

1 **Iroki: automatic customization for phylogenetic trees**

2

3 Ryan M. Moore¹, Amelia O. Harrison², Sean M. McAllister², Rachel L. Marine³, Clara S.

4 Chan², and K. Eric Wommack^{1*}

5

6 ¹Center for Bioinformatics & Computational Biology, University of Delaware, Newark, DE,

7 USA

8 ²School of Marine Science and Policy, University of Delaware, Newark, DE, USA

9 ³Department of Biological Sciences, University of Delaware, Newark, DE, USA

10

11 Corresponding author's information

12 *To whom correspondence should be addressed

13 **Address: Delaware Biotechnology Institute, 15 Innovation Way, Newark,**

14 **Delaware 19711**

15 **(Tel): (302) 831-4362**

16 **(Fax): (302) 831-3447**

17 **(E-mail): wommack@dbi.udel.edu**

18

19 **Abstract**

20 *Background*

21 Phylogenetic trees are an important analytical tool for examining species and community
22 diversity, and the evolutionary history of species. In the case of microorganisms, decreasing
23 sequencing costs have enabled researchers to generate ever-larger sequence datasets, which
24 in turn have begun to fill gaps in the evolutionary history of microbial groups. However,
25 phylogenetic analyses of large sequence datasets present challenges to extracting
26 meaningful trends from complex trees. Scientific inferences made by visual inspection of
27 phylogenetic trees can be simplified and enhanced by customizing various parts of the tree,
28 including label color, branch color, and other features. Yet, manual customization is time-
29 consuming and error prone, and programs designed to assist in batch tree customization
30 often require programming experience. To address these limitations, we developed Iroki, a
31 program for fast, automatic customization of phylogenetic trees. Iroki allows the user to
32 incorporate information on a broad range of metadata for each experimental unit
33 represented in the tree.

34

35 *Results*

36 Iroki was applied to four existing microbial sequence datasets to demonstrate its utility in
37 data exploration and presentation. Specifically, we used Iroki to highlight connections
38 between viral phylogeny and host taxonomy, explore the abundance of microbial groups
39 associated with Shiga toxin-producing *Escherichia coli* (STEC) in cattle, examine short-

40 term temporal dynamics of virioplankton communities, and to search for trends in the
41 biogeography of Zetaproteobacteria.

42

43 *Conclusions*

44 Iroki is an easy-to-use application having both command line and web-browser
45 implementations for fast, automatic customization of phylogenetic trees based on user-
46 provided categorical or continuous metadata. Iroki enables hypothesis testing through
47 improved visualization of phylogenetic trees, streamlining the process of biological sequence
48 data exploration and presentation.

49

50 *Availability*

51 Iroki can be accessed through a web browser application or via installation through
52 RubyGems, from source, or through the Iroki Docker image. All source code and
53 documentation is available under the GPLv3 license at <https://github.com/mooreryan/iroki>.
54 The Iroki web-app is accessible at www.iroki.net or through the VIROME portal
55 (<http://virome.dbi.udel.edu>), and its source code is released under GPLv3 license at
56 https://github.com/mooreryan/iroki_web. The Docker image can be found here:
57 <https://hub.docker.com/r/mooreryan/iroki>.

58

59 **Keywords**

60 Phylogeny, visualization, sequence analysis, bioinformatics, metagenomics

61

62 **Iroki: automatic customization for phylogenetic trees**

63

64 **Background**

65 Community and population ecological studies often use phylogenetic trees as a means for
66 assessing the diversity and evolutionary history of organisms. In the case of
67 microorganisms, the declining cost of sequencing has enabled researchers to gather ever-
68 larger sequence datasets from unknown microbial populations within environmental
69 samples. While large sequence datasets have begun to fill in the gaps in the evolutionary
70 history of microbial groups [1–5]; they have also posed new analytical challenges as
71 extracting meaningful trends within such highly dimensional datasets can be cumbersome.
72 In particular, scientific inferences made by visual inspection of phylogenetic trees can be
73 simplified and enhanced by customizing various parts of the tree including label and branch
74 color, branch width, and other features. Though several tree visualization packages allow
75 for manual modifications [6–9], the process can be time consuming and error prone
76 especially when the tree contains many nodes. Moreover, these packages are typically not
77 capable of batch customization without prior computer programming experience [10–13].
78
79 Iroki, a program for fast, automatic customization of phylogenetic trees, was developed to
80 address these limitations and enable users to incorporate a broad array of metadata
81 information for each experimental unit represented in the tree. Iroki is available for use
82 through a web browser interface at www.iroki.net, through the VIROME portal
83 (<http://virome.dbi.udel.edu>), and through a UNIX command line tool. Results are saved in

84 the widely used Nexus format with color metadata tailored for use with FigTree [8] (a freely
85 available and efficient tree viewer).

86

87 **Implementation**

88 Iroki enhances visualization of phylogenetic trees by coloring node labels and branches
89 according to categorical metadata criteria or numerical data such as abundance
90 information. Iroki can also rename nodes in a batch process according to user specifications
91 so that node names are more descriptive. A tree file in Newick format containing a
92 phylogenetic tree is always required. Additional required input files depend on the
93 operation(s) desired. Coloring functions require a color map or a biom [14] file. Node
94 renaming functions require a name map. The color map, name map, and biom files are
95 created by the user and, along with the Newick file, form the inputs for Iroki.

96

97 *Explicit tree coloring*

98 Iroki's principle functionality involves coloring node labels and/or branches based on
99 information provided by the user in the color map. The color map text file contains either
100 two or three tab-delimited columns depending on how branches and labels are to be
101 colored. Two columns, pattern and color, are used when labels and branches are to have
102 the same color. Three columns, pattern, label color, and branch color, are used when
103 branches and labels are to have different colors. Patterns are searched against node labels
104 either as regular expressions or exact string matches.

105

106 Entries in the color column can be any of the 657 named colors in the R programming
107 language [15] (e.g., skyblue, tomato, goldenrod2, lightgray, black) or any valid hexadecimal
108 color code (e.g., #FF78F6). In addition, Iroki provides a 19 color palette with
109 complementary colors based on Kelly's color scheme for maximum contrast [16]. Nodes on
110 the tree that are not in the color map will remain black.

111
112 Depending on user-specified options, a pattern match to node label(s) will trigger coloring of
113 the label and/or the branch directly connected to that label. Inner branches will be colored
114 to match their descendent branches if all descendants are the same color, allowing quick
115 identification of common ancestors and clades that share common metadata.

116
117 *Tree coloring based on numerical data*

118 Iroki provides the ability to generate color gradients based on numerical data, such as
119 absolute or relative abundance, from a tab-delimited biom format file [14]. Single-color
120 gradients use color saturation to illustrate numerical differences, with nodes at a higher level
121 being more saturated than those at a lower level. For example, highly abundant nodes will
122 be represented by more highly saturated colors. Two-color gradients show numerical
123 differences through both color mixing and luminosity. Additionally, the biom file may
124 specify numerical information for one group (e.g., abundance in a particular sample) or for
125 two groups (e.g., abundance in the treatment group vs. abundance in the control group).
126 For biom files with one group, single- or two-color gradients may be used. However, biom
127 files specifying two-group metadata may only use the two-color gradient.

128

129 *Renaming nodes*

130 Some packages for generating phylogenetic trees restrict the use of special characters and
131 spaces or require node names to be shorter than a specified length or (RAxML [17],
132 PHYLIP [18], etc.). Name restrictions present challenges to scientific interpretation of
133 phylogenetic trees. Iroki's renaming function uses a two-column, tab-delimited name map
134 to associate current node names, exactly matching those in the tree file, with new names.
135 The new name column has no restrictions on name length or character type. Iroki ensures
136 name uniqueness by appending integers to the ends of names, if necessary.

137

138 *Combining the color map, name map, and biom files*

139 Iroki can be used to make complex combinations of customizations by combining the color
140 map, name map, and biom files. For example, a biom file can be used to apply a color
141 gradient based on numerical data to the labels of a tree, a color map can be used to
142 separately color the branches based on user-specified conditions, and a name map can be
143 used to rename nodes in a single command or web request. Iroki follows a specific order of
144 precedence when applying multiple customizations. The color gradient inferred from the
145 biom file is applied first. Next, the color map is applied to specified labels or branches,
146 overriding the gradient applied in the previous step, if necessary. Finally, the name map is
147 used to map current names to the new names (Fig. 1).

148

149 *Output*

150 Iroki outputs the modified tree in the Nexus format. When building the phylogenetic tree,
151 FigTree uses the Nexus format file and interprets the color metadata output of Iroki.

152

153 **Results & Discussion**

154

155 *Global diversity of bacteriophage*

156 Viruses are the most abundant biological entity on Earth, providing an enormous reservoir
157 of genetic diversity, driving evolution of their hosts, influencing composition of microbial
158 communities, and affecting global biogeochemical cycles [19,20]. The viral taxonomic
159 system developed by the International Committee on Taxonomy of Viruses (ICTV) is based
160 on a suite of physical characteristics of the virion rather than on genome sequences. Noting
161 this limitation, the phage proteomic tree was created to provide a genome-based taxonomic
162 system for bacteriophage classification [21]. The phage proteomic tree was recently
163 updated to include hundreds of new phage genomes from the Phage SEED reference
164 database [22], as well as long assembled contigs from viral shotgun metagenomes (viromes)
165 collected from the Chesapeake Bay (SERC sample) [23] and the Mediterranean Sea [24].

166

167 Taxonomy and host information metadata was collected for the viral genome sequences, a
168 color map was created to assign colors based on viral family and host phyla, and Iroki was
169 used to add color metadata to branches and labels of the phage proteomic tree. Since a
170 large number of colors were required on the tree, Iroki's Kelly color palette was used to
171 provide clear color contrasts. The tree was rendered with FigTree (Fig. 2).

172
173 Adding color to the phage proteomic tree with Iroki shows trends in the data that would be
174 difficult to discern without color. Uncultured phage contigs from the Chesapeake Bay and
175 Mediterranean viromes make up a large portion of all phage sequences shown on the tree,
176 and are widely distributed among known phage. In general, viruses in the same family
177 claded together, e.g., branch coloring highlights large groups of closely related Siphoviridae
178 and Myoviridae. This label-coloring scheme also shows that viruses infecting hosts within
179 same phylum are, in general, phylogenetically similar. For example, viruses within one of
180 the multiple large groups of Siphoviridae across the tree infect almost exclusively host
181 species within the same phylum, e.g., Siphoviridae infecting Actinobacteria clade away from
182 Siphoviridae infecting Firmicutes or Proteobacteria.

183
184 *Bacterial community diversity and prevalence of E. coli in beef cattle*
185 Shiga toxin-producing *Escherichia coli* (STEC) are dangerous human pathogens that
186 colonize the lower gastrointestinal tracts of cattle and other ruminants. STEC-contaminated
187 beef and STEC shed in the feces of these animals are major sources of foodborne illness.
188 To identify possible interactions between STEC populations and the commensal cattle
189 microbiome, a recent study examined the diversity of the bacterial community associated
190 with beef cattle hide [25]. Fecal and hide samples were collected over twelve weeks and
191 SSU rRNA amplicon libraries were constructed and analyzed by Illumina sequencing [26].
192 The study indicated that the community structure of hide bacterial communities was altered
193 when the hides were positive for STEC contamination.

194

195 Iroki was used to visualize changes in the relative abundance of each cattle hide bacterial
196 OTU according to the presence or absence of STEC. A Mann-Whitney U test comparing
197 OTU abundance between STEC positive and STEC negative samples was performed, and
198 those bacterial OTUs showing a significant change in relative abundance ($p < 0.5$) were
199 placed on a phylogenetic tree according to the 16S rRNA sequence. Branches of the tree
200 were colored based on whether there was a significant change in relative abundance with
201 STEC contamination (red: $p < 0.1$, blue: $p \geq 0.1$). Node labels were colored along a
202 blue-green color gradient representing the abundance ratio of OTUs between samples with
203 STEC (blue) and without (green). Additionally, label luminosity was determined based on
204 overall abundance of each OTU (lighter: less abundant, darker: more abundant) (Fig. 3).
205 Iroki makes it clear that most OTUs on the tree showed a significant difference in
206 abundance (branch coloring) between STEC positive and STEC negative samples (node
207 coloring). Furthermore, we can see that most OTUs are at low abundance with only a few
208 highly abundant OTUs (label luminosity). The color gradient added by Iroki allows us to
209 see that the abundant OTUs were evolutionarily distant from one another and thus spread
210 out across many phylogenetic groups.

211

212 Iroki can be used to quickly test hypotheses without investing a large amount of time
213 annotating trees manually. A UPGMA tree was created based on unweighted UniFrac
214 distance [27] between 356 bacterial community profiles based on SSU rRNA amplicon
215 sequences from cattle hide and fecal samples (Fig. 4). Iroki was used to evaluate similarities

216 in sample bacterial communities according to the sampling location. Iroki colored branches
217 based on whether the sample originated from feces (blue) or from hide (red). The coloring
218 added by Iroki shows a clear partitioning of bacterial communities on the tree based on their
219 sampling location (hide or feces). However, four fecal samples grouped with hide samples,
220 and two hide samples grouped with fecal samples, highlighting the ability of Iroki to easily
221 identify good candidates for more in-depth examination. Additionally, Iroki was used to
222 illustrate a correlation between one of the most abundant bacterial families,
223 Ruminococcaceae, and sampling location. Iroki colored node labels with a color gradient
224 based on Ruminococcaceae family abundance, utilizing both a single color gradient (Fig.
225 4A) and a two color gradient (Fig. 4B). Custom trees were visualized using FigTree. Iroki's
226 automatic color gradient and ability to label branches and nodes based on different criteria
227 clearly show that Ruminococcaceae is more abundant in fecal samples than in hide
228 samples.

229

230 *Short-term dynamics of viroplankton*

231 The gene encoding Ribonucleotide reductase (RNR) is common within viral genomes and
232 thus can be used as a marker gene for studying viral diversity [23]. Moreover, RNR
233 polymorphism is predictive of some of the biological and ecological features of viral
234 populations [28]. A mesocosm experiment examined the short-term dynamics of phage
235 populations using RNR amplicon sequences, specifically, sequences of class II RNRs of
236 bacteriophages infecting cyanobacterial hosts. A phylogenetic tree was created from the
237 Cyano II RNR amplicon sequences and Iroki was used to color nodes and branches based

238 on the time point (0 h, 6 h, 12 h) at which each amplicon sequence was observed. The
239 customized tree was then visualized using FigTree (Fig. 5). Iroki's coloring showed that no
240 phylogenetic clade was dominated by OTUs observed in any particular time point; rather,
241 time points were spread relatively evenly across clades. This analysis demonstrates Iroki's
242 utility for exploring sequence datasets, allowing the researcher to quickly and easily test
243 hypotheses.

244

245 *Phylogeny of Zetaproteobacteria within a biogeographic context*

246 Biogeographical studies assess the distribution of an organism's biodiversity across space
247 and time. The extent to which microorganisms exhibit geographic distribution patterns is an
248 open question in microbial ecology. The isolated nature of the microbial communities
249 associated with deep-ocean hydrothermal vents provides an ideal system for studying the
250 biogeography of microbes. In particular, iron-oxidizing bacteria have been shown to thrive
251 in vent fluids, sediments, and iron-rich microbial mats associated with the vents. Globally,
252 iron-oxidizing bacteria make significant contributions to the iron and carbon cycles. A
253 recent study analyzed multiple SSU rRNA clone libraries to investigate the biogeography of
254 Zetaproteobacteria, a phyla containing many iron-oxidizing bacterial species, between three
255 sampling regions of the Pacific Ocean (central Pacific—Loihi seamount, western Pacific—
256 Southern Mariana Trough, and southern Pacific (Vailulu'u Seamount/Tonga Arc/East Lau
257 Spreading Center/Kermadec Arc) [29]. Sequences were aligned and a phylogeny was
258 inferred as described in [29]. Iroki was used to examine the relationship between sampling
259 location and phylotype by adding branch and label color based on geographic location and

260 renaming original node labels with OTU and location metadata. The custom tree was
261 visualized using FigTree (Fig. 6). In some cases, OTUs contained sequences from only one
262 sampling location (e.g., OTUs 12, 15, and 16), whereas other OTUs are distributed among
263 more than one sampling location (e.g., OTUs 1, 2, and 4). Often, sequences sampled from
264 the same geographic location are in the same phylotype despite being members of different
265 OTUs (e.g., OTUs 10 and 19).

266

267 *Availability and requirements*

268 A web browser version of Iroki can be accessed online at www.iroki.net or through the
269 VIROME portal (<http://virome.dbi.udel.edu/>). For users who wish to run Iroki locally, a
270 command line version of the program is installable via RubyGems, from GitHub
271 (<https://github.com/mooreryan/iroki>). A Docker image is available for users who desire the
272 flexibility of the command line tool, but do not want to install Iroki or manage its
273 dependencies (<https://hub.docker.com/r/mooreryan/iroki>). Docker is a popular software
274 container platform that allows bundling of an application with its dependencies in a
275 portable, self-contained system [30,31]. The README file, accompanying the source code,
276 provides detailed instructions for setting up and running Iroki. Further documentation and
277 tutorials can be found at the Iroki Wiki (<https://github.com/mooreryan/iroki/wiki>).

278

279 *License*

280 Iroki and its associated programs are released under the GNU General Public License
281 version 3 [32].

282

283 **Conclusions**

284 Iroki is a command line program and web browser application for fast, automatic
285 customization of large phylogenetic trees based on user specified configuration files
286 describing categorical or continuous metadata information. The output files include Nexus
287 tree files with color metadata tailored specifically for use with FigTree. Various example
288 datasets from microbial ecology studies were analyzed to demonstrate Iroki's utility. In each
289 case, Iroki simplified the processes of data exploration, data presentation, and hypothesis
290 testing. Though these examples focused specifically on applications in microbial ecology,
291 Iroki is applicable to any problem space with hierarchical data that can be represented in
292 the Newick tree format. Iroki provides a simple and convenient way to rapidly customize
293 trees, especially in cases where the tree in question is too large to annotate manually or in
294 studies with many trees to annotate.

295

296 **List of Abbreviations**

297 OTU: operational taxonomic

298 RNR: Ribonucleotide reductase

299 STEC: Shiga-toxigenic *Escherichia coli*

300

301 **Ethics approval and consent to participate**

302 Not applicable

303

304 **Consent for publication**

305 Not applicable

306

307 **Availability of data and materials**

308 Data and code used to generate figures are available on GitHub at

309 https://github.com/mooreryan/iroki_manuscript_data

310

311 **Funding**

312 This project was supported by the USDA National Institute of Food and Agriculture award

313 number 2012-68003-30155 and the National Science Foundation Advances in

314 Bioinformatics program (award number DBI_1356374).

315

316 **Competing Interests**

317 The authors declare that they have no competing interests.

318

319 **Authors' contributions**

320 RMM and SMM conceived the project. RMM wrote the manuscript and implemented Iroki.

321 AOH and RLM processed and analyzed Cyano II amplicons. All authors read, edited, and

322 approved the final manuscript.

323

324 **Acknowledgements**

325 We would like to acknowledge Daniel J. Nasko and Jessica M. Chopyk for their work on the
326 phage proteomic tree, and Barbra D. Ferrell for editing the manuscript.
327

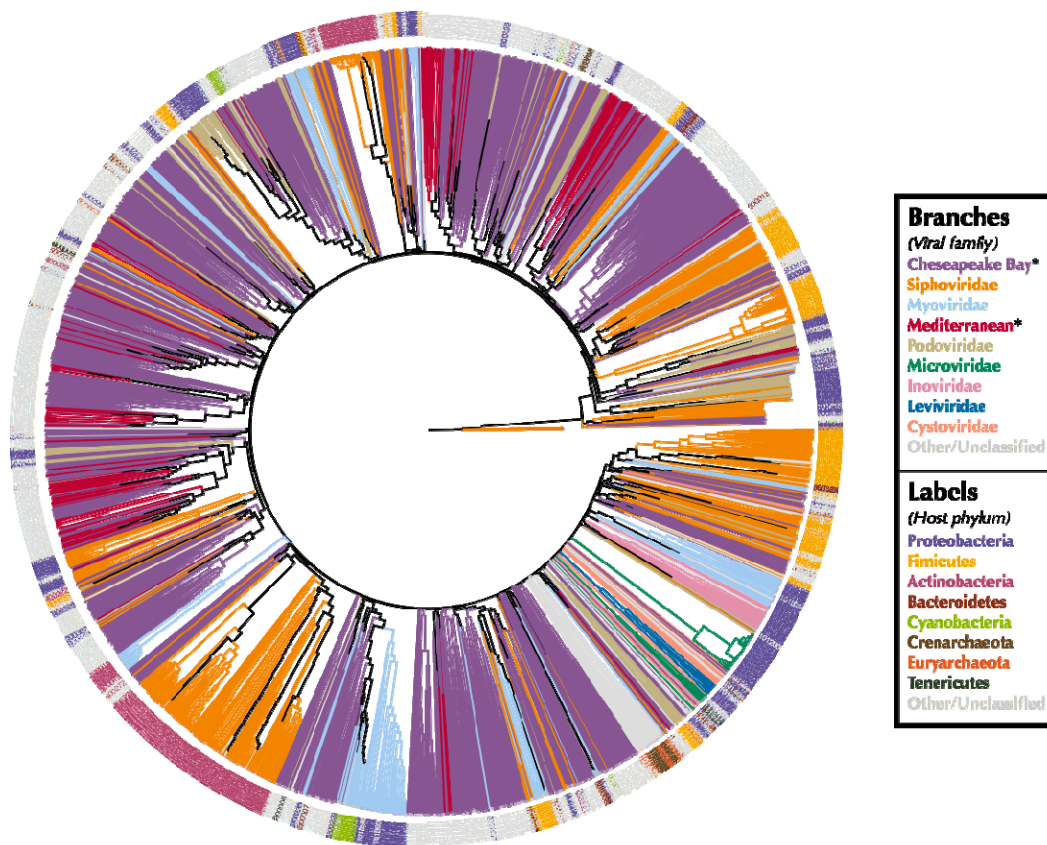
328 **References**

- 329 1. Lan Y, Rosen G, Hershberg R. Marker genes that are less conserved in their sequences
330 are useful for predicting genome-wide similarity levels between closely related prokaryotic
331 strains. *Microbiome*. 2016;4:18.
- 332 2. Larkin A a, Blinebry SK, Howes C, Lin Y, Loftus SE, Schmaus CA, et al. Niche
333 partitioning and biogeography of high light adapted *Prochlorococcus* across taxonomic
334 ranks in the North Pacific. *ISME J*. 2016;1–13.
- 335 3. Simister RL, Deines P, Botté ES, Webster NS, Taylor MW. Sponge-specific clusters
336 revisited: A comprehensive phylogeny of sponge-associated microorganisms. *Environ.*
337 *Microbiol*. 2012;14:517–24.
- 338 4. Wu Z, Yang L, Ren X, He G, Zhang J, Yang J, et al. Deciphering the bat virome catalog
339 to better understand the ecological diversity of bat viruses and the bat origin of emerging
340 infectious diseases. *ISME J*. 2016;10:609–20.
- 341 5. Müller AL, Kjeldsen KU, Rattei T, Pester M, Loy A. Phylogenetic and environmental
342 diversity of DsrAB-type dissimilatory (bi)sulfite reductases. *ISME J*. 2015;9:1152–65.
- 343 6. University W. Phylogeny Programs.
344 <http://evolution.genetics.washington.edu/phylip/software.html#Plotting>. Accessed 2016 Jul
345 21.
- 346 7. Zhang H, Gao S, Lercher MJ, Hu S, Chen WH. EvolView, an online tool for visualizing,
347 annotating and managing phylogenetic trees. *Nucleic Acids Res*. 2012;40.
- 348 8. Rambaut A. FigTree. <http://tree.bio.ed.ac.uk/software/figtree/>. Accessed 2016 Jul 21.
- 349 9. Zmasek CM. Archaeopteryx.

- 350 <https://sites.google.com/site/cmzmasek/home/software/archaeopteryx>. Accessed 2016 Jul 21.
- 351 10. Paradis E, Claude J, Strimmer K. APE: Analyses of phylogenetics and evolution in R
352 language. *Bioinformatics*. 2004;20:289–90.
- 353 11. Revell LJ. phytools: An R package for phylogenetic comparative biology (and other
354 things). *Methods Ecol. Evol.* 2012;3:217–23.
- 355 12. Huerta-Cepas J, Dopazo J, Gabaldón T. ETE: a python Environment for Tree
356 Exploration. *BMC Bioinformatics*. 2010;11:24.
- 357 13. Chen W-H, Lercher MJ, Ganfomina M, Gutierrez G, Bastiani M, Sanchez D, et al.
358 ColorTree: a batch customization tool for phylogenetic trees. *BMC Res. Notes. BioMed*
359 *Central*; 2009;2:155.
- 360 14. McDonald D, Clemente JC, Kuczynski J, Rideout J, Stombaugh J, Wendel D, et al. The
361 Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love
362 the ome-ome. *Gigascience*. 2012;1:7.
- 363 15. Ripley BD. *The R project for statistical computing*. 2001. p. 1–3.
- 364 16. Kelly KL. Twenty-two colors of maximum contrast. *Color Eng.* 1965. p. 26–7.
- 365 17. Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of
366 large phylogenies. *Bioinformatics*. 2014;30:1312–3.
- 367 18. Felsenstein J. PHYLIP. <http://evolution.gs.washington.edu/phytip.html>. Accessed 2016
368 Jul 21.
- 369 19. Suttle CA. Marine viruses — major players in the global ecosystem. *Nat. Rev. Microbiol.*
370 *Nature Publishing Group*; 2007;5:801–12.
- 371 20. Rohwer F, Thurber RV. Viruses manipulate the marine environment. *Nature. Nature*

- 372 Publishing Group; 2009;459:207–12.
- 373 21. Rohwer F, Edwards R. The Phage Proteomic Tree: a genome-based taxonomy for
374 phage. *J. Bacteriol.* 2002;184:4529–35.
- 375 22. Phage SEED. <http://www.phantome.org/PhageSeed/Phage.cgi>. Accessed 2016 Jul 21.
- 376 23. Wommack KE, Nasko DJ, Chopyk J, Sakowski EG. Counts and sequences,
377 observations that continue to change our understanding of viruses in nature. *J. Microbiol.*
378 2015;53:181–92.
- 379 24. Mizuno CM, Rodriguez-Valera F, Kimes NE, Ghai R. Expanding the marine virosphere
380 using metagenomics. *PLoS Genet. Public Library of Science*; 2013;9:e1003987.
- 381 25. Chopyk J, Moore RM, DiSpirito Z, Stromberg ZR, Lewis GL, Renter DG, et al. Presence
382 of pathogenic *Escherichia coli* is correlated with bacterial community diversity and
383 composition on pre-harvest cattle hides. *Microbiome. BioMed Central*; 2016;4:9.
- 384 26. Fadrosch DW, Ma B, Gajer P, Sengamalay N, Ott S, Brotman RM, et al. An improved
385 dual-indexing approach for multiplexed 16S rRNA gene sequencing on the Illumina MiSeq
386 platform. *Microbiome.* 2014;2:6.
- 387 27. Lozupone C, Knight R. UniFrac: a New Phylogenetic Method for Comparing Microbial
388 Communities. *Appl. Environ. Microbiol. American Society for Microbiology*; 2005;71:8228–
389 35.
- 390 28. Sakowski EG, Munsell E V., Hyatt M, Kress W, Williamson SJ, Nasko DJ, et al.
391 Ribonucleotide reductases reveal novel viral diversity and predict biological and ecological
392 features of unknown marine viruses. *Proc. Natl. Acad. Sci. National Academy of Sciences*;
393 2014;111:15786–91.

- 394 29. McAllister SM, Davis RE, McBeth JM, Tebo BM, Emerson D, Moyer CL. Biodiversity
395 and emerging biogeography of the neutrophilic iron-oxidizing Zetaproteobacteria. *Appl.*
396 *Environ. Microbiol. American Society for Microbiology (ASM)*; 2011;77:5445–57.
- 397 30. Biodocker. <http://biodocker.org/>. Accessed 2016 Jul 21.
- 398 31. Merkel D. Docker: lightweight Linux containers for consistent development and
399 deployment. *Linux J. Belltown Media*; 2014;2014:2.
- 400 32. GNU Operating System. <http://www.gnu.org/licenses/>. Accessed 2016 Jul 21.
- 401



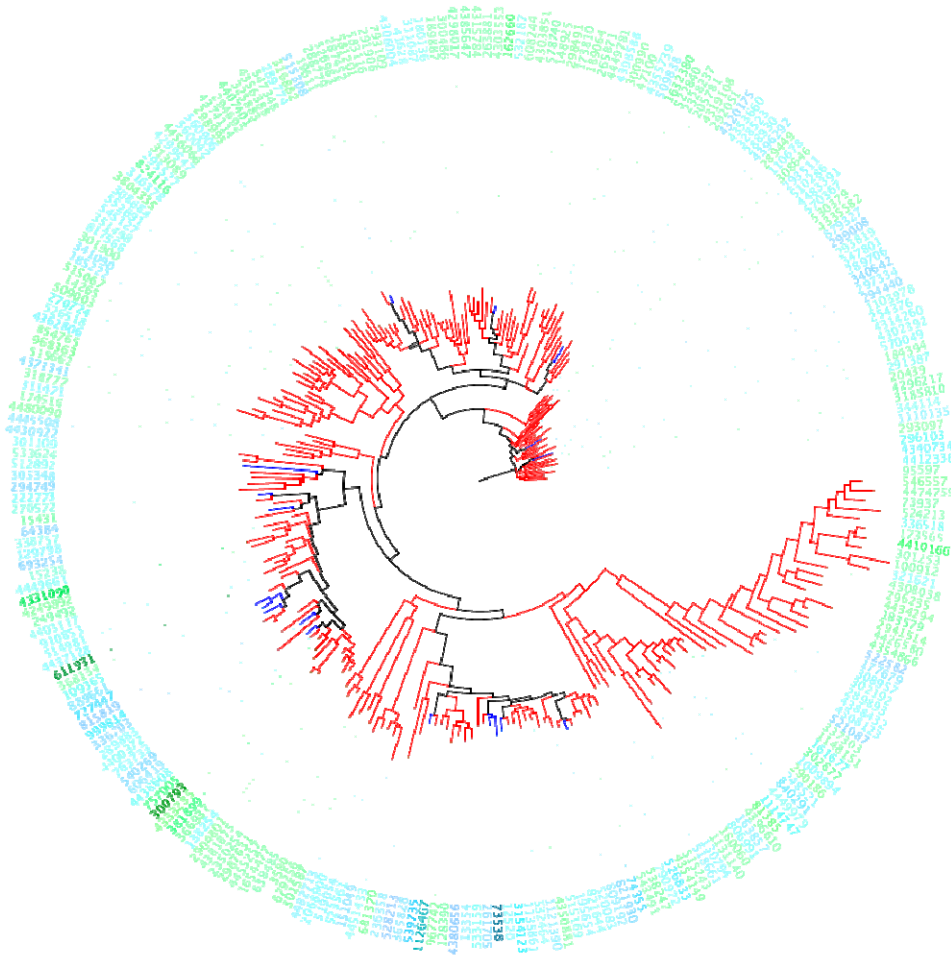
402

0.9

403 **Fig. 2: Comparing phage and their host phyla**

404 All phage genomes from Phage SEED with assembled virome contigs from the Chesapeake
405 Bay and Mediterranean Sea. Iroki highlights phylogenetic trends after coloring branches
406 according to viral family or sampling location in the case of virome contigs (marked with an
407 asterisk in the legend), and coloring node labels according to host phylum of the phage.

408



409

0.1

410 **Fig. 3: Changes in OTU abundance in two sample groups**

411 Approximate-maximum likelihood tree of OTUs that showed significant differences in
412 relative abundance between STEC positive and STEC negative cattle hide samples.

413 Branches show significance based on coloring by the p-value of a Mann-Whitney U test

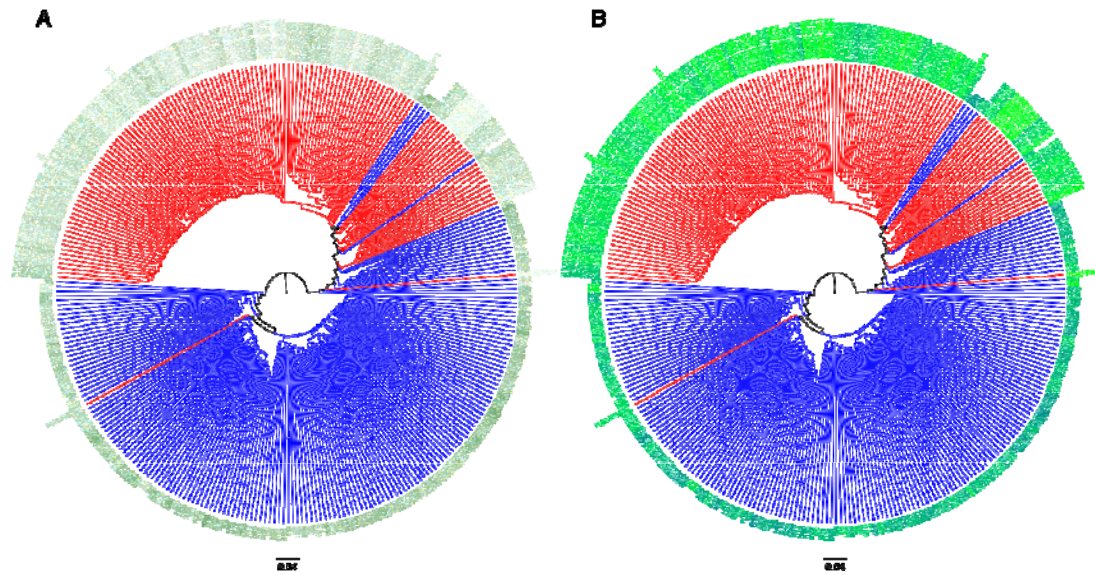
414 examining changes in abundance between samples positive for STEC ($p < 0.1$ – red) and

415 samples negative for STEC, ($p \geq 0.1$ – blue). Label color on a blue-green gradient

416 highlights OTU occurrence based on the abundance ratio between STEC positive samples

417 (green) and STEC negative samples (blue). Labels that are darker green had a higher

418 abundance in STEC positive samples, and a lower abundance in STEC negative samples.
419 For example, OTU 300793 (bottom left corner) is darker than most (indicating high overall
420 abundance) and more green than blue (indicating higher abundance in STEC positive
421 samples than in STEC negative samples). Node luminosity represents overall abundance
422 with lighter nodes being less abundant than darker nodes.

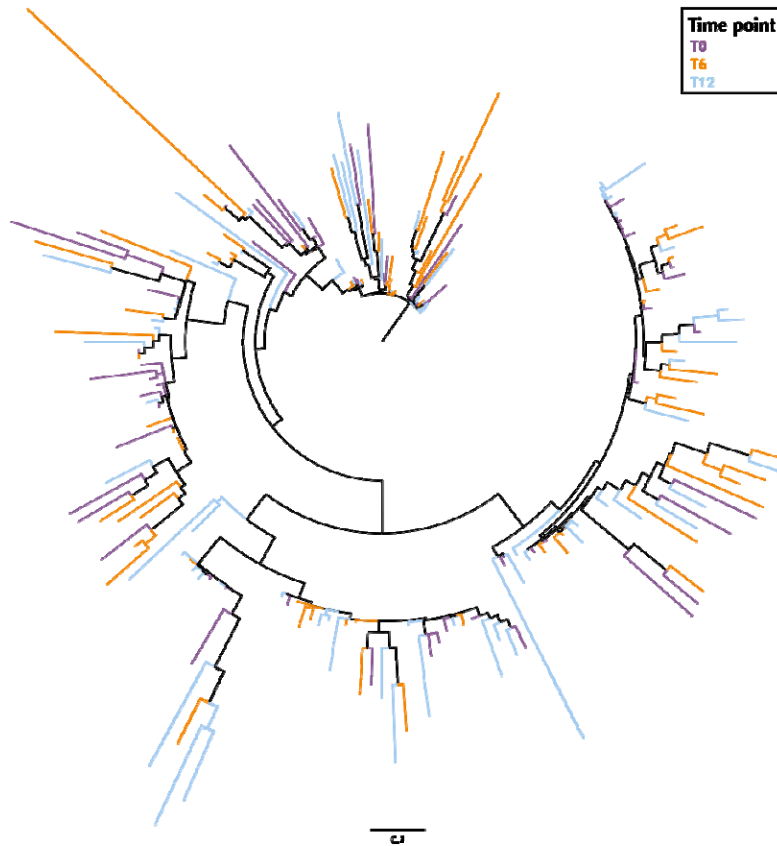


423

424 **Fig. 4: Comparing cattle fecal and hide samples and the abundance of**
425 **Ruminococcaceae**

426 Phylogeny based on UPGMA tree of pairwise unweighted UniFrac distance between 356
427 bacterial community profiles based on SSU rRNA amplicon sequences from cattle hide and
428 feces. Branches are colored by feces (blue) and hide (red). Rapid testing of the hypothesis
429 that the abundance of one of the most abundant families, Ruminococcaceae, and sample
430 origin are correlated is enabled through node label coloring by (A) a green single-color
431 gradient (color saturation increases with increasing abundance of Ruminococcaceae OTUs)
432 and (B) a light green (low abundance of Ruminococcaceae OTUs) to dark blue (high
433 abundance of Ruminococcaceae OTUs) color gradient.

434



435

436 **Fig. 5: Temporal dynamics of virioplankton populations according to Cyano II**

437 **RNR amplicon phylogeny**

438 An approximately-maximum-likelihood phylogenetic tree of 200 randomly selected class II

439 Cyano RNR representative sequences from 98% percent clusters. Iroki was used to color

440 branches by time point: zero hours – purple, six hours – orange, and twelve hours – blue.

