

1 **Next-Generation Teaching: A Template For Bringing Genomic and Bioinformatic Tools**

2 **into the Classroom**

3

4 **Supplementary Materials**

5 **Supporting Information**

6 *Bioinformatics*

7 As a first step in exploring our NGS data and to get students introduced to using a UNIX
8 command-line computing environment, we merged our forward and reverse reads (the two
9 directions of sequencing reads generated by the paired-end approach described earlier) using the
10 program FLASH (Magoč and Salzberg 2011) with settings: -m 20 -M 250, which specified the
11 minimum overlap for merging (-m) as 20 basepairs (bp) and the maximum overlap (-M) as 250
12 bp. The purpose of this step was to improve the quality and length of our sequence data by
13 merging forward and reverse reads together into one combined sequence. Next, we filtered our
14 data to remove low-quality (i.e., those reads that were unreliable due to failing the Illumina
15 “chastity filter” or having low-quality sequence scores) reads using the FASTX-Toolkit package
16 (Hannon, 2010). The FASTX-Toolkit package is a collection of command line tools, designed to
17 work with processing ‘short-read’ sequence data, the type of data generated by Illumina NGS
18 platforms. Specifically, we used the fastq_quality_filter script with settings: -q 25 -p 80, which
19 meant that 80% of all bases must have a Q33 quality score of 25 or greater.

20 After read merging and quality filtering, we moved on to the analysis of microbial diversity
21 in our data. For all analyses, we utilized the QIIME (pronounced ‘chime’) pipeline (Caporaso et
22 al. 2010) due to its popularity in the field, integrated ‘pipeline’ (meaning many modules act on
23 the output of previous components of the program), and the ease of generating output plots,

24 community statistics, and taxonomic identifications. Another benefit of the QIIME pipeline is a
25 well-documented online manual (<http://qiime.org/>) that provides overviews of specific tools, data
26 analysis tutorials, and example files included with the program source code. As a pre-packaged,
27 well-documented but still completely command-line driven framework, QIIME provides a
28 relatively easy-to-use introduction to command-line programming while also providing the tools
29 to properly analyze 16S rRNA microbial data.

30 In QIIME, we first converted our sequence data from FASTQ to FASTA format using the
31 `convert_fastaqual_fastq.py` script, as we no longer need quality scores for our data after initial
32 filtering (see an example script for this in the *.zip file included with this study's Supplementary
33 Materials). Next, we built a "mapping file" per the QIIME documentation instructions. This file
34 specified which libraries corresponded to which sampling treatments (an example mapping file is
35 also included in the *.zip file). We checked this mapping file for errors using the
36 `check_id_map.py` script. Note: This step almost always generates at least one error. If you don't
37 see more than a single error, your mapping file is likely fine.

38 Next, we added QIIME specific labels to each of our sequence libraries and combined them
39 all into one large FASTA file using `add_qiime_labels.py`. We picked operational taxonomic units
40 (OTUs) using the script `pick_de_novo_otus.py`, which compares the sequence data contained in
41 the combined FASTA file with publicly available 16S rRNA sequence databases (e.g.,
42 Greengenes) to assign taxonomy to the full dataset (and each sampling treatment). On a 2015
43 MacBook Pro, this step took approximately one hour to execute. The OTU table is generated as a
44 Biological Observation Matrix (BIOM) and we generated basic summary stats BIOM table using
45 the `summarize-table` script. In 2014 and 2015, we removed libraries from all analyses when
46 sequencing yield was less than 30K sequences for a given library. In 2016, we reduced this

47 number because we ran a “Nano” MiSeq run instead of the full, more expensive run option. We
48 elected to make this change to save money as we were no longer pursuing publication-quality
49 results, and rather, focusing on student learning objectives solely. Next, we summarized
50 communities by taxonomic composition with the `summarize_taxa_through_plots.py` script.
51 Depending on number of treatments (e.g., for 2014 our treatments were: Treatment 1 = gender,
52 Treatment 2 = location, Treatment 3 = restroom surface) this script can be run for the full dataset
53 then for as many individual treatments as necessary by invoking the `-c` flag. We computed alpha
54 and beta diversity using the scripts `alpha_rarefaction.py` and `beta_diversity_through_plots.py`.
55 The last three scripts produce a multitude of visual output in the form of bar and line graphs as
56 well as principal coordinate analyses. This output can be opened as an HTML file in a web
57 browser (e.g., Google Chrome), allowing students to interact with the data by selecting various
58 treatments, diversity metrics, or other choices. It should be noted that the outline described here
59 is only a subset of the QIIME pipeline, many additional tools for curating, analyzing, and
60 exploring the data are available.

61

62 *Removing contaminant OTUs*

63 One particularly useful component of the QIIME pipeline for education is the
64 investigation of possible contaminants OTUs. By collecting (and sequencing) the necessary
65 negative controls (e.g., field, classroom, and PCR), it is possible to identify OTUs at each stage,
66 consider them as a class (e.g., the bacterial genus *Pseudomonas* is associated with human acne
67 and a common contaminant in 16S experiments), and remove the ones that appear likely to be a
68 product of contamination. This process is outlined in Step 6 of the accompanying “rRNA
69 Instructor Analysis Guide.sh” although it should be noted that we have not incorporated into our

70 course materials. Briefly, the process begins by looking in the taxa summary data for each
71 negative control (whether PCR, classroom, or field) and identifying if particularly dominant
72 OTUs are contaminants or not. This isn't always clear and a combination of literature and
73 internet searching can typically help identify obvious ones. Next, the contaminant OTUs need to
74 be identified and copied into a list for removal. Following this step, the script
75 `filter_otus_from_otu_table.py` along with the list for removal can be combined to create a new
76 data table the contaminant OTUs removed. A follow-up summary of the output table using the
77 command 'biom summarize-table' is useful for confirming the contaminant OTUs have been
78 removed (e.g., there will be a noticeable before and after difference when comparing pre- and
79 post-contaminant filtering summaries).

80

81 **Results**

82 For UK bathrooms, when microbial diversity was compared across genders, no clear
83 pattern emerged (Figure S1A). However, when the same samples were grouped by location,
84 some loose, location-specific clustering did emerge, specifically for William T. Young Library
85 and the Student Center (Figure S1B). At the UK climbing wall, no clear patterns emerged for
86 microbial diversity across difficulty or type of holds (Figure S1C-D). One outlier sample,
87 footholds of 'Hard' routes was disconnected from the rest of the samples (Figure S1C-D). Across
88 years, a much wider variability in microbial diversity was observed for 2014 samples as the top
89 two principal components (PCs) explained 30.8% and 8.2% of the variation whereas in 2015, the
90 top two PCs explained 9.9% and 7.3% of the variability among samples. Comparisons of
91 between sample (alpha) diversity further show the difference between projects in variability
92 among treatments. In general, between treatment diversity was much lower in 2015 versus 2014

93 (Figure S2). When biodiversity metrics were compared across treatments and iterations of the
94 experience, no major patterns emerged. For Shannon diversity (SD), the highest diversity
95 observed in 2014 was for sink samples from the Main building in male bathrooms (SD = 6.61;
96 Table S3). For 2015, the highest diversity in a single treatment was observed for a foothold high
97 on an easy route (SD = 6.41; Table S4) and the five highest diversities observed were all on easy
98 routes (Table S4).

99

100 **Discussion**

101 *Additional analytical directions*

102 In this CURE, many analytical stones were left unturned. This was largely a product of
103 the course objective to be a component in a larger course that was not solely focused on genomic
104 education. With more time, or in a post-course independent research setting, a number of
105 complementary analyses could be explored. Chief among these is the question of contamination.
106 Built into the QIIME pipeline are tools for bioinformatically removing OTUs called from
107 negative or control samples from the broader OTU database. For instance, in 2015, it would be
108 possible to remove those OTUs associated with ‘clean’ climbing hold from all of the OTU
109 libraries linked to treatment holds and thereby clarify diversity with this common library of
110 control OTUs removed.

111 Furthermore, QIIME provides a mechanism for exploring specific taxonomic diversity at
112 any level of organization from kingdom to species. Taking a taxonomic-focused perspective on
113 the data, perhaps with a human health slant, would provide students with another outlet for
114 understanding where microorganisms occur and the frequency with which putative pathogens,
115 for instance, are present in the environment.

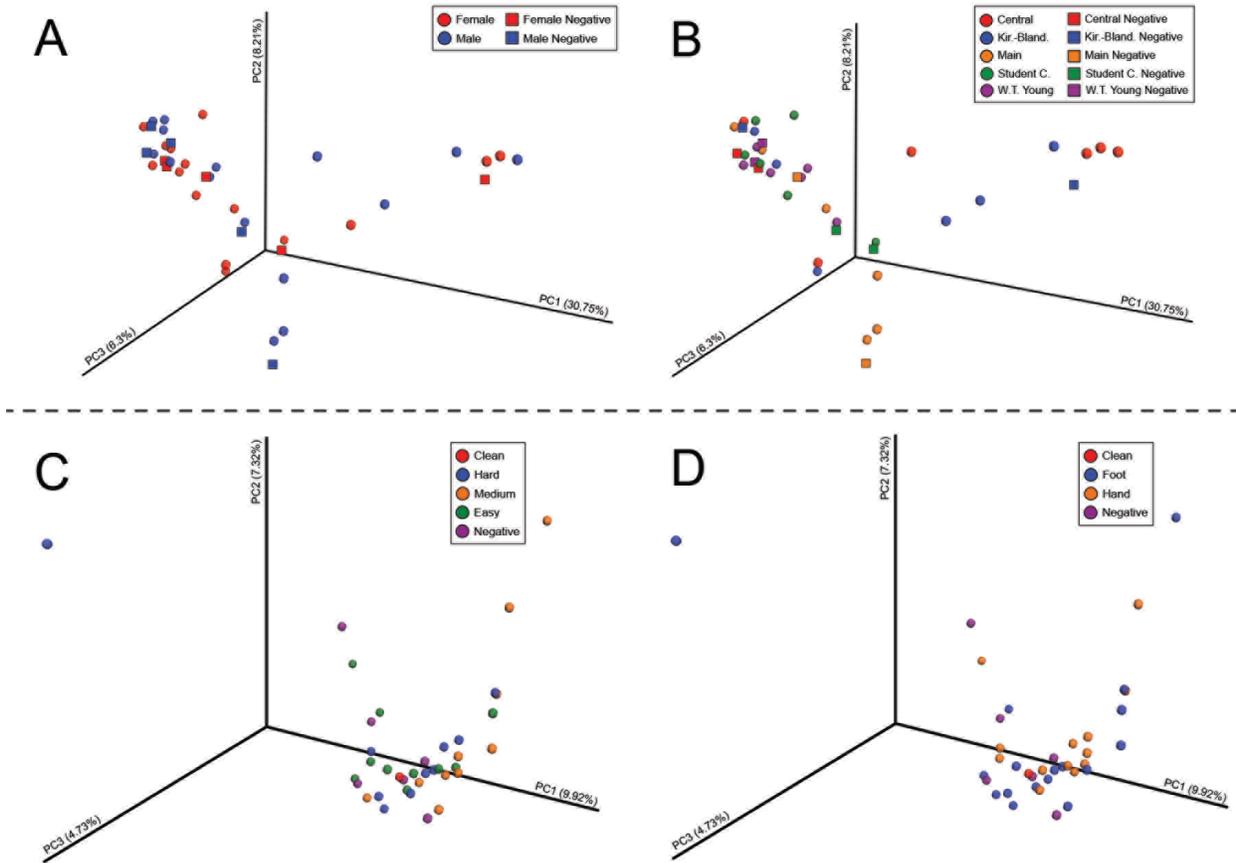
116 **References**

- 117 1. Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., *et al.* (2010).
118 QIIME allows analysis of high-throughput sequencing data. *Nature Methods*, 7, 335-336.
- 119 2. Hannon, G. (2010) FASTX-Toolkit. http://hannonlab.cshl.edu/fastx_toolkit/
- 120 3. Magoč, T., Salzberg, S.L. (2011) FLASH: fast length adjustment of short reads to
121 improve genome assemblies. *Bioinformatics*, 27, 2957-63.

122 **Supplementary Figures**

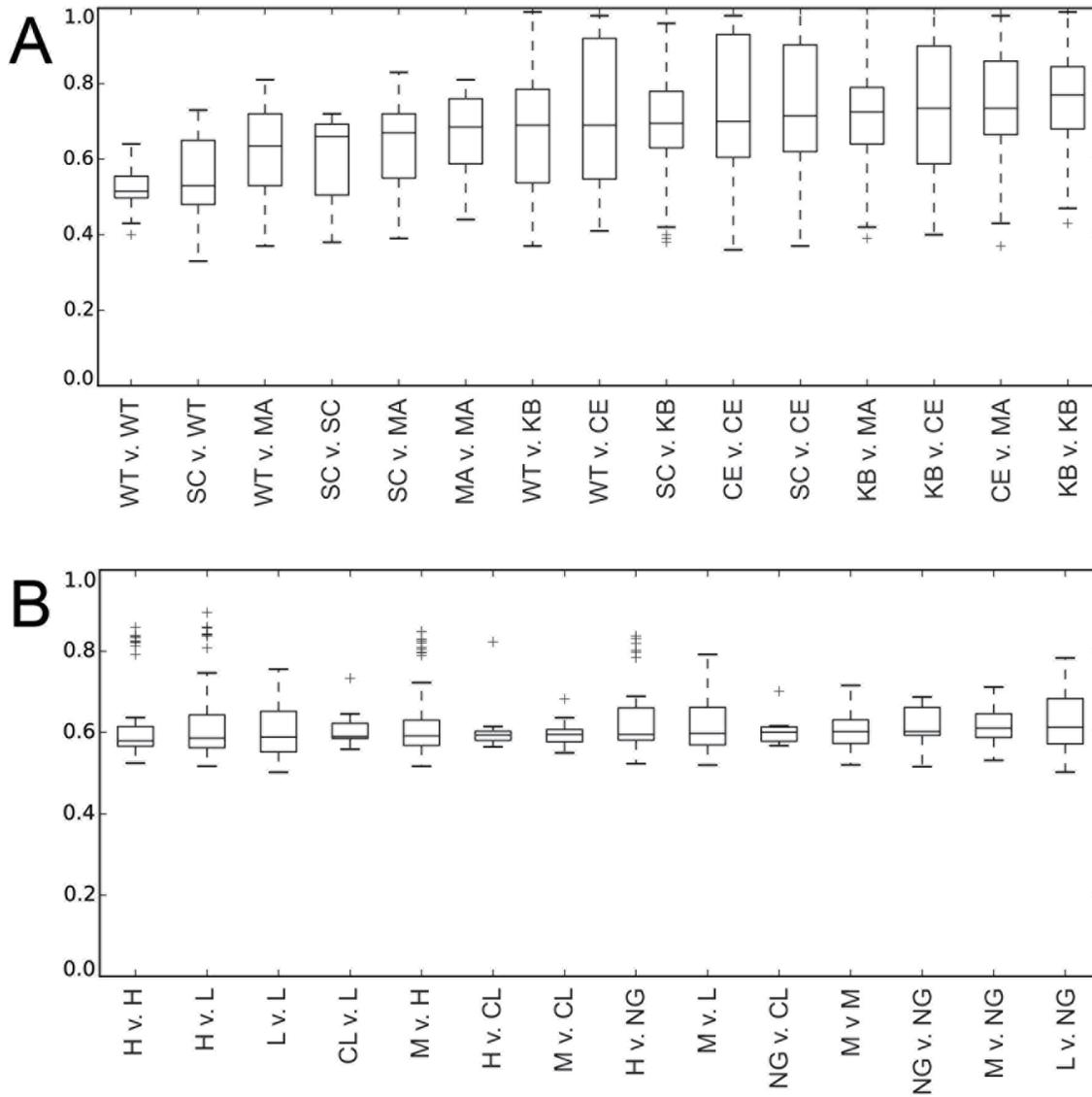
123

124 FIGURE S1. Three-dimensional principal coordinate analysis plots of prokaryotic beta diversity
125 for two years of the course-based undergraduate research experience. While plots appear static in
126 this representation, they are interactive for students and can be rotated, scaled, or modified in a
127 number of ways. (A) 2014: A comparison of prokaryotic diversity across genders in University
128 of Kentucky (UK) bathrooms. (B) 2014: UK bathroom prokaryotic diversity by location.
129 2015: Prokaryotic diversity of the climbing wall at the UK Johnson Center with samples colored
130 by route difficulty. (D) 2015: Prokaryotic diversity of the climbing wall at the UK Johnson
131 Center with samples colored by hold type.
132
133



134

135 FIGURE S2. Within and between treatment prokaryotic beta diversity comparisons for two years
 136 of the course-based undergraduate research experience. Hash marks represent individual samples.
 137 (A) 2014: location comparisons (WT = William T. Young Library, SC = Student Center, MA =
 138 Main Building, KB = Kirwan and Blanding Dormitories, CE = Central Dormitory). (B) Route
 139 height comparisons (H = high, M = middle, L = low, CL = clean hold, NG = negative).
 140



141

142 **Supplementary Tables**

143

144 TABLE S1. Weekly schedule of laboratory equipment we used for a section of the course with
145 ~20 students.

146

Item	Quantity	Notes
<i>Week 2 – Sampling and DNA extraction</i>		
Sterilizing tubes	1 per group	Benchtop heated tubes for sterilizing scissors
Scissors	5+ pairs	Metal with no coating (so they can be sterilized)
<i>Week 3 – PCR amplification</i>		
Thermocycler	1 per lab	
PCR primers	1 pair per library	Earth microbiome primers. See Bates et al. (2011) <i>ISME</i> and Gilbert et al. (2014) <i>BMC Biology</i> . http://press.igsb.anl.gov/earthmicrobiome/emp-standard-protocols/16s/
<i>Week 4 – Gel electrophoresis</i>		
Centrifuge	1 per group	
Gel electrophoresis rig and power supply	1 per group	
<i>Weeks 5-7 – Bioinformatics</i>		
Computer	1 per student	We used MacBook Pros
<i>Miscellaneous equipment</i>		
Containers for ice	1 per group	For reagents and extracted DNA
Micropipettes	1 set per 2 groups	10 μ L, 100 μ L, 1000 μ L
Kimwipes	1 box per group	
Tube racks	1 per group	

147

148 TABLE S2. Budget and item descriptions for resources needed for implementing this course-
 149 based undergraduate research experience in prokaryotic diversity. All costs are in USD.
 150

Item	Quantity	Cost	Catalog number or link
Illumina MiSeq	1 lane	\$1371 ^a	
Sterile culture swabs	100-pack	\$36.74/each	B4320115
Extract N Amp PCR Ready mix	1.2 mL bottle	\$116.50/each	E3004
Extraction buffer	24 mL bottle	\$86.10/each	E7526-24mL
Dilution buffer	12 mL bottle	\$33.60/each	D5688-12mL
EZ Vision DNA dye	1 kit	\$139.12/kit	N313-KIT
			Earth microbiome primers. See Bates et al. (2011) <i>ISME</i> and Gilbert et al. (2014) <i>BMC Biology</i> .
515f/806r primers	1 pair per library	\$43.50 per forward and reverse pair	http://press.igsb.anl.gov/earthmicrobiome/emp-standard-protocols/16s/ https://press.igsb.anl.gov/earthmicrobiome/emp-standard-protocols/primer-ordering-and-resuspension/
2 mL screw cap tubes	1 bag (500)		
PCR strip tubes and lids			
1 kb ladder			
1.75% agarose gels			
Various software packages	n/a	Free	All software used is freely available
MacBook Pro	30	\$1,499	Provided by the UK Department of Biology for use in a wide array of courses

151 ^a in-house rate at the University of Kentucky as of Fall 2015

152 TABLE S3. Prokaryotic alpha diversity across treatments and metrics for the 2014 iteration of
 153 the course-based undergraduate research experience. Treatments are sorted by Shannon diversity
 154 in descending order. Observed = total observed OTUs.
 155

Location	Surface	Gender	Shannon	Observed	chaol
Main	Sink	M	6.61	923.8	1740.2
Student Center	Negative	F	6.59	732.6	1555.1
Student Center	Sink	F	6.24	722.7	1132.8
Student Center	Negative	M	6.19	618	1059.1
W.T. Young	Faucet	M	6.1	872.3	1819.8
Main	Negative	M	6.01	622.9	1224.4
W.T. Young	Sink	M	6	633.4	987.6
Central	Door	F	5.97	653.9	1197.9
Kirwan-Blanding	Door	F	5.97	963.3	2842.6
Main	Door	M	5.95	599.8	1045.6
Central	Faucet	M	5.91	658.4	1082.1
Kirwan-Blanding	Sink	M	5.85	688.7	1239.4
Kirwan-Blanding	Sink	F	5.77	624.3	1204.9
Main	Door	F	5.71	642.5	975.9
Main	Sink	F	5.56	675.4	1032.6
W.T. Young	Door	M	5.56	633.2	1008.2
Central	Negative	F	5.46	590.8	900.0
Student Center	Door	F	5.45	691.9	1088.3
Student Center	Sink	M	5.4	618	953.2
Main	Faucet	F	5.39	683.6	1029.9
W.T. Young	Sink	F	5.37	625.8	1036.0
Central	Door	M	5.35	696.7	1129.3
Main	Negative	F	5.31	495.1	999.3
Central	Sink	M	5.3	577.6	1764.3
W.T. Young	Door	F	5.3	635.8	974.3
Kirwan-Blanding	Faucet	F	5.23	554.4	894.9
Kirwan-Blanding	Faucet	M	5.17	633.9	1501.7
Kirwan-Blanding	Negative	M	5.17	517.8	874.3
Student Center	Door	M	5.13	655.1	994.8
Kirwan-Blanding	Door	M	5.09	707.7	1099.7
W.T. Young	Faucet	F	5.07	608.4	972.4
W.T. Young	Negative	F	5.05	500	1052.3
Central	Negative	M	5.04	423.4	798.3
Kirwan-Blanding	Negative	F	5.04	692	1607.8
Student Center	Faucet	F	5.04	632.2	1043.9
Central	Sink	F	5.02	573.9	1830.0
Central	Faucet	F	5	547.7	1280.0
W.T. Young	Negative	M	4.98	435.1	1115.5
Student Center	Faucet	M	4.92	544.4	853.2
Main	Faucet	M	n/a	n/a	n/a

156

157 TABLE S4. Prokaryotic alpha diversity across treatments and metrics for the 2015 iteration of
 158 the course-based undergraduate research experience. Treatments are sorted by Shannon diversity
 159 in descending order. Observed = total observed OTUs.
 160

Difficulty	Height	Hold type	Shannon	Observed	chaol
Easy	High	Foot	6.41	950.7	1706.9
Easy	Low	Hand	5.61	756.9	1238.2
Easy	Middle	Hand	4.93	556.7	819.9
Easy	Middle	Foot	4.83	651.3	939.0
Easy	Middle	Foot	4.71	551.8	804.2
Medium	Low	Foot	4.67	620.7	933.5
Easy	Middle	Negative	4.67	556.9	798.9
Clean	Clean	Clean	4.65	411.3	593.3
Easy	Low	Negative	4.65	491.1	696.4
Easy	Low	Hand	4.50	537.6	759.8
Hard	Low	Foot	4.49	543.8	804.1
Hard	Low	Foot	4.48	162.5	415.9
Medium	Middle	Hand	4.45	466.2	678.9
Hard	Low	Hand	4.39	348.2	564.0
Hard	Middle	Hand	4.35	429.6	657.2
Easy	Low	Negative	4.31	365.7	584.5
Hard	Middle	Foot	4.30	311.3	477.9
Hard	High	Foot	4.26	427.4	631.0
Medium	Low	Hand	4.25	497.4	748.2
Hard	Low	Hand	4.24	467.5	642.8
Medium	Middle	Foot	4.17	268	404.7
Easy	High	Hand	4.11	453.1	697.0
Medium	High	Foot	4.10	378.3	568.7
Easy	High	Negative	4.06	340.9	599.0
Hard	Middle	Hand	4.06	475	734.9
Hard	Middle	Foot	4.03	340.8	499.2
Easy	Middle	Negative	3.98	310.8	595.4
Medium	High	Hand	3.96	386.2	548.5
Hard	High	Hand	3.92	401.3	586.6
Easy	Low	Foot	3.92	537.3	822.7
Easy	Low	Foot	3.73	302.2	425.3
Easy	High	Negative	3.68	304.1	672.2
Easy	High	Foot	3.65	357.6	562.9
Medium	High	Hand	3.60	341.7	579.9
Medium	Middle	Foot	2.75	360.7	578.7
Medium	Low	Foot	1.71	311.3	676.2
Medium	Middle	Hand	1.68	294	641.3

161

162 **Supplementary Materials – Instructor Guide and Example Scripts.zip:**

163

164 For a step-by-step run through of the analyses performed with more computational detail, see the
165 provided document ‘rRNA Instructor Analysis Guide.sh’. We recommend working in
166 TextWrangler (<http://www.barebones.com/products/textwrangler/>) or an equivalent text editor
167 (not Microsoft Word).

168

169 We have also included several template scripts.

170 **Supplementary Materials – Handouts and Protocols:**

171 The following are all of the handouts and protocols used in 2016^{1,2}.

172

173 ¹ There is no Week 1 handout because the first was a discussion of the project and experimental
174 design. Students were tasked with developing hypotheses in small groups for the class to vote on.

175

176 ² Though we did not require students to turn in answers to handout questions, we have provided
177 general thoughts for what we were steering them towards in red.

178 **Week 2 – Sampling and DNA Extraction.**

179

180 It is the third BIO 198 16S rRNA prokaryotic next-generation sequencing project and this year
181 we conducting a project to address hypotheses of prokaryotic species diversity in elevators in
182 buildings on the UK campus.

183

184 We will be able to address a number of hypotheses with this study. Some of these include:

185

186 1. The originally proposed hypotheses is that taller buildings will facilitate greater elevator
187 use, relative to shorter buildings, and that the “up” elevator buttons in the taller buildings
188 will have greater prokaryotic diversity and abundance compared to shorter buildings.

189

190 2. In addition, buildings on campus vary in their age, with some very new buildings having
191 less human traffic and use, relative to older buildings of the same height. We may expect
192 elevators in the older buildings to have greater prokaryotic diversity and abundance.

193

194 3. And there are a number of additional ways in which we may be able to compare samples
195 drawn from within and between elevators on campus. These will not be enumerated here,
196 but I bet everyone can identify the possibilities from the sampling design.

197

198 To test these hypotheses, we will perform environmental sampling from elevators in three
199 different buildings on campus: The Jacobs Science Building (JSB), Patterson Office Tower POT,
200 and the Thomas Hunt Morgan (THM) Building. The elevators across each of these buildings
201 vary in their controls, and so we won't sample exactly the same surfaces, but there will be
202 overlap.

203

204 Each lab section will be divided up into five groups. Individual group sampling breaks down like
205 this:

206

207 **Group 1**

208 JSB north elevator

209 Built in 2015

210

211 *Sample 1* Ground floor outside up button

212 *Sample 2* Inside 1st floor button

213 *Sample 3* Inside 2nd floor button

214 *Sample 4* Inside 3rd floor button

215 *Sample 5* Inside 1st floor push bar

216 *Sample 6* Inside 2nd floor push bar

217 *Sample 7* Inside 3rd floor push bar

218 *Sample 8* Alarm button

219 *Sample 9* Clean swab (negative)

Group 2

JSB service elevator

Built in 2015

Sample 1 Ground floor outside up button

Sample 2 Inside 1st floor button

Sample 3 Inside 2nd floor button

Sample 4 Inside 3rd floor button

Sample 5 Inside Basement button

Sample 6 Inside Penthouse button

Sample 7 Alarm button

Sample 8 Clean swab (negative)

220 **Group 3**
221 POT Elevator A
222 Built in 1968
223

224 **Sample 1** Ground floor outside up button
225 **Sample 2** Left panel, 1st floor button
226 **Sample 3** Left panel, 6th floor button
227 **Sample 4** Left panel, 12th floor button
228 **Sample 5** Left panel, 18th floor button
229 **Sample 6** Mezzanine button
230 **Sample 7** Alarm button
231 **Sample 8** Clean swab (negative)
232
233

234 **Group 5**
235 THM south elevator
236 Built in 1974
237

238 **Sample 1** Ground floor outside up button
239 **Sample 2** 1st floor button
240 **Sample 3** 2nd floor button
241 **Sample 4** 3rd floor button
242 **Sample 5** Alarm button
243 **Sample 6** Clean swab (negative)
244
245

246 This results in a total of **39 samples** that will be collected during each lab section. Each lab
247 section will collect the same set of samples, and perform DNA extractions and PCR for these
248 samples. This will result in three replicate samples of each surface environment, which will be
249 pooled prior to sequencing.

Group 4
POT Elevator D
Built in 1968

Sample 1 Ground floor outside up button
Sample 2 Left panel, 1st floor button
Sample 3 Left panel, 6th floor button
Sample 4 Left panel, 12th floor button
Sample 5 Left panel, 18th floor button
Sample 6 Mezzanine button
Sample 7 Alarm button
Sample 8 Clean swab (negative)

250 **Surface Sampling Protocol**

251

252 Throughout this work, handle materials and perform work while always aiming to minimize the
253 potential for contaminating samples. Wear clean gloves when handling materials and collecting
254 samples. Avoid touching non-target surfaces (and yourself) once you have put on your gloves.

255

256 ***Preparation***

257

258 1) All students will perform sampling in each section, with each student working within their
259 group to help sample the different buttons or bars associated with their elevator. Before
260 heading over to the different buildings and elevators, we will assemble the following materials
261 in the lab. We will wear gloves when collecting these items. Each group needs enough of
262 these for each student.

- 263 ● Clean gloves to be used when you have arrived at your sampling location. Place these in
264 a Ziploc bag.
- 265 ● A set of sterile swabs (contained in sterile tubes). Collect enough of these for the number
266 of samples to be collected from your elevator.
- 267 ● Sharpie for labeling tubes.

268

269 2) Descriptively label sterile swabs to account for all of the individual surface samples to be
270 collected. Include section number (i.e., 001, 002, or 003), elevator in abbreviated format
271 (e.g., JSB north elevator = JSBN), and sample number (1, 2, 3,).

272

273 **For example: 001, JSBN, 1**

274

275 3) Note that one tube in each group will be used as a negative control. This swab will be opened
276 when you get to the elevator, and then simply placed back in its tube without actually
277 swabbing anything.

278

279 4) Assemble all of the materials that will be used for each elevator sampling. Place into a bag.

280

281 5) Head to your elevator with Professor X or a TA!

282

283 ***Note: Remember that sample collection steps will involve both targeted samples and a negative***
284 ***control!***

285 ***Sampling***

286

- 287 1) When you have reached your elevator, carefully put on a clean set of gloves.
288 ● You should have a plastic bag with you containing multiple sets of clean gloves.
289 ● Identify one person in your group to serve as the button pusher (to open the door) and
290 door holder.

291

- 292 2) Next, collect one negative control sample. To do so, open one sterile swab and remove it
293 from its tube, then replace it back into the tube.

- 294 ● Carefully remove sampling swab from its container to maintain sterility.
295

296

- 297 3) Next, sample from the outside “up” button. Use the appropriate swab to gently brush the
298 button. You do not need to press super hard to pick up cells.

- 299 ● Carefully remove sampling swabs from their container to maintain their sterility.

- 300 ● Carefully place each individual swab in its proper sample tube and apply cap.

301

- 302 4) Hail the elevator and step inside to sample the buttons inside the elevator. Repeat Step 3 for a
303 the remaining button or push bar surfaces.

- 304 ● Each student can be responsible for 1-2 samples within their elevator.

305

- 306 5) Once you have finished sampling, head back to the lab. If you’re one of the POT groups,
hustle!

307 **DNA Extraction**

308

309 1) When everyone has returned to the lab, glove up!

310

311 2) Each group will obtain:

312

- 1-9 sterile 2.0 ml screw-cap tubes, depending on how many samples you collected

313

- Scissors for cutting swab tips.

314

- 2.0 ml tubes containing Extract-N-Amp Plant **Extraction** Solution. These tubes will be labeled specifically for each group and will contain enough solution to use in your extractions. Be sure to obtain the correct tube.

315

- 2.0 ml tubes tube containing Extract-N-Amp Plant **Dilution** Solution. These tubes will be labeled specifically for each group and will contain enough solution to use in your extractions. Be sure to obtain the correct tube.

316

317

318

319

- Sharpie for labeling tubes.

320

321

322 3) Label the 2.0 ml screw-cap tubes to correspond with each surface sample that your group is

323 working with.

324

- Use the same labeling style that you used for your sterile swabs.

325

- Also note on the tube that it contains DNA.

326

- It is often a good idea to label both the cap and the tube, in case one of them is accidentally removed (e.g., with EtOH).

327

328

329 3) Prepare your work surface by cleaning with a light application of 95% Ethanol (EtOH). Use

330 care when working to limit contamination of surface and supplies

331

332 4) Sterilize your scissors using a heat sterilizer.

333

334 5) Carefully remove one of your sterile swabs and place it into its corresponding 2.0 ml tube.

335 Then using your scissors, cut off most of the stick so that the end with the cotton tip fits in

336 the tube with the cap closed.

337

- Before you cut the stick, pull the swab up slightly so that it is not touching the bottom of the tube, which will allow it to drop a bit when you make your cut.

338

- Discard the other end of the stick in your waste container.

339

- Use care to make sure you are transferring swabs to their appropriately labeled tubes.

340

341

342 6) Repeat steps 4 and 5 for all of your sample swabs, including the negative controls.

343

- Pay attention to using good aseptic (sterile) technique. Always be thinking about how you can avoid cross-contaminating tools, reagents, samples, etc.

344

345

346 *When you have transferred all of your sample swabs to their 2.0 ml screw-cap tubes:*

347

348 7) Add 250 μ L of Extract-N-Amp Plant **Extraction** Solution to each tube and re-cap the tubes.

349

- Be sure to use a new tip for every tube! Avoid cross contamination.

350

- We suggest you do each tube individually, recapping before moving on to the next one.

351

- 352 8) When all groups have completed step 7, place them in a floating tube rack and heat in a 95°C
353 water bath for 10 minutes. Use a timer to monitor this time.
354
- 355 9) Centrifuge your tubes at 2500 g for 5 minutes. Ask for assistance if you have never loaded a
356 centrifuge before. Everything should be properly balanced.
357 • Use care when removing your tubes from the centrifuge so that you do not splash liquid
358 around inside the tube. We want everything to remain at the bottom so that you can
359 remove the lid without anything coming out of the tube.
360
- 361 10) Add 250 µL of Extract-N-Amp Plant **Dilution** Solution to each tube. When adding the
362 solution, keep your pipette tip in the liquid at the bottom of your tube and slowly pipette up
363 and down 2-3 times. Do this delicately so as to not splash liquid up out of the tube.
364 • Again, you should do this one tube at a time, closing each tube before you move on to the
365 next one.
366 • Use a new tip each time.
367
- 368 11) Collect all your tubes in your tube rack, label the rack with a piece of tape/sharpie with your
369 group initials, location, and date. Give your samples to an instructor to be stored at 4°C until
370 next week.

371 **Week 3 - PCR Amplification of the v4 region of prokaryotic 16S rRNA**

372

373 This week you will work in the same groups as in Week 4 when you collected and prepared your
374 DNA samples.

375

376 Our goal this week is to amplify copies of an ~250 bp segment of the 16S rRNA gene from the
377 various prokaryotic genomes in your DNA sample. This segment is referred to as the v4 region
378 because it is one of nine variable regions of the 16S rRNA gene in prokaryotes. We will use the
379 variation in sequence data from the v4 region to discriminate among the different prokaryotic
380 species that we have sampled in our bathroom surfaces.

381

382 From each elevator, you collected a range of surface samples. We will perform PCR on each
383 sample using uniquely barcoded (also referred to “indexed”) PCR primers. So, for example, PCR
384 of the north elevator in the Thomas Hunt Morgan Building will have the following reactions:
385

001 Group 5 samples	# of PCRs
Sample 1, ground floor up button	1
Sample 2, 1st floor button	1
Sample 3, 2nd floor button	1
Sample 4, 3rd floor button	1
Sample 5, Alarm button	1
Sample 6, Clean swab (negative)	1
PCR negative	1
Total	7

386

Table 1. Example of the PCR work to be performed in today’s lab, using the THM north elevator samples.

387

388

389 Let’s pay attention to a number of aspects of our sampling design here and our PCR work on
390 these samples.

391

392 First, each section will be performing a replicate of our overall elevator sampling design,
393 meaning that each surface sample will be PCRd three times. **Can you think of a good reason
394 for us to perform these PCR replicates across sections?**

395

396 **SUGGESTED ANSWER:** We performed PCR replicates across sections incase individual PCR
397 replicates failed.

398

399 Second, note that we are performing PCR on our clean swab samples (the ones you opened and
400 then closed at your particular elevator. **Why do we perform PCR for these samples?**

401

402 **SUGGESTED ANSWER:** We perform PCR for our negative swabs to amplify any possible
403 contaminants so that we can sequence them (and see what they are) later in the project.

404

405 Third, also note that we will perform a negative control for your PCR. This is separate from the
406 PCR we perform on your clean swab. **What is a PCR negative control? Why do we do this?**

407

408 **SUGGESTED ANSWER:** A PCR negative control is an identical reaction with water
409 substituted for DNA. This provides a control for in case reagents (and not our DNA sample) are
410 contaminated. Amplification in the negative control would indicate contamination.

411

412 Fourth, note that a different number of PCRs will be performed across groups within a section.
413 Make sure you pay attention to setting up PCR in a way that accounts for your group's samples.

414

415 Finally, remember that for each surface, we will want to generate PCR amplicons (amplified
416 stretches of DNA for our target gene region) for our different surface samples that can be
417 distinguished from each other in our Illumina sequence data. This will require us to use PCR
418 primers that have different barcode (indexed) sequences for each set of surface samples. We will
419 use the same index sequences for our PCR replicates of a particular sample across sections. You
420 can reference Table 2 (see below) to see which sets of PCR primers you need to use for each of
421 your surface samples.

422

423 **Table 2.** Information for 16S rRNA indexed primers. The Primer Mix name is the label given to
424 the tube of combined forward and reverse indexed primers. This tube will be used in preparing
425 the PCR cocktails. The individual forward and reverse primer names are listed for each primer
426 mixture and refer to the primer names provided in Flores et al. (2013).

427

Group/elevator/Treatment	Primer Mix Name	Index Primer Forward	Index Primer Reverse
G1 / JSB north / Sample 1	Primer Index 1	v4.SA501	v4.SA701
G1 / JSB north / Sample 2	Primer Index 2	v4.SA502	v4.SA701
G1 / JSB north / Sample 3	Primer Index 3	v4.SA503	v4.SA701
G1 / JSB north / Sample 4	Primer Index 4	v4.SA504	v4.SA701
G1 / JSB north / Sample 5	Primer Index 5	v4.SA505	v4.SA701
G1 / JSB north / Sample 6	Primer Index 6	v4.SA506	v4.SA701
G1 / JSB north / Sample 7	Primer Index 7	v4.SA507	v4.SA701
G1 / JSB north / Sample 8	Primer Index 8	v4.SA508	v4.SA701
G1 / JSB north / Sample 9	Primer Index 9	v4.SA501	v4.SA702
G1 / JSB north / PCR Neg.	Primer Index 10	v4.SA502	v4.SA702
G2 / JSB service / Sample 1	Primer Index 11	v4.SA503	v4.SA702

G2 / JSB service / Sample 2	Primer Index 12	v4.SA504	v4.SA702
G2 / JSB service / Sample 3	Primer Index 13	v4.SA505	v4.SA702
G2 / JSB service / Sample 4	Primer Index 14	v4.SA506	v4.SA702
G2 / JSB service / Sample 5	Primer Index 15	v4.SA507	v4.SA702
G2 / JSB service / Sample 6	Primer Index 16	v4.SA508	v4.SA702
G2 / JSB service / Sample 7	Primer Index 17	v4.SA501	v4.SA703
G2 / JSB service / Sample 8	Primer Index 18	v4.SA502	v4.SA703
G2 / JSB service / PCR Neg.	Primer Index 19	v4.SA503	v4.SA703
G3 / POT A / Sample 1	Primer Index 20	v4.SA504	v4.SA703
G3 / POT A / Sample 2	Primer Index 21	v4.SA505	v4.SA703
G3 / POT A / Sample 3	Primer Index 22	v4.SA506	v4.SA703
G3 / POT A / Sample 4	Primer Index 23	v4.SA507	v4.SA703
G3 / POT A / Sample 5	Primer Index 24	v4.SA508	v4.SA703
G3 / POT A / Sample 6	Primer Index 25	v4.SA501	v4.SA704
G3 / POT A / Sample 7	Primer Index 26	v4.SA502	v4.SA704
G3 / POT A / Sample 8	Primer Index 27	v4.SA503	v4.SA704
G3 / POT A / PCR Neg.	Primer Index 28	v4.SA504	v4.SA704
G4 / POT D / Sample 1	Primer Index 29	v4.SA505	v4.SA704
G4 / POT D / Sample 2	Primer Index 30	v4.SA506	v4.SA704
G4 / POT D / Sample 3	Primer Index 31	v4.SA507	v4.SA704
G4 / POT D / Sample 4	Primer Index 32	v4.SA508	v4.SA704
G4 / POT D / Sample 5	Primer Index 33	v4.SA501	v4.SA705
G4 / POT D / Sample 6	Primer Index 34	v4.SA502	v4.SA705
G4 / POT D / Sample 7	Primer Index 35	v4.SA503	v4.SA705
G4 / POT D / Sample 8	Primer Index 36	v4.SA504	v4.SA705
G4 / POT D / PCR Neg.	Primer Index 37	v4.SA505	v4.SA705
G5 / THM north / Sample 1	Primer Index 38	v4.SA506	v4.SA705
G5 / THM north / Sample 2	Primer Index 39	v4.SA507	v4.SA705
G5 / THM north / Sample 3	Primer Index 40	v4.SA508	v4.SA705
G5 / THM north / Sample 4	Primer Index 41	v4.SA501	v4.SA706
G5 / THM north / Sample 5	Primer Index 42	v4.SA502	v4.SA706
G5 / THM north / Sample 6	Primer Index 43	v4.SA503	v4.SA706
G5 / THM north / PCR Neg.	Primer Index 44	v4.SA504	v4.SA706

429 Now let's look at what goes into each PCR. They will each be conducted in a 20 μ L reaction that
430 includes:

431

- 432 ● 10 μ L of Extract-N-Amp Ready Mix (Sigma-Aldrich)
- 433 ● 1 μ L of a mixture containing both the forward and reverse primers, each at 5 μ M
- 434 concentration. See Table 2 for guidance on which primer mixtures to use.
- 435 ● 5 μ L of PCR-grade dH₂O
- 436 ● 4 μ L of your prepared DNA sample (from Week 4).

437

438 ***Let's get started!***

439

440 1) Prepare to take lab notes for all information you think is important while working through
441 this lab protocol and exercise.

442

443 2) Obtain one strip of twelve 200 μ L PCR tubes. You will use these tubes to cover:

444

445

446 ● Each group will have left-over tubes that will not be used. No worries.

447

448 3) Label your tubes using a fine Sharpie industrial marker. Make sure you provide enough
449 information on each tube so that they are clearly distinguished from each other, and that your
450 tubes can be identified from other groups.

451

452 ● These tubes are small and have limited space for writing. It will be best if you develop a
453 reduced labeling scheme that will allow you to distinguish different tubes (e.g., numbers)
454 and record the details of each tube in your lab notes. Make sure to include information for
455 section number and group!

455

456 ● It is also best if you put a label on both the side of the tube AND the top, just in case one
457 of them rubs off.

457

458 ● Because your tubes will be connected to each other, you don't need to fully label every
459 tube with all of this information.

458

459 ● Finally, place an X on the empty tubes at the end.

460

461 4) Collect ice using your ice bucket, then obtain the following reagents:

462

463 ● A microcentrifuge tube containing Extract-N-Amp PCR Ready Mix. Note that your
464 group will obtain a specific tube of this Ready Mix with enough volume to perform your
465 PCRs. Be sure to obtain the correct tube.

465

466 ● 0.5 mL microcentrifuge tube containing the primer mixes needed for the different surface
467 samples that your group is working with. Be sure to retrieve the proper primer mix for
468 each surface sample!

468

469 ● A 1.5 mL microcentrifuge tube containing 60 μ L of PCR-grade dH₂O.

469

470 Keep all of your reagents on ice. Some of these reagents are temperature sensitive.

471

472 Furthermore, the Taq polymerase within your Extract-N-Amp PCR Ready Mix should be

- 473 kept inactive to prevent unwanted polymerase activity, which is done by keeping it nice and
474 cold.
475
- 476 5) Obtain your DNA samples from last week. Using a centrifuge, spin them at high speed for a
477 short amount of time (a few seconds) to make sure all of the liquid is at the bottom of the
478 tube.
- 479 ● This is always a good idea before you open any of your tubes, including reagents. It is
480 particularly important for your DNA tubes.
 - 481 ● Keep your DNA samples on ice.

482 **Preparing PCR: An overview.**

483

484 To set up your PCRs, you will not pipette all reagents for each reaction individually. **Why not?**

485

486 **SUGGESTED ANSWER:** This answer is two-fold. First, we will be pipetting some very small
487 volumes. It's easy to make an error with small volumes so it's best to avoid pipetting them when
488 possible. Second, by combining reagents into a master cocktail or mix, we save significant time by
489 reducing the number of individual pipettes needed.

490

491 Instead, a more appropriate way to prepare your PCR mixture is to create a cocktail of reagents
492 that can be dispensed as aliquots to your individual PCR tubes. This cocktail will be prepared for
493 all PCRs using reagents that are exactly the same, and omitting the things that make them
494 different. In this case, each PCR will differ in the primers being used, and the surface-sample
495 DNA. We will be adding the PCR primers and DNA separately to each tube. **How many**
496 **different cocktails will you need for today?**

497

498 **SUGGESTED ANSWER:** If the PCR primers and DNA will be added separately, then
499 everything else will be the same. Therefore, only one cocktail needs to be made for this
500 approach.

501

502 Example: Let's say your instructor is performing PCR to amplify the same gene region from a
503 set of **five** DNA samples and he isn't concerned with separately barcoding these samples. He
504 would prepare a cocktail that will have enough reagents for:

505

- a PCR for each of his five DNAs.
- a PCR to serve as a negative control (to test for contamination of PCR reagents).
- one extra PCR that allows for pipetting error. **Why have this extra?**

506

507

508

509 **SUGGESTED ANSWER:** We add an extra sample (or two) to make sure that we have enough
510 for every reaction. Some volume is lost to normal pipetting (e.g., some liquid sticks to the pipette
511 tip) so it's important to have a buffer when preparing the cocktail.

512

513 Therefore, he would prepare a cocktail for **seven** PCRs using the following reagents and
514 volumes.

515

516 10 μL x 7 = 70 μL Extract-N-Amp Ready Mix

517 1 μL x 7 = 7 μL Primer Mix

518 5 μL x 7 = 35 μL PCR-grade dH_2O

519 = **112 μL total volume of cocktail**

520

521 Using this 112 μL cocktail, he would dispense 16 μL to each of his six PCR tubes (five DNA
522 tubes and one negative control). There will probably then be a little cocktail left over in the tube.
523 No big deal.

524
525 He would then add 4 μL of the appropriate DNA to each tube to complete the preparation of a 20
526 μL PCR volume.

527
528 For the PCR negative control, the addition of 4 μL of PCR-grade dH_2O will be used to bring the
529 PCR volume up to 20 μL .

530
531 **What if Professor X needs to do a similar PCR with PCR primers that separately barcoded**
532 **each sample?**

533
534 **SUGGESTED ANSWER:** In this case, Dr. X would make a single cocktail that would not
535 include PCR primers in it. The primers would be added separately to each reaction.

536 **Week 4 - Gel electrophoresis of PCR amplicons**

537

538 This week you will continue to work in the same groups that you were in to perform surface
539 sampling and PCR.

540

541 We'll be using gel electrophoresis to determine if your PCR reactions successfully amplified the
542 expected v4 segment of the 16S rRNA gene.

- 543 ● In addition, this gel work will be used to determine if there is any potential contamination
544 in your surface sampling (based on PCR of your swab negative controls).
- 545 ● We'll also be able to check if there is any contamination in your PCR, either in the
546 reagents, or in the preparation of the reactions (based on your PCR negative controls).

547

548 *A little background on gel electrophoresis....*

549

550 Electrophoresis of DNA is typically performed using an agarose gel. Agarose is a polysaccharide
551 extracted from seaweed. In a powdered form, agarose will dissolve into a liquid buffer using heat
552 (typically using a microwave). The melted agarose-buffer mixture can be poured into a mold to
553 form a gel. As it cools, the gel will solidify; kinda like jello.

- 554 ● The liquid buffer is usually a mixture of dH₂O, Tris, Acetic Acid or Boric Acid, and
555 Ethylenediaminetetraacetic acid (EDTA). These buffers are usually referred to by the
556 acronyms TAE and TBE, depending on their composition.
- 557 ● A comb is used to create wells in the gel.
- 558 ● DNA (PCR product in your case) will be pipette loaded into these wells.

559

560 The solidified gel is submerged in the same buffer (TAE or TBE) that was used to make the gel.
561 This is done in a gel rig that contains a positive electrode (cathode) at one end and a negative
562 electrode (anode) at the other end. DNA has an electric charge to it. When an electric current is
563 applied to the gel rig (filled with TAE or TBE), the DNA will migrate towards the positive
564 electrode.

- 565 ● The TAE/TBE buffer has ions in it that facilitate the transmission of electrons from
566 negative to positive end. Don't ever put dH₂O in the gel rig. Without ions, no electric
567 charge will move through the gel.

568

569 The percentage of agarose in your gel will determine how fast DNA molecules will migrate
570 through the gel matrix. In general, shorter DNA molecules will migrate faster than larger DNA
571 fragments. If you have small DNA fragments, you may want the agarose percentage to be a bit
572 higher so that they don't run too fast.

- 573 ● Our gels will be made at 1.5% agarose to account for the expected 350 bp size of our
574 PCR fragment.

575

576 To load DNA/PCR product into a well, it is first mixed with a loading buffer. This usually has
577 some kind of dense molecule (e.g., sucrose) that helps everything fall to the bottom of the well.

578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620

A dye is also included that will give the loading buffer some color, and which has a molecule that will migrate with the DNA. This will allow you to track the movement of your DNA on the gel.

In addition, a DNA ladder is also loaded into a well on the gel, typically to one side or the other of the samples in the run. This contains a series of DNA fragments of known size. It allows you to compare your PCR results to known DNA sizes to assess whether your PCR was successful.

Finally, a molecule is added to your DNA/PCR material prior to loading that interacts with your DNA and allows it to be visualized under UV light.

- Historically, this has been Ethidium Bromide (or commonly written as ‘EtBr’), which is also a known mutagen and carcinogen.
- In today’s lab, we’ll be using a non-toxic DNA dye provided in the EZ-Vision loading buffer.

Let’s get started!

- 1) As usual, take notes throughout your lab activity today.
- 2) Obtain ice for your ice tray/bucket and the following items:
 - Your PCR tubes from last week. Remember, you should have three sets of strip tubes, one for each of your surface samples.
 - A tube containing 1kb DNA size ladder.
 - A tube of EZ-Vision Three DNA loading dye.
 - A tube of gel-loading buffer.
 - Two new strips of PCR tubes.
- 3) Centrifuge your PCR tubes to make sure everything is at the bottom of the tube. Remember, you just need to provide a quick spin. Make sure to balance your tubes!

First, we’re going to make sure everyone gets a little practice with pipetting into the wells of agarose gels.

- 4) Obtain an agarose gel and place it into your gel rig.
 - Note that this gel has two rows of 12 wells.
- 5) Add 1X TAE buffer to the gel rig until it just covers the top of your gel.
 - This doesn’t have to be exact, just make sure your gel is completely submerged.
 - If any air bubbles are present in the wells, gently push them out with a yellow pipette tip. Be careful not to puncture the gel.

- 621 6) In one of your empty strip-tubes, add 10 μL of gel loading buffer to each of the 12 tubes.
622 • Be sure to use the gel loading buffer, and NOT the EZ-Vision DNA dye.
623
- 624 7) Next, pipette 9 μL of the loading buffer from the first tube and practice loading it into the
625 first well on bottom of the gel.
626 • This requires precision and a steady hand. And it definitely requires practice.
627 • An instructor can help you with technique.
628 • Be very careful not to puncture the gel with your tip. This is especially true for the
629 bottom of the well. A punctured well bottom will result in most of your loading material
630 leaking out.
631 • Avoid having air in the end of your pipette tip before you submerge it into the buffer. If
632 you do have a small amount of air, slightly press on your pipette plunger to expel the air.
633
- 634 8) Each student will practice loading two to three wells in the gel.
635 • After you perform the practice loadings we'll move on to the next step. We can leave the
636 loading buffer in the wells for now.
637

638 *Now we'll move on to loading and running your PCR product on a gel.*
639

- 640 9) Obtain a clean strip-tube and number the tubes according to your PCR numbers moving left
641 to right. Label one additional tube with an "L" for ladder.
642 • We will use these tubes to mix the EZ-Vision dye with a sample of our PCR.
643
- 644 10) Pipette 2 μL of EZ-Vision Three DNA dye into each of the newly labeled tubes.
645
- 646 11) Next, using the PCRs you performed last week, pipette 8 μL of each into the tubes
647 containing the EZ-Vision Three DNA dye. There now should be a total of 10 μL (8 of PCR
648 product and 2 of dye) in each of the labeled tubes.
649 • Do this pipetting moving left to right, and match each PCR tube to the corresponding
650 tubes containing your EZ-Vision Three DNA dye.
651
- 652 14) Pipette 5 μL of the DNA ladder into the corresponding tube containing EZ-Vision Three
653 DNA dye.
654
- 655 15) Now, pipette 9.5 μL of each of your PCR-dye mixtures and load them into the wells of the
656 top row of your agarose gel. It's OK to leave a tiny amount left in the tube.
657 • Use the top row! Not the bottom row that you used for practice loadings.
658 • Again, avoid having air in the end of your pipette tip before placing it into the TAE
659 buffer when loading.
660
- 661 16) Load 6.5 μL of your ladder in the next available well of your gel.
662
- 663 17) Place the lid on your gel rig. Check to make sure the connectors are properly secured to the
664 power source. And check to make sure the bottom of your agarose gel is closest to the
665 positive node of the gel rig.

666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693

- 18) Turn on your power source and set the voltage to 75.
 - 19) ***Keep an eye on the tracking dye on your gel.*** The EZ-Vision Three DNA dye contains three tracking dyes that migrate at different rates.
 - The light blue dye (Xylene cyanol) runs at about 4000 bp. Big.
 - The darker blue dye runs at about 400 bp. Just a bit bigger than our expected fragment size.
 - The reddish band runs at about 10 bp. Smaller than one of our primers.
 - 19) When your gel has finished running, turn off the power supply, remove the lid to your gel rig, and bring your gel to the gel visualization equipment.
 - This consists of a UV light box, which will be used to see DNA within the gel.
 - 20) With the help of an instructor, check your gels using the UV system to assess the success of your PCR and the presence/absence of amplification in your negative controls.
 - If your PCR was successful, you should see a band, or a series of bands, that are ~350 bp in size.
 - Hopefully, you will not detect amplification in your negative controls. But, it might be there. We can discuss more.
 - 21) After lab, an instructor will take your gels and photograph them under UV and produce a digital photo of your gel
- We will summarize the results across the entire class and (fingers crossed) prepare to send our PCR product to the Advanced Genetic Technology Center (AGTC) on campus for sequencing on the Illumina MiSeq.***

694 **Week 5 - Intro to sequence data and command line analysis (Bioinformatics, Part 1)**
695

- 696 1) We'll start with an introduction to the data. Look for the raw sequencing data in the
697 **16SrRNA_data** folder on the **Desktop**. This is the entire set of sequence data generated from
698 our sequencing run. Identify the set of sequences corresponding to your particular group.
699
- 700 2) You will see that each set of sequences is coded as a **.gz** file. This is a compressed format.
701 We will need to unpack the sequence file from its compressed format. You could just double
702 click on the file and have your computer do this in an automated way, but we're going to
703 learn a better way to do this using your terminal and command lines.
704
- 705 3) Now, go to the dock on your computer and locate the **Terminal** app. Open it up. It will open
706 a window that has a command prompt. This typically ends in a **\$** sign and has a flashing
707 cursor. The command prompt is a powerful thing, and gives you the ability to do way more
708 with a computer than what you normally do with a point and click mouse and a GUI
709 (graphical user interface) interface.
710
- 711 4) We will start with a simple step to check our software. In your terminal prompt type **gcc**
712 **--version** and hit **enter**. This checks to see if we have command line tools installed. This
713 should return approximately four lines of information. If it doesn't, let us know.
714
- 715 5) Next, we will start learning how to navigate the computer via Terminal. To start, type the
716 command **ls** and press **enter**. This is a list command and will show you all of the files in your
717 current directory. When you open your terminal, the default starting point is your home
718 directory, so you should be seeing the files in your home directory, including the **Desktop**
719 folder (also called a directory).
720
- 721 6) We want to move into the **Desktop** directory. To do this type the command **cd Desktop**. This
722 is a change directory command (hence, **cd**).
723
- 724 7) Use the **ls** command again, and you should see everything on the **Desktop**.
725
- 726 8) Change directories again to move into the **16SrRNA_data** directory (**cd 16SrRNA_data**).
727 Then use another **ls** command to list everything in the directory.
728
- 729 9) You should see a list of compressed data files, each in a **.gz** format. Identify the files
730 associated with your group. You'll see that the each file begins with a group number
731 designation along with a abbreviation of the building and sample number (e.g.,
732 G1_JSBN_S1).
733
- 734 10) You can also identify an **Index#** in each file, which refers to a index combination that were
735 in the PCR primers for a particular sample.

	Index	Sample description
736		
737	Index1	JSB north elevator, sample 1
738	Index2	JSB north elevator, sample 2
739	Index3	JSB north elevator, sample 3
740	Index4	JSB north elevator, sample 4
741	Index5	JSB north elevator, sample 5
742	Index6	JSB north elevator, sample 6
743	Index7	JSB north elevator, sample 7
744	Index8	JSB north elevator, sample 8
745	Index9	JSB north elevator, sample 9
746	Index10	JSB north elevator, PCR negative
747	Index11	JSB service elevator, sample 1
748	Index12	JSB service elevator, sample 2
749	Index13	JSB service elevator, sample 3
750	Index14	JSB service elevator, sample 4
751	Index15	JSB service elevator, sample 5
752	Index16	JSB service elevator, sample 6
753	Index17	JSB service elevator, sample 7
754	Index18	JSB service elevator, sample 8
755	Index19	JSB service elevator, PCR negative
756	Index20	POT elevator A, sample 1
757	Index21	POT elevator A, sample 2
758	Index22	POT elevator A, sample 3
759	Index23	POT elevator A, sample 4
760	Index24	POT elevator A, sample 5
761	Index25	POT elevator A, sample 6
762	Index26	POT elevator A, sample 7
763	Index27	POT elevator A, sample 8
764	Index28	POT elevator A, PCR negative
765	Index29	POT elevator D, sample 1
766	Index30	POT elevator D, sample 2
767	Index31	POT elevator D, sample 3
768	Index32	POT elevator D, sample 4
769	Index33	POT elevator D, sample 5
770	Index34	POT elevator D, sample 6
771	Index35	POT elevator D, sample 7
772	Index36	POT elevator D, sample 8
773	Index37	POT elevator D, PCR negative
774	Index38	THM north elevator D, sample 1
775	Index39	THM north elevator D, sample 2
776	Index40	THM north elevator D, sample 3
777	Index41	THM north elevator D, sample 4
778	Index42	THM north elevator D, sample 5
779	Index43	THM north elevator D, sample 6
780	Index44	THM north elevator D, PCR negative
781	Index45	JSB service elevator, sample 1, PCR redo
782	Index46	THM north elevator D, sample 3, PCR redo
783	Index47	THM north elevator D, sample 5, PCR redo
784	Index48	Mystery sample

- 785 11) Now we will unpack some of the **.gz** files. Each group should be able to identify the set of
786 samples that they worked on. From your group samples, choose to work on a particular
787 sample. Note that some samples did not yield much sequence data. Choose to work with one
788 of the larger files, and not one with a file size <10 KB. To see the file sizes in your terminal,
789 use the command **ls -lh** to provide a detailed list of what is inside your current directory.
790
- 791 12) Note that each surface sample has two different **.gz** files, one with an **R1** in the file name, the
792 other with an **R2** in the file name. These correspond to our paired-end sequences. Read 1 is
793 sequence that starts from one end of our PCR fragment and Read 2 is sequence that starts
794 from the other end. Each 250 bp sequence was generated on an Illumina MiSeq sequencer.
795
- 796 13) When you have identified your particular sequence set to work with, give the following
797 command in your terminal: **gunzip yoursamplename_R1.fastq.gz -k** This will unpack the
798 Read 1 sequences from your particular sequences and produce a text file with a similar
799 filename (without the **.gz** extension). The **-k** flag tells the program to keep a copy of the
800 original gzip file.
801
- 802 14) Use the **ls** command to see the unpacked sequence. It should be similar to the **.gz** file, except
803 it won't have that extension. Rather, its extension will now just be **.fastq**.
804
- 805 15) Now do the same for your Read 2 sequences: **gunzip yoursamplename_R2.fastq.gz -k**
806
- 807 16) Using your Finder (not your terminal), open the **16SrRNA_data** folder on your Desktop.
808 From your dock (at the bottom of your screen), open TextWrangler. Now, drag your recently
809 created R1 fastq file into the TextWrangler window. This will open your file.
810
- 811 17) Let's take some time to look at one of our raw data files. It may look like a bunch of cryptic
812 information, but we can make sense of it. Sequences are grouped in blocks of 4 lines. The
813 first line of each sequence looks something like this:
814
- 815 • **@** - Each sequence identifier line starts with **@**.
 - 816 • **InstrumentID** - unique identifier of the sequencer (M00329)
 - 817 • **RunNumber** - Run number on instrument (129).
 - 818 • **Flowcell_ID** - ID of flowcell (000000000-ABTD7).
 - 819 • **Lane Number** - positive integer, currently 1-8 (1)
 - 820 • **Tile Number** - positive integer (1)
 - 821 • **X** - x coordinate of the spot. Integer which can be negative (14930)
 - 822 • **Y** - y coordinate of the spot. Integer which can be negative (1700)
 - 823 • **Read Number** - 1 for single reads; 1 or 2 for paired ends (1)
 - 824 • **Whether it is filtered** - NB: Y if the read is filtered out, not in the fastq file, N otherwise
825 (N)
 - 826 • **Control number** - 0 when none of the control bits are on, otherwise it is an even number
827 (0)
 - 828 • **Index Combination** - The barcode sequences from that particular sequence.
- 829
- 830 18) The second line contains the actual sequence data generated for that read.

- 831
 832 19) The third line contains a + sign, indicating an upcoming quality score line.
 833
 834 20) The fourth line contains the quality score for each corresponding base call. You'll notice that
 835 these are not actually numbers, but are instead letters and characters. This is a space issue
 836 here, as it's easier for the Sequencer to note and store a single ASCII character, than it is a
 837 larger number. Here's what the translation might look like:
 838

char	value	Q	Error Prob.
!	33	0	1.000000e+00
"	34	1	7.943282e-01
#	35	2	6.309573e-01
\$	36	3	5.011872e-01
%	37	4	3.981072e-01
&	38	5	3.162278e-01
'	39	6	2.511886e-01
(40	7	1.995262e-01
)	41	8	1.584893e-01
*	42	9	1.258925e-01
+	43	10	1.000000e-01
,	44	11	7.943282e-02
-	45	12	6.309573e-02
.	46	13	5.011872e-02
/	47	14	3.981072e-02
0	48	15	3.162278e-02
1	49	16	2.511886e-02
2	50	17	1.995262e-02
3	51	18	1.584893e-02
4	52	19	1.258925e-02
5	53	20	1.000000e-02
6	54	21	7.943282e-03
7	55	22	6.309573e-03
8	56	23	5.011872e-03
9	57	24	3.981072e-03
:	58	25	3.162278e-03
;	59	26	2.511886e-03
<	60	27	1.995262e-03
=	61	28	1.584893e-03
>	62	29	1.258925e-03
?	63	30	1.000000e-03
@	64	31	7.943282e-04
A	65	32	6.309573e-04
B	66	33	5.011872e-04
C	67	34	3.981072e-04
D	68	35	3.162278e-04

877 | E | 69 | 36 | 2.511886e-04 |
878 | F | 70 | 37 | 1.995262e-04 |
879 | G | 71 | 38 | 1.584893e-04 |
880 | H | 72 | 39 | 1.258925e-04 |
881 | I | 73 | 40 | 1.000000e-04 |

882

883 21) Next let's consider our paired-end sequence data. Read 1 and Read 2 are both about 250 bp
884 in length, and our 16S rRNA PCR fragment is only expected to be somewhere in the 250-300
885 bp size range. This means that our R1 and R2 sequence reads probably substantially overlap.
886 We can use some software to put them together and create a consensus sequence. We'll do
887 this with the program FLASH (Fast Length Adjustment of SHort Reads).

888

889 22) Do a Google search for **FLASH read merger**. Your top hit should be the program website
890 run through Johns Hopkins University. Go to the **FLASH** website and click the **sourceforge**
891 link near the bottom of the page. Next click on the **Download FLASH-1.2.11.tar.gz** link to
892 download the current version of FLASH.

893

894 23) Locate the downloaded **FLASH-1.2.11.tar.gz** file in your **Downloads** folder and move it to
895 your **Desktop**.

896

897 24) Open your terminal window and move to the Desktop directory. If you are still in the
898 **16SrRNA_data** directory use the command **cd ..** to move back one level to the **Desktop**. If
899 you are unsure of your current directory, you can always use the command **pwd** (print
900 working directory) to identify your current location in directory space.

901

902 25) Use an **ls** command to verify the **FLASH-1.2.11.tar.gz** file is in your Desktop directory.

903

904 26) Now we'll unpack the **FLASH-1.2.11.tar.gz** file. This one is a little more complex than
905 the .gz files we unpacked earlier. It's referred to as a Tar file, and we use the following
906 command: **tar -zxvf FLASH-1.2.11.tar.gz**

907

908 27) This creates a new **FLASH-1.2.11** directory in your Desktop. You can migrate into this
909 directory using the **cd FLASH-1.2.11** command.

910

911 28) Use the **ls** command to take a look inside the directory you just unpacked. You'll see a
912 number of different files. Many of these have a **.c** or **.h** extension. These are the actual source
913 code for FLASH, but in their current format they don't function as an actual executable
914 program (we often refer to software as an executable). Other files give us some information
915 about using FLASH.

916

917 29) Use the command **cat README** to look inside the README file. The cat command is a
918 quick way to look at the text in a file. The README file gives us a bunch of basic
919 information about FLASH. Notice that it tells us to use the command **make** to compile
920 FLASH. Many programs distributed as source code contain a make file that tells the
921 computer how to compile the software.

922

923 30) Use the command **make** to compile FLASH.
924

925 31) In your terminal, use the **ls -lh** command to see the file list in the FLASH directory, along
926 with the dates when files were created (among other things). You should see a file simply
927 labeled **flash**. This is the compiled executable. Use the command **./flash --help** to see a list of
928 options to use with FLASH. We'll take a look at a few of the options here.
929

- 930 • First, take a look at the **-m** flag. When we use different settings in an executable, we refer
931 to these as flags, which are typically preceded by a dash. The **-m** flag sets the minimum
932 overlap for our R1 and R2 reads. The default is 10 bp, meaning that if the 2 reads do not
933 overlap by at least 10 bp, they will not be merged. We'll leave this at the default setting.
934
- 935 • Next, look at the **-M** flag. This set the maximum overlap between reads. Any reads that
936 have greater overlap than this setting will not be merged. The default is 65 bp, which is
937 too low for us given that we have 250bp reads covering an ~250 bp fragment. We'll use
938 **-M 250**.
939
- 940 • And now look at the **-o** flag. This sets the prefix for your output files.
941
- 942 • There are a number of other flags to consider, but we'll be good using the default settings
943 for the rest of the parameters of the analysis.
944

945 32) Before we run FLASH on your data, let's first move your specific data files (the ones you
946 unpacked that ends in **.fastq**) to your FLASH-1.2.11 directory. To do this, you need to be in
947 your **16SrRNA_data** folder. Navigate there by typing **cd**
948 **/Users/username/Desktop/16SrRNA_data** and hit enter. Note that this uses a path to guide
949 you to the **16SrRNA_data** folder.
950

951 33) Now let's copy the files you will use over to the FLASH directory. This is a good
952 demonstration on how to use the copy (**cp**) command. It will be something like this...
953

```
954     cp yourfilename_R1.fastq yourfilename_R2.fastq
955 /Users/username/Desktop/FLASH.1.2.11
```

956

957 34) Next, let's run FLASH on your data. First, we need to move back to our FLASH-1.2.11
958 directory. To get there, type **cd /Users/username/Desktop/FLASH-1.2.11**. Now, use the
959 following command. Note: This command is all one line separated by spaces and
960 'yourfilename' is only being used as a template!
961

```
962     ./flash -M 250 yourfilename_R1.fastq yourfilename_R2.fastq -o yourfilename
```

- 963 35) Assuming that flash worked for you (you should see screen output indicating it worked), use
964 an **ls** command to look at the new files that have been created.
965
- 966 36) Use **cat Your_prefix.hist** to see a summary of the merged reads by length.
967
- 968 37) Next, look at the plot for these values using **cat Your_prefix.histogram**.
969
- 970 38) Now that we have merged reads for your data, let's take a little time to explore how we might
971 learn about their prokaryotic source. In TextWrangler, open your file that ends in
972 ***.extendedFrag.fastq**.
973
- 974 39) Open a web browser and go to the following website: **http://www.ncbi.nlm.nih.gov/** This is
975 the NCBI homepage, a clearinghouse for genetic information and resources.
976
- 977 40) Now click on the **DNA & RNA** link. This will list a variety of genetic databases.
978
- 979 41) Next click on the **Genbank** link. This is a database containing all public DNA sequence data.
980 Most journals require you to publish your DNA sequence data prior to publication. It usually
981 is submitted here.
982
- 983 42) Now click on the **Search Genbank** link. We're going to explore some of our sequence data.
984
- 985 43) Click on the **BLAST** link. Then, click on the **nucleotide blast** link.
986
- 987 44) OK, here we are, ready to go. Let's BLAST!

988 **Week 6 – Bioinformatics, Part 2: QIIME.**

989

990 ***Installing MacQIIME***

991 1) Locate the file labeled **MacQIIME_1.9.1-20150604_OS10.7.tgz** on the desktop.

992

993 2) Open your terminal and change directory (**cd**) to the Desktop. Use an **ls** command to make
994 sure the MacQIIME tar file is available on the desktop.

995

996 3) Use the command **tar -xvf MacQIIME_1.9.1-20150604_OS10.7.tgz** to unpack the tar file.
997 This may take a few minutes as it unpacks all of the files. When your command prompt
998 returns, use an **ls** command to make sure there is a new **MacQIIME_1.9.1-**
999 **20150604_OS10.7** folder on your desktop.

1000

1001 4) In your terminal, change directory into the **MacQIIME_1.9.1-20150604_OS10.7**
1002 directory. Use an **ls** command to make sure you have an **install.s** file in this folder. This is
1003 similar to a make file and will perform the compiling of MacQIIME.

1004

1005 5) Next, use the command **./install.s** to install a copy of MacQiime on your computer. This
1006 installs MacQIIME in the root directory of the computer, meaning we can use it from
1007 within any folder that the terminal is in.

1008

1009 6) Use the command **macqiime** to source it, making it functional and ready to be used. If it
1010 installed properly, a number of lines should output to the screen and your user prompt will
1011 now start with the term “MacQIIME.”

1012

1013 7) As a further test, use the command **print_qiime_config.py -t** to see if a bunch of
1014 information is output that tells you about the various programs and configurations that
1015 make up QIIME.

1016 **Converting fastq to fasta**

1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055

Our fastq files contain lots of extra stuff in them that we no longer need, so we are going to convert them to a simpler file format termed fasta. We will use QIIME to do this for us with a **convert_fastaqual_fastq.py** command. Because we are doing this for multiple files, we can use a shell script to repeat this operation.

- 8) Using TextWrangler, open the **fastq_fasta_198.sh** file in the **files_for_week_9** folder. Notice that there are two lines in this file corresponding to two different fastq files. We will modify this .sh file to account for the five fastq data files that we'll process today. These are:

- G1_JSBN_S1.final.fastq**
- G2_JSBS_S1.final.fastq**
- G3_POTA_S1.final.fastq**
- G4_POTD_S1.final.fastq**
- G5_THM_S1.final.fastq**

- These are the sequence files that have resulted from the FLASH read merging work we did last week. In addition, we used the software package FASTX-Toolkit to filter sequences by quality. In each set of sequences, this removed approximately 25-30 sequences that were of generally low quality.

- 9) Place a copy of your newly modified **fastq_fasta_198.sh** file into the **week_9_data** folder located on the Desktop where the five fastq files are located.
- 10) Using your terminal, use a change directory to move into the **week_9_data** folder.
- 11) Check to make sure that you have sourced **macqiime**. Your command prompt should begin with the term **MacQIIME**.
- 12) Use the command **chmod 777 fastq_fasta_198.sh**.
- This is necessary to provide the correct set of permissions to this newly created file. Without doing this your computer will likely not have permission to run the script.
- 13) Now use the command **./fastq_fasta_198.sh** to execute the script that you have made.
- This will create a new folder labeled **fasta_files** on your Desktop that you specified in your script.
 - It'll probably take about 5-6 minutes to complete this job.

1056 **Building a QIIME mapping file**

1057

1058 14) Next, find the **rRNA_2016_Mapping.txt** file in the **files_for_week_9** folder on the
1059 Desktop and open this in TextWrangler.

- 1060 • This file links the individual fasta files to their particular treatment variables.
- 1061 • It can also contain information for the barcode and Primer sequences, but ours have
1062 been removed from the sequences, so these fields are left empty.

1063

1064 15) Note that the underscores in the SampleID names for our fasta files have been replaced with
1065 periods. This is because QIIME doesn't recognize the underscore character. Change the file
1066 names of the newly created fasta files in the **fasta_files** folder to account for this.

1067

1068 16) Move the **rRNA_2016_Mapping.txt** file to the **Desktop**.

1069

1070 17) In your terminal, change directory to the **fasta_files** folder.

1071

1072 18) In your terminal, use the following long command:

```
1073 validate_mapping_file.py -m  
1074 /Users/homedirectory/Desktop/rRNA_2016_Mapping.txt -o  
1075 /Users/homedirectory/Desktop/check_id_output -p -b
```

1076

- 1077 • This will check your mapping file to make sure that it matches your data files and is
1078 ready to be used in further analysis.

1079

1080 19) A new folder will be created labeled **check_id_output**. Open this in your finder and open
1081 the **rRNA_2016_Mapping.log** file in TextWrangler or click on the **.html** file and open in a
1082 web browser.

- 1083 • QIIME probably said there were errors in the mapping file when used the validation
1084 command, but your file should not list any. If you think there are errors, let us know.

1085

1086 20) Return to the terminal and give the following long command:

1087

```
1088 add_qiime_labels.py -m /Users/homedirectory/Desktop/rRNA_2016_Mapping.txt -i  
1089 /Users/homedirectory/Desktop/week_9_data/fasta_files/ -c SampleID -o  
1090 /Users/homedirectory/Desktop/QIIME_Labeled
```

1091

- 1092 • This adds a new folder to your **Desktop** labeled **QIIME_Labeled** and then creates a
1093 new sequence file labeled **combined_seqs.fna** that combines all of your individual
1094 fasta files, and adds QIIME-specific labeling to each sequence to keep track of which

1095 sequences belong to which sample. If you would like, open this in TextWrangler to
1096 look at the full data set.

1097

1098 **Assigning your sequences to Operational Taxonomic Units (OTUs)**

1099

1100 21) Next, we will use the following command to go through the process of OTU assignment:

1101

1102 **pick_de_novo_otus.py -i**

1103 **/Users/homedirectory/Desktop/QIIME_Labeled/combined_seqs.fna -o**

1104 **/Users/homedirectory/Desktop/QIIME_OTUs**

1105

- 1106 • This may take a few minutes. We've got lots of data.
- 1107 • This command makes QIIME go through the entire sequence file, which contains over
1108 60,000 sequences, and identify which sequences have high sequence similarity. The
1109 default setting here is to lump all sequences with at least 97% sequence similarity into a
1110 single representative sequence.
- 1111 • Here, an OTU is a proxy for species, assuming that all species are divergent in their
1112 sequence by at least 3%. It's not a perfect metric, but it'll work for us here.

1113

1114 **Summarize OTUs**

1115

1116 22) Use this long command to summarize all of our OTU information into a simplified table:

1117

1118 **biom summarize-table -i**

1119 **/Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -o**

1120 **/Users/homedirectory/Desktop/QIIME_OTUs/table_summary.txt**

1121

- 1122 • This doesn't do a whole heck of a lot, but it's one of the steps we need to take in the
1123 process of using QIIME to summarize our data.

1124

1125 **Generate heatmap of OTUs for each sample**

1126

1127 23) Here, we'll have QIIME generate a heatmap illustrating the diversity and abundance of
1128 different prokaryotic OTUs across our give samples. Use the following command:

1129

1130 **make_otu_heatmap.py -i**

1131 **/Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -o**

1132 **/Users/homedirectory/Desktop/QIIME_OTUs/heatmap.pdf**

1133

- 1134 • This will make a PDF file of the heatmap inside the **QIIME_OTUs** folder. Open it up.

1135 **Build summary plots of taxa (OTUs)**

1136

1137 24) Use this command to generate plots of taxonomic diversity for our samples at different
1138 hierarchical levels (e.g., class, order, family, etc.):

1139

1140 **summarize_taxa_through_plots.py -i**

1141 **/Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -o**

1142 **/Users/homedirectory/Desktop/QIIME_OTUs/taxa_summary -m**

1143 **/Users/homedirectory/Desktop/rRNA_2016_Mapping.txt**

1144

1145 • Look inside the **taxa_summary** folder and the **taxa_summary_plots** folder to find
1146 two interactive files labeled **area_charts.html** and **bar_charts.html**. You can open
1147 these in Firefox to browse a summary of results across samples.

1148 • Do any of the samples look like they are different from the others?

1149

1150 **Alpha diversity and rarefaction**

1151

1152 25) Here, we'll calculate measures of prokaryotic diversity for our samples and treatments. In
1153 addition, we'll perform some analyses to ask if we have sampled enough to accurately
1154 capture the diversity within our samples. Use the following command:

1155

1156 **alpha_rarefaction.py -i /Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom**

1157 **-m /Users/homedirectory/Desktop/rRNA_2016_Mapping.txt -o**

1158 **/Users/homedirectory/Desktop/QIIME_OTUs/arare/ -t**

1159 **/Users/homedirectory/Desktop/QIIME_OTUs/rep_set.tre**

1160

1161 • This operation will take a little time to complete. Maybe five minutes or so.

1162 • When completed, open the arare folder and we'll take a look at some of the results.

1163 • **What is Alpha diversity? What is a rarefaction plot?**

1164

1165 **SUGGESTED ANSWER:** Alpha diversity is the taxonomic diversity of each site (or
1166 sample). A rarefaction plot in this context refers to a subsampling plot where a number of
1167 sequences (e.g., 1000) is subsampled and biodiversity is calculated. This is done iteratively
1168 (at least 10x or so) and for different levels of sequencing (1000, 2000, 3000, etc.). The final
1169 plot shows how much diversity is discovered at each level of sequencing and provides a
1170 way to assess how appropriate the amount of sequencing was for each sample.

1171 **Beta diversity and differentiation among samples**

1172

1173 26) Now we'll assess a different measure of diversity, beta diversity. **What is this?** Use the
1174 following command:

1175

1176 **SUGGESTED ANSWER:** Beta diversity is the level of dissimilarity among samples or
1177 treatments. Essentially, how similar is sample A to samples B, C, etc.

1178

1179 **beta_diversity_through_plots.py -i**

1180 **/Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -m**

1181 **/Users/homedirectory/Desktop/rRNA_2016_Mapping.txt -o**

1182 **/Users/homedirectory/Desktop/QIIME_OTUs/betadiv_even/ -t**

1183 **/Users/homedirectory/Desktop/QIIME_OTUs/rep_set.tre -e 8000 -f**

1184

1185 We'll take a look at the results from these analyses and discuss.

1186 **Week 7 - Bioinformatics, Part 3: Quantitative analysis and revisiting our prokaryotic**
1187 **diversity hypotheses**

1188
1189 This week we will use our larger sequence data set to measure prokaryotic diversity and
1190 statistically test hypotheses about how this diversity varies across our sampled environments. We
1191 will do this in QIIME. This will be a tractable exercise in our two-hour lab because the very
1192 time-consuming steps in the analysis (e.g., the `pick_de_novo_otus.py` step) have already been
1193 performed for you.

1194
1195 One thing to note going into these analyses: some samples have been removed from the analysis
1196 due to low sequence output. This is particularly true for a handful of samples (8) for which we
1197 recovered < 40 sequences. To provide somewhat reasonable estimates of prokaryotic diversity,
1198 we removed any sample with < 5000 sequences (merged reads). Some of these were PCR
1199 negative controls. In total, we will still analyze a collection of 32 different sequenced samples. A
1200 list of all samples, the number of merged sequences per sample, and which ones were left out of
1201 analysis can be found in the Google Drive link that has been emailed to everyone.

- 1202
- 1203 1) On your Desktop, check for the **16SrRNA_data_week_10** folder and look inside to make
1204 sure it contains a **fasta_files** folder with 48 (**.fna**) files.
1205
 - 1206 2) Now check for the **QIIME_Labeled** folder on your Desktop. We have performed this step in
1207 the analysis in preparation for today's lab, which is saving us computation time.
 - 1208 • This contains the combined sequences in fasta format from all of the samples to be
1209 analyzed.
 - 1210
 - 1211 3) And, check for the **QIIME_OTUs** folder on your Desktop. Here, we have already performed
1212 the following steps:
 - 1213 • `pick_de_novo_OTUs.py`
 - 1214 • `biom summaize-table`
 - 1215 • `make_otu_heatmap.py`
 - 1216
 - 1217 • Remember, sequences that were < 3% divergent from each other were binned together
1218 into the same "species" or OTU category. This step also made a record of the number of
1219 sequences identified for each OTU, which provides a measure of abundance, or
1220 frequency, of that particular OTU.
 - 1221 • You may find the heatmap interesting, or not.
 - 1222
 - 1223 4) Next, locate on your Desktop the **Full_data_processing_files** folder, which contains a
1224 number of scripts and analysis files used in the analysis of the complete data set.
 - 1225 • Note that some of these have already been run on the data (e.g., read merging and quality
1226 filtering scripts). We provide them here so you have access to everything that has been
1227 done with the data.
 - 1228
 - 1229 5) Move the **rRNA_2016_48_Mapping.txt** file to your Desktop.

- 1230 6) Open the **rRNA_2016_48_Mapping.txt** in Textwrangler. You should see lines for the 32
 1231 input files that we are including in our analysis. Note that we currently have only one
 1232 **Treatment** category defined here, specifying if the building sampled was a tall (**high**) or
 1233 short (**low**) building.
 1234
- 1235 7) First, lets generate some histograms of prokaryotic diversity for the different samples using
 1236 the following command:
 1237
- ```
1238 summarize_taxa_through_plots.py -i

 1239 /Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -o

 1240 /Users/homedirectory/Desktop/QIIME_OTUs/taxa_summary_sample -m

 1241 /Users/homedirectory/Desktop/rRNA_2016_48_Mapping.txt
```
- 1242
- 1243 • Note that a copy of this command is located in the **Long\_commands\_week10.txt** file  
 1244 found inside the **Full\_data\_processing\_files** folder.
- 1245
- 1246 8) Locate these results in the **QIIME\_OTUs/taxa\_summary\_sample/taxa\_summary\_plots**  
 1247 folder and browse your results in the **bar\_charts.html** file. This should open in your web  
 1248 browser.  
 1249
- 1250 9) Now lets generate the same histograms summarized by our first treatment using the  
 1251 following command:  
 1252
- ```
1253 summarize_taxa_through_plots.py -i  

  1254 /Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -o  

  1255 /Users/homedirectory/Desktop/QIIME_OTUs/taxa_summary_treatment -m  

  1256 /Users/homedirectory/Desktop/rRNA_2016_48_Mapping.txt -c Treat1 -f
```
- 1257
- 1258 10) Find your results and see if there appear to be any differences.
 1259
- 1260 11) Next, lets calculate measures of alpha diversity for samples and treatment, and perform
 1261 rarefaction analysis of the data using the following command:
 1262
- ```
1263 alpha_rarefaction.py -i /Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -

 1264 m /Users/ homedirectory /Desktop/rRNA_2016_48_Mapping.txt -o /Users/

 1265 homedirectory /Desktop/QIIME_OTUs/arare/ -t /Users/ homedirectory

 1266 /Desktop/QIIME_OTUs/rep_set.tre -f
```
- 1267
- 1268 12) Locate the results of this analysis in the **QIIME\_OTUs/arare/alpha\_rarefaction\_plots**  
 1269 folder.  
 1270
- 1271 13) This analysis will provide three measures of alpha diversity. What are they and how do they  
 1272 differ in how they calculate alpha diversity? You may need to do some research on these to  
 1273 understand what each measure is telling you.  
 1274

1275 **SUGGESTED ANSWER:** This answer will depend on the metric being calculated. Some  
1276 common defaults are. Shannon = Shannon's diversity index, which is a quantitative measure that  
1277 reflects how many different types (species or otherwise) are in a sample. Observed OTUs is  
1278 simply the total number of OTUs observed, purely from a count perspective. Documentation of  
1279 the various metrics that can be used can be seen here: [http://scikit-](http://scikit-bio.org/docs/latest/generated/skbio.diversity.alpha.html)  
1280 [bio.org/docs/latest/generated/skbio.diversity.alpha.html](http://scikit-bio.org/docs/latest/generated/skbio.diversity.alpha.html)

1281  
1282 14) What can you say about alpha diversity for the various samples? What about diversity when  
1283 comparing tall vs. short buildings?  
1284

1285 15) Based on the rarefaction plots for various measures of alpha diversity, what can you say  
1286 about the sufficiency of our data to estimate alpha diversity for our samples or treatments?  
1287 • Note that you can save a copy of the **arare** folder and its contents for later reference. A  
1288 flash drive may work best for this step. You'll need to copy the entire folder to be able to  
1289 access the results on another computer.  
1290

1291 16) OK, now calculate measures of beta diversity and a principal component plot for your  
1292 samples and treatment using the following command:  
1293

```
1294 beta_diversity_through_plots.py -i
1295 /Users/homedirectory/Desktop/QIIME_OTUs/otu_table.biom -m
1296 /Users/homedirectory/Desktop/rRNA_2016_48_Mapping.txt -o
1297 /Users/homedirectory/Desktop/QIIME_OTUs/betadiv_even/ -t
1298 /Users/homedirectory/Desktop/QIIME_OTUs/rep_set.tre -f
1299
```

1300 • You will probably receive a minor warning about a negative eigenvalue, which is a  
1301 statistic used in the principal component analysis. It's OK.  
1302

1303 17) Open the **unweighted\_unifrac\_dm.txt** file located in the **QIIME\_OTUs/betadiv\_even**  
1304 folder. You'll want to use Textwrangler for this.

1305 • Here you are seeing a matrix of distance measures between samples (i.e., beta diversity),  
1306 measured using a calculation termed UniFrac, which takes into account phylogenetic  
1307 distance between samples. What is the difference between a weighted and unweighted  
1308 measure? Is there a good reason to use, or not to use, the weighted measure of beta  
1309 diversity for our study?  
1310

1311 **SUGGESTED ANSWER:** Weighted versus unweighted distance measures refer to whether  
1312 OTU abundances are accounted for (weighted) or ignored (unweighted) when comparing  
1313 distances among samples. Therefore, if the goal is simply to assess what taxa are present,  
1314 unweighted distance should be used.  
1315

1316 18) It's not easy to interpret all of those pairwise comparisons of beta diversity between samples.  
1317 There are just too many combinations to identify a clear pattern. So, instead, you've also  
1318 performed a principal component analysis (PCA) on the data. Open the **index.html** file in the

1319 **QIIME\_OTUs/betadiv\_even folder**. You'll need to open this in a web browser to visualize  
1320 the PCA plot.

- 1321 • PCA is a way of taking complex data (like the beta diversity matrix) and condensing it  
1322 down to a smaller number of dimensions (PC axes) that explain all of that variation. It  
1323 provides a much more visually easy way to look for patterns in your results.
- 1324 • Here's a cool interactive website that provides an explanation of PCA:  
1325 <http://setosa.io/ev/principal-component-analysis/>
- 1326 • Do you see any beta diversity patterns in that would suggest differences among samples  
1327 or treatment?

1328  
1329 19) Now let's ask if we have a statistical difference between our two treatment groups. We'll do  
1330 this by testing for a significant difference in beta diversity. Use the following command:

```
1331 compare_categories.py --method anosim -i
1332 /Users/homedirectory/Desktop/QIIME_OTUs/betadiv_even/unweighted_unifrac_dm.tx
1333 t -m /Users/homedirectory/Desktop/rRNA_2016_48_Mapping.txt -c Treat1 -o
1334 /Users/homedirectory/Desktop/QIIME_OTUs/anosim_out -n 99
```

1335  
1336  
1337 20) Locate the results file in the **QIIME\_OTUs/anosim\_out** folder. Do the results support the  
1338 hypothesis that taller buildings on campus have greater prokaryotic diversity compared to  
1339 shorter buildings?

1340  
1341 21) OK, so you've run through the analysis of the full data set and tested the main hypothesis that  
1342 was proposed for this project. What other hypotheses can be tested with our sampling design  
1343 and data? Other differences exist among the buildings and elevators we sampled, and we  
1344 collected a range of samples within elevators.

1345  
1346 For the remainder of the lab, you will work in your groups to develop additional hypotheses  
1347 to be tested. You will then address these hypotheses using your knowledge of how to  
1348 calculate measures of diversity and perform a statistical test.

1349  
1350 This will require you to modify the **rRNA\_2016\_48\_Mapping.txt** in Textwrangler to add in  
1351 the new treatments. We will be available to help you with this and deal with any errors that  
1352 may arise.

1353 **Supplementary Materials – Survey Questions.**

1354

1355 **Course survey** – A pre-course survey was distributed to gauge student interest, understanding,  
1356 and prior knowledge of various topics. We performed a follow-up, post-course survey at the end  
1357 of the semester in order to gauge how the topics discussed in the course addressed knowledge  
1358 gaps and peaked student interest.

1359

1360 The survey was developed and administered using the online tool Survey Monkey  
1361 ([www.surveymonkey.com](http://www.surveymonkey.com)) and was completely anonymous. Questions marked with an asterisk  
1362 were additions for the 2016 CURE.

1363 **Pre-course Survey:** Questions from the 2016 iteration of the course.

1364

1365 Please use the rating scale below to answer the following questions:

1366

1367 Rating scale: (a) None, (b) little experience/comfort, (c) intermediate experience/Some comfort,  
1368 or (d) highly experienced/comfort

1369

1370 1. Please rate your previous experience in biology classes with the following subjects:

1371 a. Scientific method (e.g., deductive logic)

1372 b. Hypothesis generation

1373 c. Experimental design

1374 d. Use of molecular lab equipment

1375 e. Data analysis

1376 f. Problem solving

1377 g. Communicating science (via oral or written reports)

1378 h. Reading peer review literature

1379

1380 2. For any topics that you rated either intermediate experience or highly experienced please  
1381 provide a specific example.\*

1382

1383 3. Please rate your general knowledge of the following topics

1384 a. PCR

1385 b. Genetic analysis

1386 c. Statistics

1387 d. Primary literature

1388 e. Bioinformatic analysis

1389

1390 4. Please rate your comfort level with the following areas (e.g., how comfortable would you  
1391 be with conducting the following on your own)

1392 a. Generating a hypothesis from a given observation

1393 b. Designing an experiment to test a hypothesis

1394 c. Analyzing data for patterns

1395 d. Laboratory techniques

1396 e. Figuring out the next step in an experiment

1397 f. Ability to work independently

1398 g. Ability to work with others

1399 h. Communicating my work

1400 i. Time management

1401 j. Understanding peer review literature

1402

1403 5. What about biology excites you?\*

1404

1405 6. Please describe what a biologist is to you. What type of skills would you expect them to  
1406 have? \*

1407

1408

- 1409      7. Please briefly describe your long-term career goals  
1410  
1411      8. Please briefly describe what you hope to gain from this course

1412 **Post-course Survey:**

1413

1414 Please use the rating scale below to answer the following questions:

1415

1416 Rating scale: (a) No improvement, (b) improved slightly, (c) moderately improved, or (d) greatly  
1417 improved

1418

1419 1. Please rate how the overall content and experience in the course has influenced your  
1420 knowledge of the following topics:

1421 a. Scientific method (e.g., deductive logic)

1422 b. Hypothesis generation

1423 c. Experimental design

1424 d. Use of molecular lab equipment

1425 e. Data analysis

1426 f. Problem solving

1427 g. Communicating science (via oral or written reports)

1428 h. Reading peer review literature

1429

1430 2. Please rate how the 16s rRNA project contributed in your knowledge in the following  
1431 topics:

1432 a. Scientific method (e.g., deductive logic)

1433 b. Hypothesis generation

1434 c. Experimental design

1435 d. Use of molecular lab equipment

1436 e. Data analysis

1437 f. Problem solving

1438 g. Communicating science (via oral or written reports)

1439 h. Reading peer review literature

1440

1441 3. Please provide a specific example of a topic that you moderately or greatly improved on\*

1442

1443 4. Please rate how this course has influenced your general knowledge of the following  
1444 topics:

1445 a. PCR

1446 b. Genetic analysis

1447 c. Statistics

1448 d. Primary literature

1449 e. Bioinformatic analysis

1450

1451 5. Please rate how the overall content and experience in this course has changed your  
1452 comfort level (e.g., how comfortable would you be conducting the following on your  
1453 way) with the following topics:

1454 a. Generating a hypothesis from a given observation

1455 b. Designing an experiment to test a hypothesis

1456 c. Analyzing data for patterns

1457 d. Laboratory techniques

- 1458 e. Figuring out the next steps in an experiment  
1459 f. Ability to work independently  
1460 g. Ability to work with others  
1461 h. Communicating my work  
1462 i. Time management  
1463 j. Understanding peer-review literature  
1464
- 1465 6. Please rate how the 16s rRNA project influenced your comfort level with the following  
1466 areas:
- 1467 a. Generating a hypothesis from a given observation  
1468 b. Designing an experiment to test a hypothesis  
1469 c. Analyzing data for patterns  
1470 d. Laboratory techniques  
1471 e. Figuring out the next steps in an experiment  
1472 f. Ability to work independently  
1473 g. Ability to work with others  
1474 h. Communicating my work  
1475 i. Time management  
1476 j. Understanding peer-review literature  
1477
- 1478 7. Regarding question 6, please provide a specific example of a topic that you moderately or  
1479 greatly improved on\*
- 1480
- 1481 8. Please briefly describe how this course has or hasn't influenced your view of what a  
1482 biologist does\*
- 1483
- 1484 9. Please provide any feedback about the 16s RNA project