksactie (GOA) from
620 Ghent University (BOF15/GOA/006).

621 **References**

- 622 1. Lennon JT, Jones SE. Microbial seed banks: the ecological and evolutionary implications
623 of dormancy. *Nat Rev Microbiol* 2011; **9**: 119–30.
- 624 2. Carini P, Marsden PJ, Leff JW, Morgan EE, Strickland MS, Fierer N. Relic DNA is
625 abundant in soil and obscures estimates of soil microbial diversity. *Nat Microbiol* 2016; **2**:
626 16242.
- 627 3. Widder S, Allen RJ, Pfeiffer T, Curtis TP, Wiuf C, Sloan WT, et al. Challenges in
628 microbial ecology: Building predictive understanding of community function and
629 dynamics. *ISME J* 2016; **10**: 2557–2568.
- 630 4. Gasol JM, Del Giorgio PA. Using flow cytometry for counting natural planktonic bacteria
631 and understanding the structure of planktonic bacterial communities. *Sci Mar* 2000; **64**:
632 197–224.
- 633 5. Vives-Rego J, Lebaron P, Caron Nebe-von. Current and future applications of flow
634 cytometry in aquatic microbiology. *FEMS Microbiol Rev* 2000; **24**: 429–448.
- 635 6. Wang Y, Hammes F, De Roy K, Verstraete W, Boon N. Past, present and future
636 applications of flow cytometry in aquatic microbiology. *Trends Biotechnol* 2010; **28**: 416–
637 424.
- 638 7. Gasol JM, Zweifel UL, Peters F, Jed A, Zweifel ULI, Fuhrman JEDA. Significance of
639 Size and Nucleic Acid Content Heterogeneity as Measured by Flow Cytometry in Natural
640 Planktonic Bacteria Significance of Size and Nucleic Acid Content Heterogeneity as
641 Measured by Flow Cytometry in Natural Planktonic Bacteria. *Appl Environ Microbiol*
642 1999; **65**: 4475–4483.
- 643 8. Lebaron P, Servais P, Agogué H, Courties C, Joux F. Does the High Nucleic Acid Content
644 of Individual Bacterial Cells Allow Us to Discriminate between Active Cells and Inactive
645 Cells in Aquatic Systems? *Appl Environ Microbiol* 2001; **67**: 1775–1782.
- 646 9. Bouvier T, Del Giorgio PA, Gasol JM. A comparative study of the cytometric
647 characteristics of High and Low nucleic-acid bacterioplankton cells from different aquatic
648 ecosystems. *Environ Microbiol* 2007; **9**: 2050–2066.
- 649 10. Wang Y, Hammes F, Boon N, Chami M, Egli T. Isolation and characterization of low
650 nucleic acid (LNA)-content bacteria. *ISME J* 2009; **3**: 889–902.
- 651 11. Lebaron P, Servais P, Baudoux a.-C, Bourrain M, Courties C, Parthuisot N. Variations of
652 bacterial-activity with cell size and nucleic acid content assessed by flow cytometry.
653 *Aquat Microb Ecol* 2002; **28**: 131–140.
- 654 12. Servais P, Casamayor EO, Courties C, Catala P, Parthuisot N, Lebaron P. Activity and
655 diversity of bacterial cells with high and low nucleic acid content. *Aquat Microb Ecol*
656 2003; **33**: 41–51.

- 657 13. Morán X, Bode A, Suárez L, Nogueira E. Assessing the relevance of nucleic acid content
658 as an indicator of marine bacterial activity. *Aquat Microb Ecol* 2007; **46**: 141–152.
- 659 14. Servais P, Courties C, Lebaron P, Troussellier M. Coupling bacterial activity
660 measurements with cell sorting by flow cytometry. *Microb Ecol* 1999; **38**: 180–189.
- 661 15. Bowman JS, Amaral-zettler LA, Rich JJ, Luria CM, Ducklow HW. Bacterial community
662 segmentation facilitates the prediction of ecosystem function along the coast of the
663 western Antarctic Peninsula. *ISME J* 2017; **11**: 1460–1471.
- 664 16. Morán XAG, Ducklow HW, Erickson M. Single-cell physiological structure and growth
665 rates of heterotrophic bacteria in a temperate estuary (Waquoit Bay, Massachusetts).
666 *Limnol Oceanogr* 2011; **56**: 37–48.
- 667 17. Read DS, Gweon HS, Bowes MJ, Newbold LK, Field D, Bailey MJ, et al. Catchment-
668 scale biogeography of riverine bacterioplankton. *ISME J* 2015; **9**: 516–526.
- 669 18. Sherr EB, Sherr BF, Longnecker K. Distribution of bacterial abundance and cell-specific
670 nucleic acid content in the Northeast Pacific Ocean. *Deep Res Part I Oceanogr Res Pap*
671 2006; **53**: 713–725.
- 672 19. Jochem FJ, Lavrentyev PJ, First MR. Growth and grazing rates of bacteria groups with
673 different apparent DNA content in the Gulf of Mexico. *Mar Biol* 2004; **145**: 1213–1225.
- 674 20. Arnoldini M, Heck T, Blanco-Fernández A, Hammes F. Monitoring of Dynamic
675 Microbiological Processes Using Real-Time Flow Cytometry. *PLoS One* 2013; **8**: e80117.
- 676 21. Ramseier MK, von Gunten U, Freihofer P, Hammes F. Kinetics of membrane damage to
677 high (HNA) and low (LNA) nucleic acid bacterial clusters in drinking water by ozone,
678 chlorine, chlorine dioxide, monochloramine, ferrate(VI), and permanganate. *Water Res*
679 2011; **45**: 1490–1500.
- 680 22. Schattenhofer M, Wulf J, Kostadinov I, Glöckner FO, Zubkov M V., Fuchs BM.
681 Phylogenetic characterisation of picoplanktonic populations with high and low nucleic
682 acid content in the North Atlantic Ocean. *Syst Appl Microbiol* 2011; **34**: 470–475.
- 683 23. Proctor CR, Besmer MD, Langenegger T, Beck K, Walser J-C, Ackermann M, et al.
684 Phylogenetic clustering of small low nucleic acid-content bacteria across diverse
685 freshwater ecosystems. *ISME J* 2018.
- 686 24. Vila-Costa M, Gasol JM, Sharma S, Moran MA. Community analysis of high- and low-
687 nucleic acid-containing bacteria in NW Mediterranean coastal waters using 16S rDNA
688 pyrosequencing. *Environ Microbiol* 2012; **14**: 1390–1402.
- 689 25. Meinshausen N, Bühlmann P. Stability selection. *J R Stat Soc Ser B Stat Methodol* 2010.
- 690 26. Kursu MB, Rudnicki WR. Feature Selection with the Boruta Package. *J Stat Softw* 2010;
691 **36**: 1–13.

- 692 27. Shade A, Jones SE, Caporaso JG, Handelsman J, Knight R, Fierer N, et al. Conditionally
693 rare taxa disproportionately contribute to temporal changes in microbial diversity. *MBio*
694 2014; **5**: e01371-14.
- 695 28. Herren CM, McMahon KD. Keystone taxa predict compositional change in microbial
696 communities. *Environ Microbiol* 2018; 1–34.
- 697 29. Biddanda BA. Global Significance of the Changing Freshwater Carbon Cycle Emerging
698 Role of Freshwater in the Global Theater. *Eos (Washington DC)* 2017; **98**: 1–5.
- 699 30. Ma J, Prince AL, Bader D, Hu M, Ganu R, Baquero K, et al. High-fat maternal diet during
700 pregnancy persistently alters the offspring microbiome in a primate model. *Nat Commun*
701 2014; **5**: 1–11.
- 702 31. Chen J, Chia N, Kalari KR, Yao JZ, Novotna M, Soldan MMP, et al. Multiple sclerosis
703 patients have a distinct gut microbiota compared to healthy controls. *Sci Rep* 2016; **6**: 1–
704 10.
- 705 32. Revell LJ, Harmon LJ, Collar DC. Phylogenetic signal, evolutionary process, and rate.
706 *Syst Biol* 2008; **57**: 591–601.
- 707 33. Pagel M. Inferring the historical patterns of biological evolution. *Nature* 1999; **401**: 877–
708 884.
- 709 34. Lin W, Shi P, Feng R, Li H. Variable selection in regression with compositional
710 covariates. *Biometrika* 2014; **101**: 785–797.
- 711 35. Baxter NT, Zackular JP, Chen GY, Schloss PD. Structure of the gut microbiome
712 following colonization with human feces determines colonic tumor burden. *Microbiome*
713 2014; **2**: 1–11.
- 714 36. Schubert AM, Rogers M a M, Ring C, Mogle J, Petrosino JP, Young VB, et al.
715 Microbiome Data Distinguish Patients with *Clostridium difficile* Infection and Non- *C* .
716 *difficile* -Associated Diarrhea from Healthy. *MBio* 2014; **5**: 1–9.
- 717 37. Li H. Microbiome, Metagenomics, and High-Dimensional Compositional Data Analysis.
718 *Annu Rev Stat Its Appl* 2015; **2**: 73–94.
- 719 38. Strobl C, Boulesteix A-L, Kneib T, Augustin T, Zeileis A. Conditional variable
720 importance for random forests. *BMC Bioinformatics* 2008; **9**: 307.
- 721 39. Zaura E, Brandt BW, Prodan A, Teixeira De Mattos MJ, Imangaliyev S, Kool J, et al. On
722 the ecosystemic network of saliva in healthy young adults. *ISME J* 2017; **11**: 1218–1231.
- 723 40. Degenhardt F, Seifert S, Szymczak S. Evaluation of variable selection methods for
724 random forests and omics data sets. *Brief Bioinform* 2017; 1–12.
- 725 41. McCarthy A, Chiang E, Schmidt ML, Denev VJ. RNA Preservation Agents and Nucleic

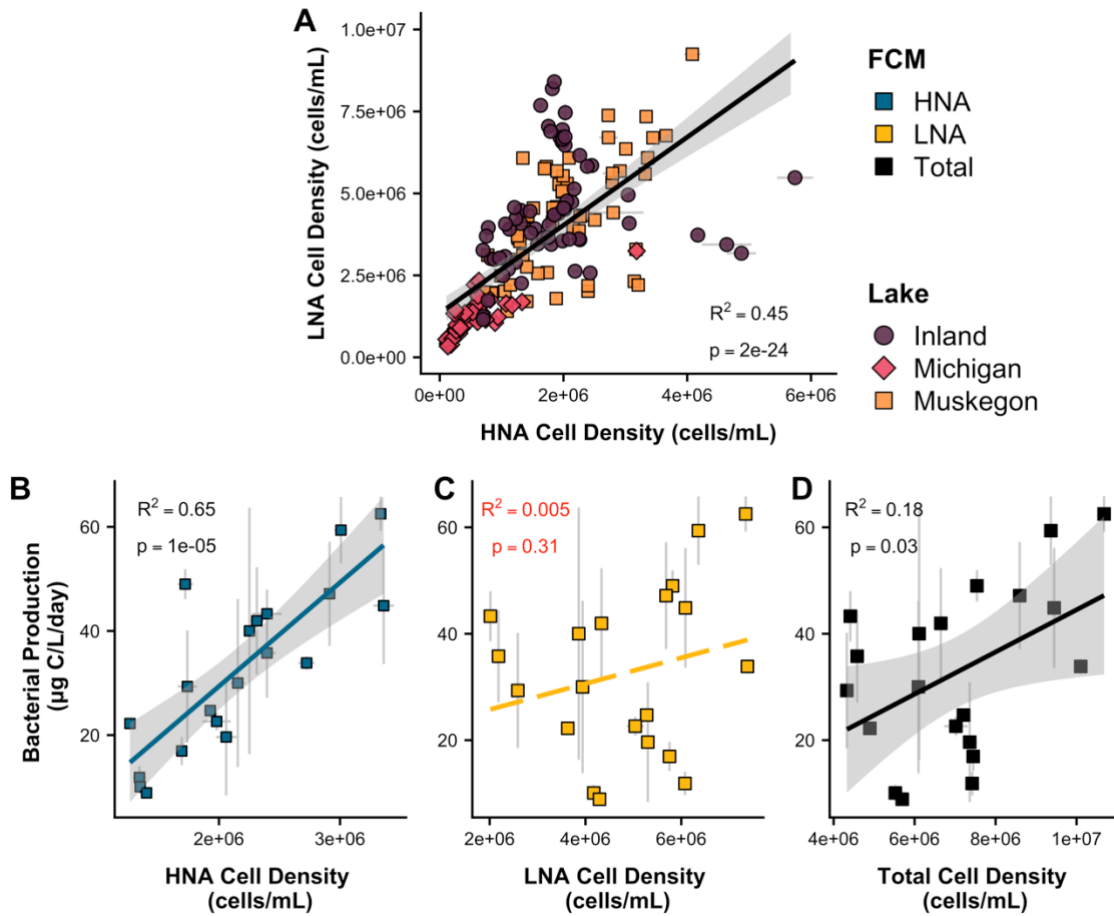
- 726 Acid Extraction Method Bias Perceived Bacterial Community Composition. *PLoS One*
727 2015; **10**: e0121659.
- 728 42. Louca S, Doebeli M, Parfrey LW. Correcting for 16S rRNA gene copy numbers in
729 microbiome surveys remains an unsolved problem. *Microbiome* 2018; **6**: 1–12.
- 730 43. Van der Gucht K, Cottenie K, Muylaert K, Vloemans N, Cousin S, Declerck S, et al. The
731 power of species sorting: local factors drive bacterial community composition over a wide
732 range of spatial scales. *Proc Natl Acad Sci U S A* 2007; **104**: 20404–20409.
- 733 44. Adams HE, Crump BC, Kling GW. Metacommunity dynamics of bacteria in an arctic
734 lake: The impact of species sorting and mass effects on bacterial production and
735 biogeography. *Front Microbiol* 2014; **5**: 1–10.
- 736 45. Chiang E, Schmidt ML, Berry MA, Biddanda BA, Burtner A, Johengen TH, et al.
737 Verrucomicrobia are prevalent in north- temperate freshwater lakes and display class-
738 level preferences between lake habitats. *PLoS One* 2018; **13**: 1–20.
- 739 46. Jones SE, Lennon JT. Dormancy contributes to the maintenance of microbial diversity.
740 *Proc Natl Acad Sci* 2010; **107**: 5881–5886.
- 741 47. Zimmerman R, Iturriaga R, Becker-Birck J. Simultaneous determination of the total
742 number of aquatic bacteria and the number thereof involved in respiration. *Appl Environ*
743 *Microbiol* 1978; **36**: 926–935.
- 744 48. Aanderud ZT, Vert JC, Lennon JT, Magnusson TW, Breakwell DP, Harker AR. Bacterial
745 dormancy is more prevalent in freshwater than hypersaline lakes. *Front Microbiol* 2016;
746 **7**: 1–13.
- 747 49. Jia X, Dini-Andreote F, Falcão Salles J. Community Assembly Processes of the Microbial
748 Rare Biosphere. *Trends Microbiol* 2018; **xx**: 1–10.
- 749 50. Amy PS, Morita RY. Starvation-survival patterns of sixteen freshly isolated open-ocean
750 bacteria. *Appl Environ Microbiol* 1983; **45**: 1109–1115.
- 751 51. Corno G, Jürgens K. Direct and indirect effects of protist predation on population size
752 structure of a bacterial strain with high phenotypic plasticity. *Appl Environ Microbiol*
753 2006; **72**: 78–86.
- 754 52. Coleman ML, Sullivan MB, Martiny AC, Steglich C, Barry K, Delong EF, et al. Genomic
755 Islands and the Ecology and Evolution of Prochlorococcus. *Science (80-)* 2006; **311**:
756 1768–1770.
- 757 53. Hunt DE, David L a, Gevers D, Preheim SP, Alm EJ, Polz MF. Resource partitioning and
758 sympatric differentiation among closely related bacterioplankton. *Science (80-)* 2008;
759 **320**: 1081–1085.
- 760 54. Denev VJ, Kalnejais LH, Mueller RS, Wilmes P, Baker BJ, Thomas BC, et al.

- 761 Proteogenomic basis for ecological divergence of closely related bacteria in natural
762 acidophilic microbial communities. *Proc Natl Acad Sci U S A* 2010; **107**: 2383–2390.
- 763 55. Shapiro BJ, Polz MF. Ordering microbial diversity into ecologically and genetically
764 cohesive units. *Trends Microbiol* 2014; **22**: 235–247.
- 765 56. Newton RJ, Jones SE, Eiler A, McMahon KD, Bertilsson S. A guide to the natural history
766 of freshwater lake bacteria. *Microbiology and molecular biology reviews* . 2011.
- 767 57. Woodhouse JN, Kinsela AS, Collins RN, Bowling LC, Honeyman GL, Holliday JK, et al.
768 Microbial communities reflect temporal changes in cyanobacterial composition in a
769 shallow ephemeral freshwater lake. *ISME J* 2016; **10**: 1337–1351.
- 770 58. Tada Y, Suzuki K. Changes in the community structure of free-living heterotrophic
771 bacteria in the open tropical Pacific Ocean in response to microalgal lysate-derived
772 dissolved organic matter. *FEMS Microbiol Ecol* 2016; **92**: 1–13.
- 773 59. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al.
774 Introducing mothur: open-source, platform-independent, community-supported software
775 for describing and comparing microbial communities. *Appl Environ Microbiol* 2009; **75**:
776 7537–7541.
- 777 60. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal
778 RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids*
779 *Res* 2013; **41**: 590–596.
- 780 61. Rohwer RR, Hamilton JJ, Newton RJ, McMahon KD. TaxAss: Leveraging Custom
781 Databases Achieves Fine-Scale Taxonomic Resolution. *bioRxiv* 2017; 214288.
- 782 62. Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, et al. Correlation
783 detection strategies in microbial data sets vary widely in sensitivity and precision. *ISME J*
784 2016; **10**: 1669–1681.
- 785 63. McMurdie PJ, Holmes S. phyloseq: An R Package for Reproducible Interactive Analysis
786 and Graphics of Microbiome Census Data. *PLoS One* 2013; **8**: e61217.
- 787 64. Kirchman D, K'nees E, Hodson R. Leucine incorporation and its potential as a measure of
788 protein synthesis by bacteria in natural aquatic systems. *Appl Environ Microbiol* 1985; **49**:
789 599–607.
- 790 65. Simon M, Azam F. Protein content and protein synthesis rates of planktonic marine
791 bacteria. *Mar Ecol Prog Ser* 1989; **51**: 201–213.
- 792 66. Props R, Schmidt ML, Heyse J, Vanderploeg HA, Boon N, Denev VJ. Flow cytometric
793 monitoring of bacterioplankton phenotypic diversity predicts high population-specific
794 feeding rates by invasive dreissenid mussels. *Environ Microbiol* 2017; **00**.
- 795 67. Prest EI, Hammes F, Köttsch S, van Loosdrecht MCM, Vrouwenvelder JS. Monitoring

- 796 microbiological changes in drinking water systems using a fast and reproducible flow
797 cytometric method. *Water Res* 2013; **47**: 7131–7142.
- 798 68. R Core Team. R: A Language and Environment for Statistical Computing. 2018. Vienna,
799 Austria.
- 800 69. Paliy O, Shankar V. Application of multivariate statistical techniques in microbial
801 ecology. *Mol Ecol* 2016; **25**: 1032–1057.
- 802 70. Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome datasets are
803 compositional: And this is not optional. *Front Microbiol* 2017; **8**: 1–6.
- 804 71. Quinn TP, Erb I, Richardson MF, Crowley TM. Understanding sequencing data as
805 compositions: an outlook and review. *Bioinformatics* 2018; 1–9.
- 806 72. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn:
807 Machine Learning in Python. *J Mach Learn Res* 2011; **12**: 2825–2830.
- 808 73. Probst P, Wright M, Boulesteix A-L. Hyperparameters and Tuning Strategies for Random
809 Forest. *arXiv* 2018; preprint.
- 810 74. Guyon I, Weston J, Barnhill S, Vapnik V. Gene Selection for Cancer Classification using
811 Support Vector Machines. *Mach Learn* 2002; **46**: 389–422.
- 812 75. Nilsson R, Peña JM, Björkegren J, Tegnér J. Consistent Feature Selection for Pattern
813 Recognition in Polynomial Time. *J Mach Learn Res* 2007; **8**: 589–612.
- 814 76. Roberts DR, Bahn V, Ciuti S, Boyce MS, Elith J, Guillera-Aroita G, et al. Cross-
815 validation strategies for data with temporal, spatial, hierarchical, or phylogenetic structure.
816 *Ecography (Cop)* 2017; **40**: 913–929.
- 817 77. Seabold S, Perktold J. Statsmodels: Econometric and Statistical Modeling with Python.
818 *Proc 9th Python Sci Conf* 2010; 57–61.
- 819 78. Price MN, Dehal PS, Arkin AP. FastTree 2 - Approximately maximum-likelihood trees
820 for large alignments. *PLoS One* 2010; **5**.
- 821 79. Harmon LJ, Weir JT, Brock CD, Glor RE, Challenger W. GEIGER: Investigating
822 evolutionary radiations. *Bioinformatics* 2008; **24**: 129–131.
- 823 80. Revell LJ. phytools: An R package for phylogenetic comparative biology (and other
824 things). *Methods Ecol Evol* 2012; **3**: 217–223.
- 825 81. Jombart T, Balloux F, Dray S. adephylo: New tools for investigating the phylogenetic
826 signal in biological traits. *Bioinformatics* 2010; **26**: 1907–1909.

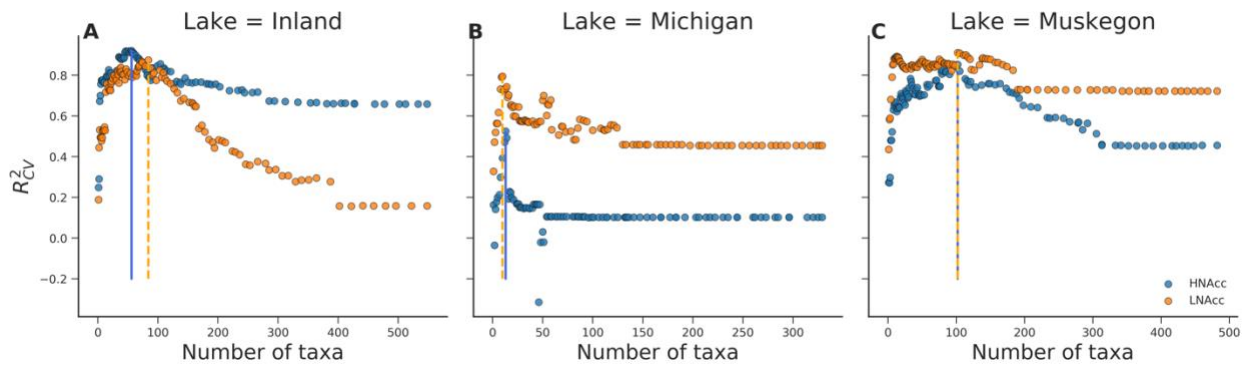
827

828 **Figure 1:** (A) Correlation between HNA cell counts and LNA cell counts across the three
829 freshwater lake ecosystems. (B-D) Muskegon Lake bacterial heterotrophic production and its
830 correlation with (B) HNA cell counts, (C) LNA cell counts, and (D) total cell counts. The grey
831 area in plots A, B, and D represents the 95% confidence intervals.



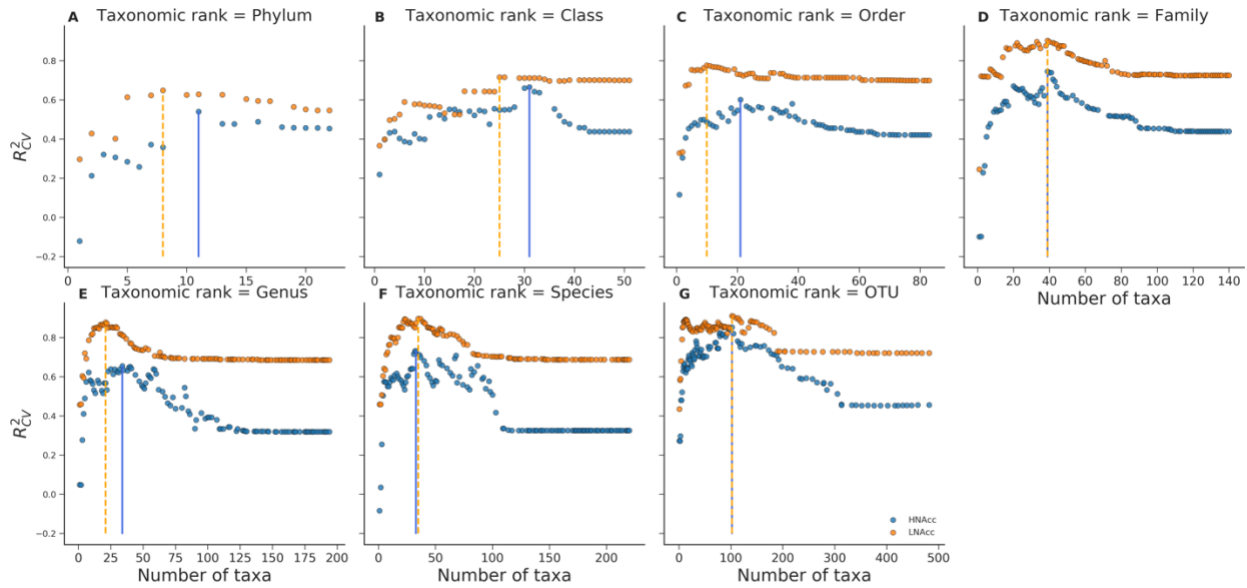
832

833 **Figure 2:** R_{CV}^2 in function of the number of OTUs, which were iteratively removed based on the
834 RL score and evaluated using the Lasso at every step. The solid (HNA) and dashed (LNA)
835 vertical lines corresponds to the threshold (i.e., number of OTUs) which resulted in a maximal
836 R_{CV}^2 . **(A)** Inland system ($R_{CV,max}^2 = 0.92$), HNAcc; **(B)** Lake Michigan ($R_{CV,max}^2 = 0.53$),
837 HNAcc; **(C)** Muskegon lake, HNAcc ($R_{CV,max}^2 = 0.85$); **(D)** Inland system, LNAcc (
838 $R_{CV,max}^2 = 0.87$); **(E)** Lake Michigan, LNAcc ($R_{CV,max}^2 = 0.79$); **(F)** Muskegon lake, LNAcc (
839 $R_{CV,max}^2 = 0.91$).



840

841 **Figure 3:** Evaluation of HNAcc and LNAcc predictions using the Lasso at all taxonomic levels
842 for the Muskegon lake system, expressed in terms of R_{CV}^2 , using different subsets of taxonomic
843 variables. Subsets were determined by iteratively eliminating the lowest-ranked taxonomic
844 variables based on the RL score.



845

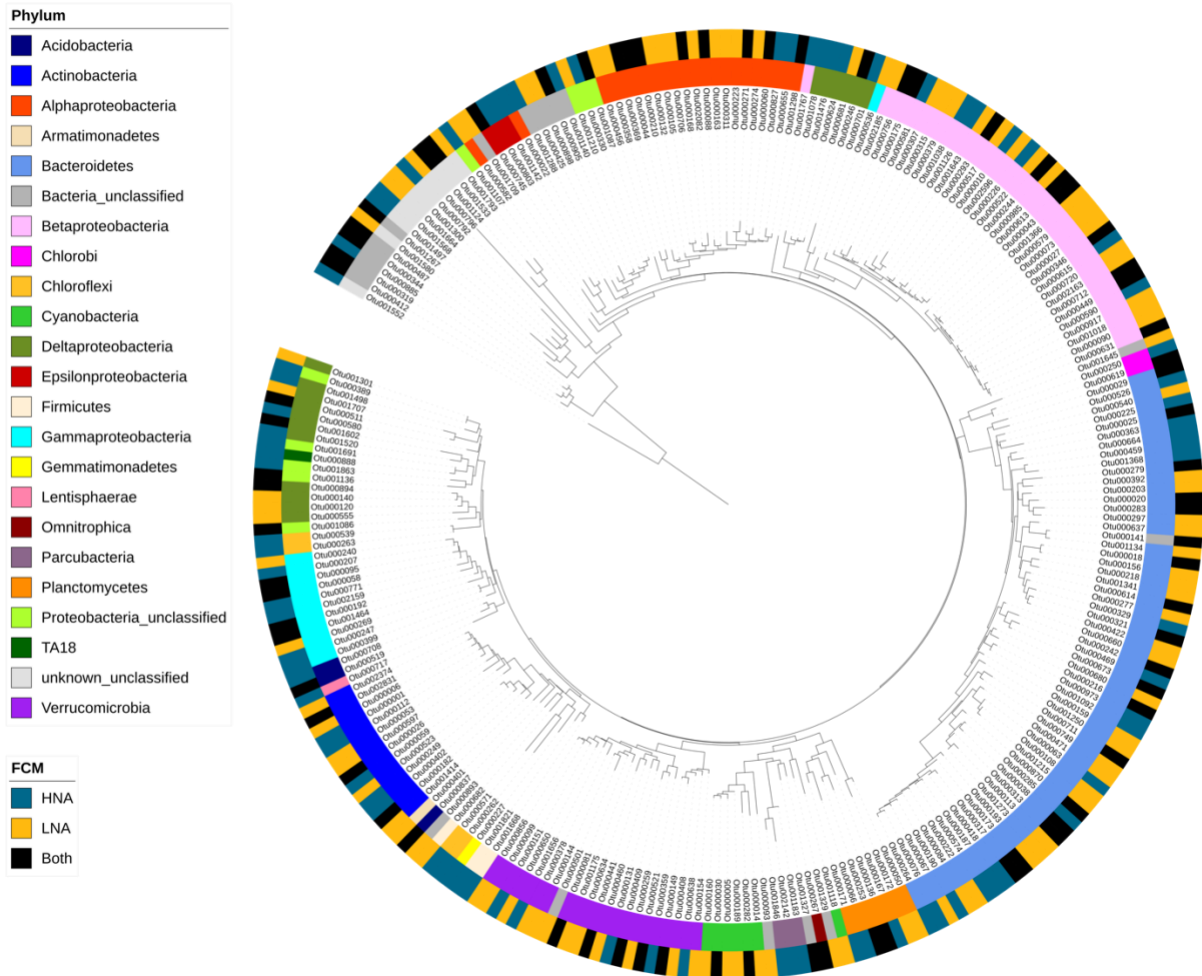
846 **Figure 4:** Hierarchical clustering of the RL score for the top 10 selected OTUs within each lake
847 system and FCM functional groups with the selected OTU (rows) across HNA and LNA groups
848 within the three lake systems (columns).



849

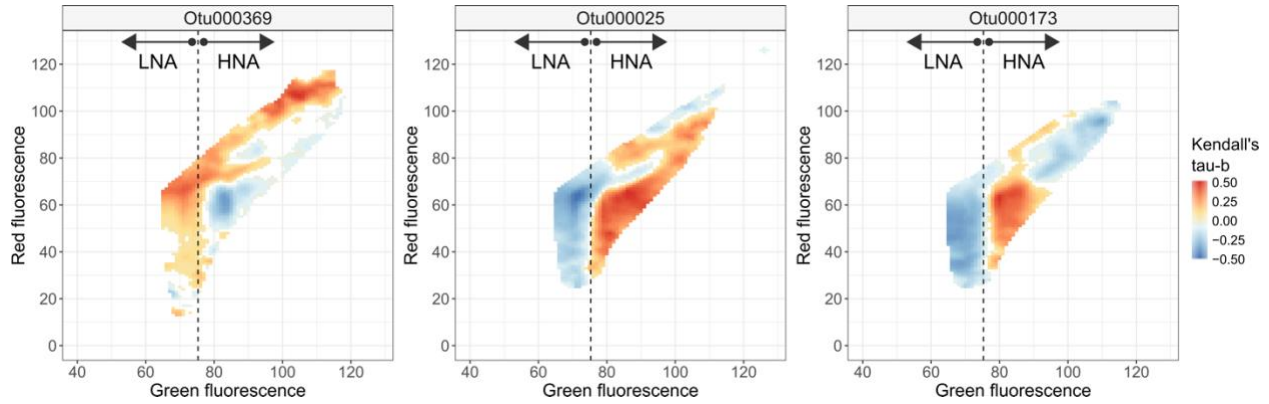
850 **Figure 5:** Phylogenetic tree with all HNA and LNA selected OTUs from each of the three lake
851 systems with their phylum level taxonomic classification and association with HNA, LNA or to
852 both groups based on the RL score threshold values.

Tree scale: 0.1



853

854 **Figure 6:** Correlation (Kendall's tau-b) between the relative abundances of the top three OTUs
855 selected by the RL and the densities in the cytometric fingerprint. The fluorescence threshold
856 used to define HNA and LNA populations is indicated by the dotted line.



857

858 **Table 1:** Top scored OTUs according to the RL per functional population and lake ecosystem.
859 Selection according to the Boruta algorithm is given in addition to the RL score. Descriptive
860 statistics by means of the Kendall rank correlation coefficient (KRCC) have been added with
861 level of significance in function of the HNA/LNA population. Full taxonomy of the OTUs is
862 given in **Table S2**.
863

Lake system	Functional group	OTU	RL score	Boruta selected	Kendall's tau (HNA)	P-value (HNA)	Kendall's tau (LNA)	P-value (LNA)
Inland	HNA	OTU369	0.382	yes	-0.43	<0.001	-0.28	0.0012
	LNA	OTU555	0.384	no	0.089	N.S.	0.22	0.011
Michigan	HNA	OTU025	0.362	yes	0.46	<0.001	0.41	<0.001
	LNA	OTU168	0.428	yes	0.26	0.0092	0.4	<0.001
Muskegon	HNA	OTU173	0.462	yes	0.5	<0.001	0.2	0.019
	LNA	OTU029	0.568	no	0.26	0.0029	0.49	<0.001

