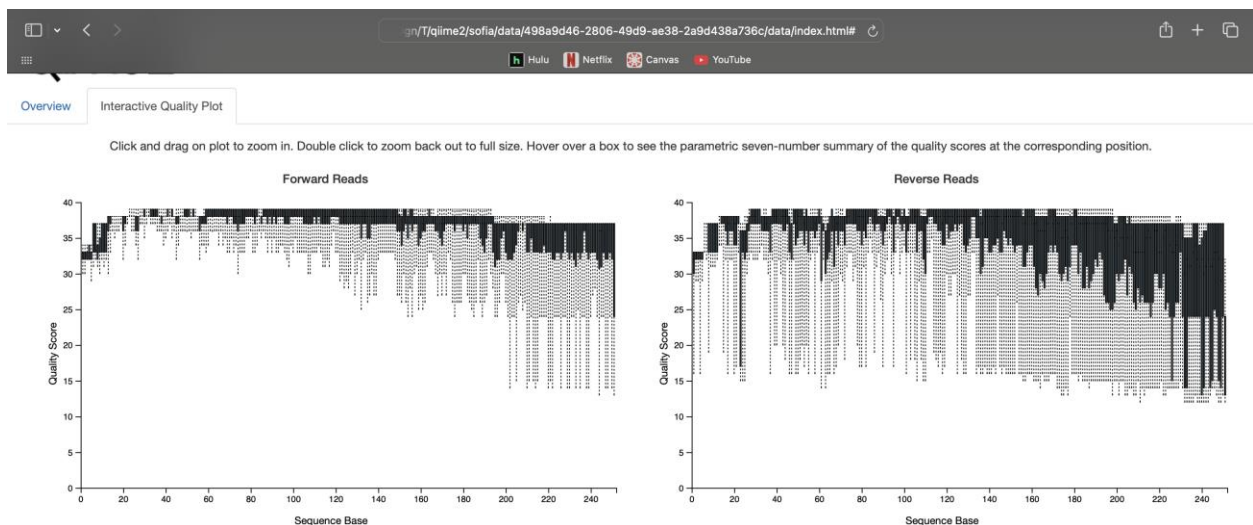


1) Examine the metadata file. What columns do you see that you think might be useful for later when doing alpha and beta diversity metrics? HINT: there are 3

The columns that might be useful for later when doing alpha and beta diversity metrics are the population category, sex category, and flock category.

2) Include a screenshot of your interactive quality plot. Based on this plot, what values would you choose for **--p-trunc-len** and **--p-trim-left** for both the forward and reverse reads? Why have you chosen those numbers? HINT: If you trim and truncate too much, you will lose too many of your reads, making your downstream analysis not useful. Think about it scientifically and only trim and truncate where the overall quality averages below 25-30.



Based on the plot, the values I would choose for **--p-trunc-len** forward is 240 and for the reverse is 200. The value for **--p-trim-left** forward is 0 and for reverse its 0. I chose these numbers because both graphs have good quality data for the left sequences so there is no need to trim. However, there is a need to truncate because some of the qualities drop below 30. For the forward reads it drops below 30 around 240 and for the reverse reads it's around 200.

3) Include a screenshot of the table summary from visualizing your table and a screenshot of the sequence length statistics from the rep-seqs file. Remember, we may eventually want to cut any samples with less than 10,000 reads. Do you see any in the interactive sample detail that might need to be cut? If so, which ones?

## Table summary

Summary Statistic	Value
Number of samples	24
Number of features	5,948
Total frequency	526,012

There are 8 frequencies under 10,00: 84\_S61\_L001, 254\_S69\_L001, 168\_S37\_L001, 265\_S133\_L001, 100\_S359\_L001, 125\_S13\_L001, 104\_S93\_L001, 189\_S23\_L001.

## Sequence Length Statistics

Download sequence-length statistics as a TSV

Sequence Count	Min Length	Max Length	Mean Length	Range	Standard Deviation
5948	240	426	252.37	186	9.65

4) Once you have generated your taxonomy visualization, sort it by confidence. What are your top hits? What about if you sort by taxon? What hits do you see?

When you visualize the taxonomy file and sort it by confidence the top hits are:

k\_\_Bacteria; p\_\_Proteobacteria; c\_\_Deltaproteobacteria; o\_\_Myxococcales  
k\_\_Bacteria; p\_\_Proteobacteria; c\_\_Alphaproteobacteria; o\_\_Ellin329; f\_\_ ; g\_\_ ; s\_\_

k\_\_Bacteria; p\_\_Bacteroidetes; c\_\_Cytophagia; o\_\_Cytophagales; f\_\_Cytophagaceae;  
g\_\_Spirosoma; s\_\_

k\_\_Bacteria; p\_\_Proteobacteria; c\_\_Gammaproteobacteria; o\_\_Acidithiobacillales;  
f\_\_Acidithiobacillaceae; g\_\_Acidithiobacillus; s\_\_.

When you visualize the taxonomy file and sort it by taxon the top hits are:

k\_\_Archaea; p\_\_Crenarchaeota; c\_\_Thaumarchaeota; o\_\_Cenarchaeales; f\_\_SAGMA-X; g\_\_;  
s\_\_

k\_\_Archaea; p\_\_Crenarchaeota; c\_\_Thaumarchaeota; o\_\_Nitrososphaerales;  
f\_\_Nitrososphaeraceae; g\_\_Candidatus Nitrososphaera; s\_\_gargensis

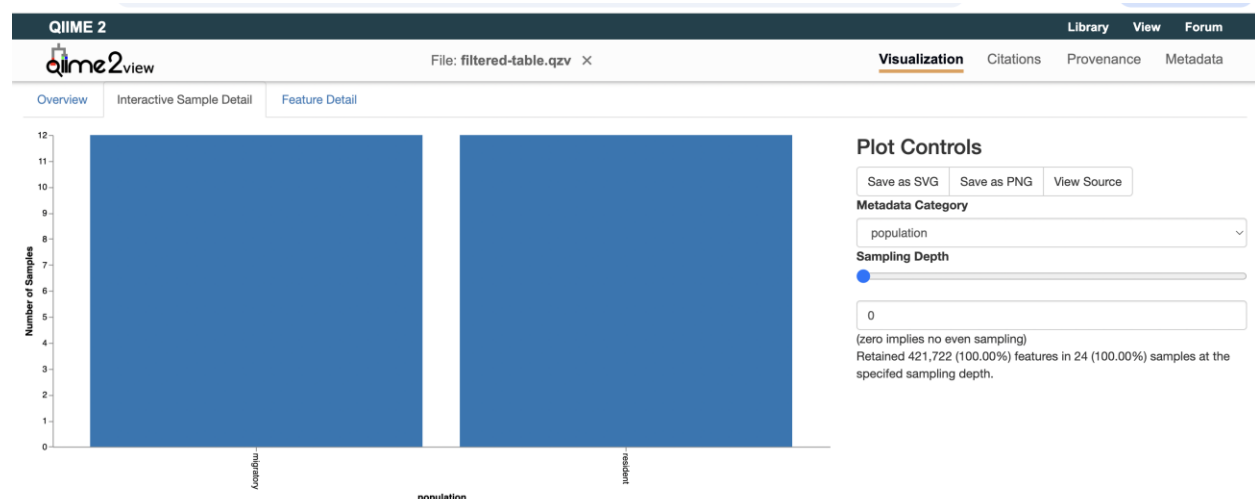
k\_\_Archaea; p\_\_Crenarchaeota; c\_\_Thaumarchaeota; o\_\_Nitrososphaerales;  
f\_\_Nitrososphaeraceae; g\_\_Candidatus Nitrososphaera; s\_\_gargensis

k\_\_Archaea; p\_\_Euryarchaeota; c\_\_Halobacteria; o\_\_Halobacteriales; f\_\_Halobacteriaceae

### 5) When you visualize level 3 of taxonomy, what level is this? Do you see any trends as you sort by various metadata categories?

Level 3 of taxonomy represents class. At level 3 taxonomy you can see all the species are from the Bacteria kingdom and the phylum/class differ from species to species. When the chart is sorted by population, we can see some samples have higher level of richness and evenness while others are dominated by few bacterial groups. When the chart is sorted by sex, we can see the compositions stay consistent between sexes with no major shift on either side. When the chart is sorted by flock, we can see the species compositions vary greatly between flocks.

### 6) After visualizing your filtered-table.qzv, what cutoff value will you use for generating alpha and beta diversity? Why? Include a screenshot of the interactive sample view to help justify your reasoning.



Sample ID	population	Frequency
106_S98_L001		48,643
13_S95_L001		43,345
18_S71_L001		29,814
174_S146_L001		27,837
212_S94_L001		27,406
4_S157_L001		26,985
309_S47_L001		25,639
364_S22_L001		25,451
78_S46_L001		25,288
128_S36_L001		20,229
366_S45_L001		18,684
307_S70_L001		17,596
163_S60_L001		14,482
245_S122_L001		12,809
61_S109_L001		12,596
385_S170_L001		11,151
254_S69_L001		8,285
168_S37_L001		7,538
84_S61_L001		7,273
265_S133_L001		4,432
100_S359_L001		2,609
125_S13_L001		2,097
104_S93_L001		1,128
189_S23_L001		405

Based on the data, the cutoff value for generating an alpha and beta diversity would be 10000. Once the data reaches the 10000 mark there seems to be a greater decline in frequencies. This implies that those data points have possible sequencing failure so, by choosing 10000 as the cutoff I keep 20 good samples and remove 8 moderate ones. By removing these samples, the strong statistical power of the data is kept

**7) The first metric we will analyze is alpha diversity. In your own words, what is alpha diversity and what are the differences between the two types of alpha diversity we will analyze (Shannon and Observed features)?**

Alpha diversity measures the diversity of species in a single sample/environment. This can show us the different types of species in an environment and how they are dispersed. A Shannon Alpha diversity measures the richness and evenness of the species in an environment. This differs from Observed features alpha diversity because it accounts for both the number of species in the environment and how evenly they are distributed. The Observed features alpha diversity measures the richness of species in an environment. It counts the number of different species in an environment. This differs from the Shannon alpha diversity because it doesn't consider distribution of the population.

**8) Since you are looking at two metrics of diversity for 3 metadata categories, it would be helpful to make a table of the significant values. Are any of your comparisons significant? For one of the metadata data columns, there are actually 4 options. Include a screenshot of**

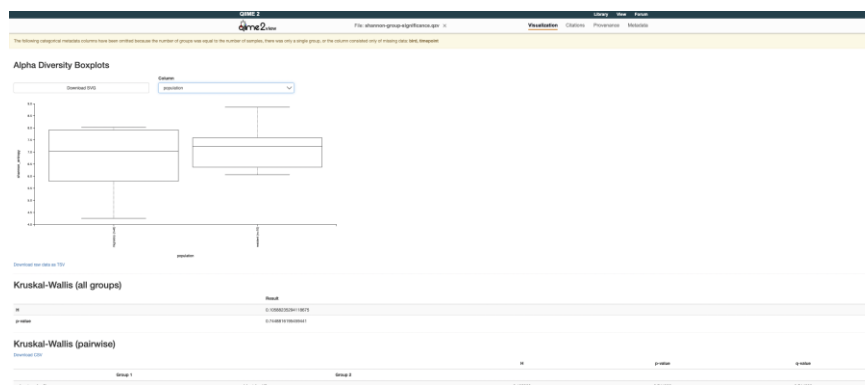
**the pairwise comparisons for Shannon and Observed. Are there any significant comparisons? (HINT: look at the q value)**

When comparing the q values for the Shannon file and observed features file there seems to be no significance between the categories. For the populations there is no significant difference between the results with both being greater than 0.05. For the sex category there is no significant difference with both values being greater than 0.05. Finally, for the flock category there is a difference between the two graphs with there being a significance for the Shannon data and no significance for the observed features.

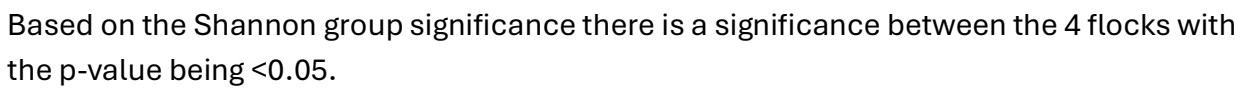
	H Shannon	P-Value Shannon	q-Value Shannon	H Observed	P-Value Observed	Q-Value Observed
Population	0.105882	0.744881	0.744882	0.294118	0.587594	0.587594
Sex	0.002801	0.957791	0.957791	0.473389	0.491432	0.49143
Flock	0.428571 3.750000 0.060000 3.84000	0.064078 0.052808 0.806496 0.050044	0.090960 0.806496 0.090960 0.090960	3.428571 3.750000 0.060000 0.960000	0.064078 0.052806 0.806496 0.327187	0.128155 0.806496 0.490780 0.557643

### Shannon-group-significance

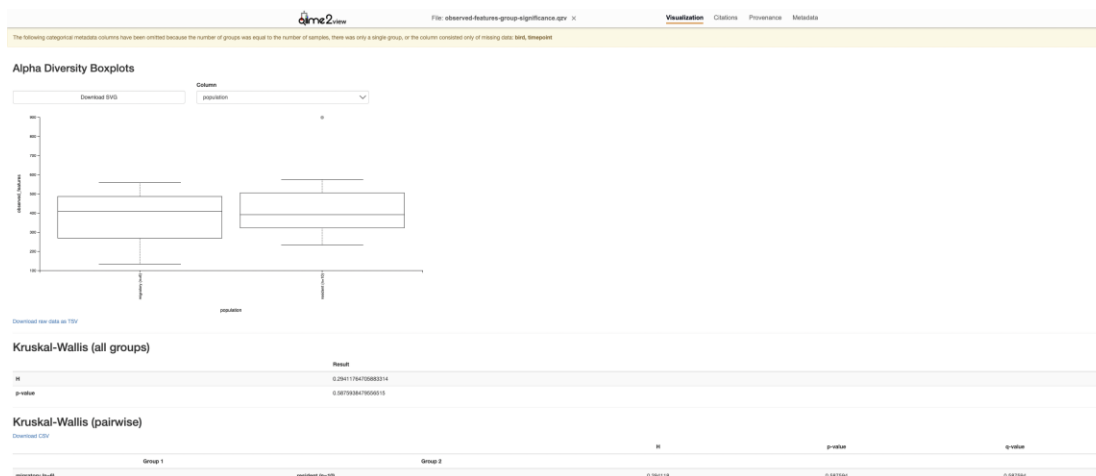
Based on the Shannon group significance there is no significance between the populations with the p value being  $g>0.05$ .



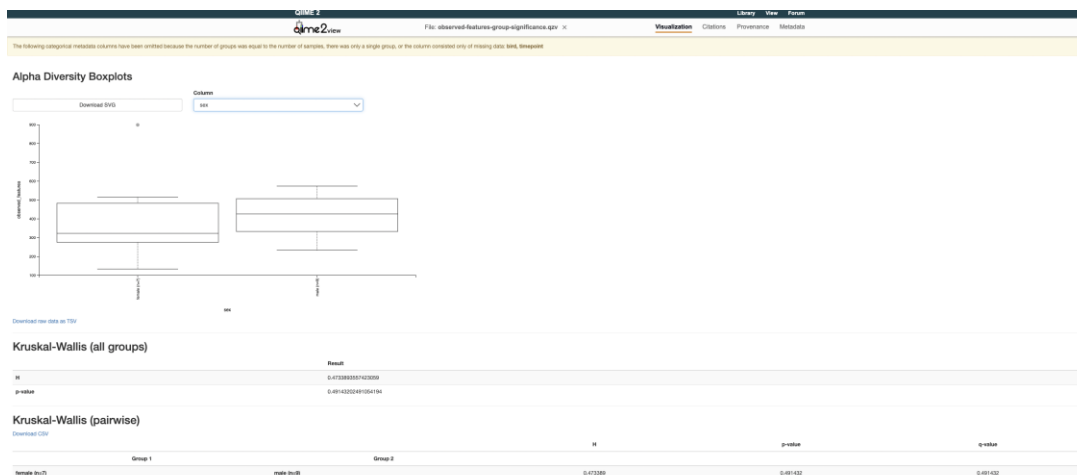
Based on the Shannon group significance there is no significance between the two sexes with the p-value being  $>0.05$ .



For the overserved features group significance table there is no significance between the two populations with the p-value being  $>0.05$ .



Based on the observed features group significance table there is no significance between the two sexes with the p-value being  $>0.05$ .



Based on the observed features group significance there is no significance with the overall group p-value being greater than 0.05.

Using the bray-Curtis visualization for the population data there is significance with the p-value being less than 0.05.



QIIME 2		Library	View	Forum
		File: bray-curtis-sex-significance.qzv	X	
		Visualization	Citations	Provenance
				Metadata
Overview				
PERMANOVA results				
method name	PERMANOVA			
test statistic name	pseudo-F			
sample size	16			
number of groups	2			
test statistic	1.249632			
p-value	0.103			
number of permutations	999			


Using the Bray-Curtis visualization for the sex data there is no significance with the P-value being greater than 0.05.

QIIME 2		Library	View	Forum
		File: bray-curtis-flock-significance.qzv	X	
		Visualization	Citations	Provenance
				Metadata
Overview				
PERMANOVA results				
method name	PERMANOVA			
test statistic name	pseudo-F			
sample size	16			
number of groups	4			
test statistic	1.530614			
p-value	0.002			
number of permutations	999			

Using the Bray-Curtis visualization for the flock data there is significance with the P-value being less than 0.05.

QIIME 2		Library	View	Forum
		File: weighted-unifrac-population-significance.qzv	X	
		Visualization	Citations	Provenance
				Metadata
Overview				
PERMANOVA results				
method name	PERMANOVA			
test statistic name	pseudo-F			
sample size	16			
number of groups	2			
test statistic	2.294866			
p-value	0.038			
number of permutations	999			

Using the Weighted-Unifrac visualization for the population data there is significance with the P-value being less than 0.05.

QIIME 2		Library	View	Forum
		File: weighted-unifrac-sex-significance.qzv	X	
		Visualization	Citations	Provenance
				Metadata
Overview				
PERMANOVA results				
method name	PERMANOVA			
test statistic name	pseudo-F			
sample size	16			
number of groups	2			
test statistic	0.99437			
p-value	0.456			
number of permutations	999			

Using the Weighted-Unifrac visualization for the sex data there is no significance with the P-value being greater than 0.05.

qime2view

File: weighted-unifrac-flock-significance.qzv x

Visualization

Citations

Provenance

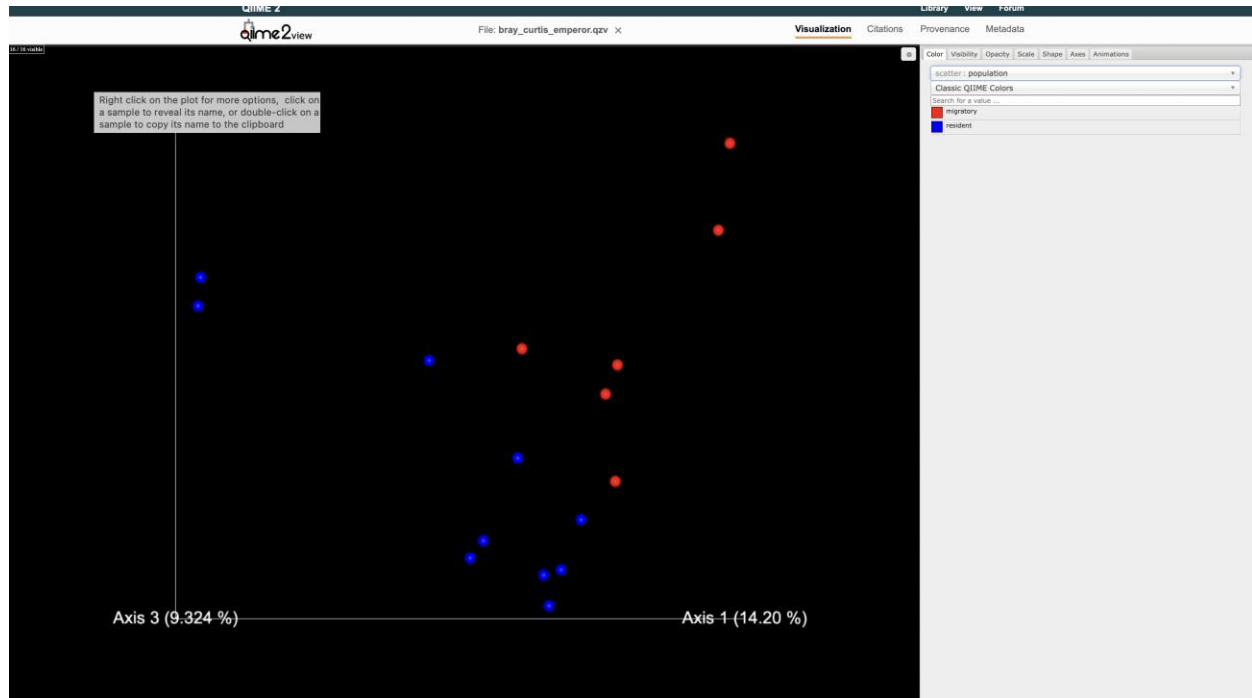
Metadata

Overview

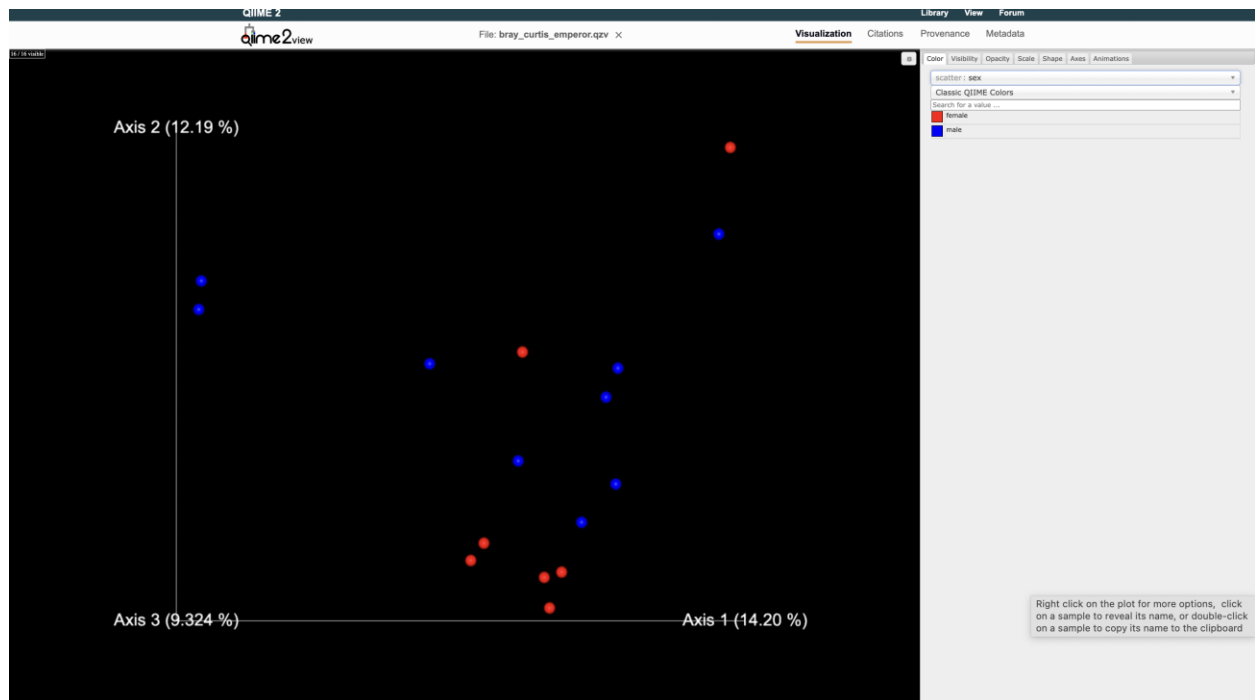
	PERMANOVA results
method name	PERMANOVA
test statistic name	pseudo-F
sample size	16
number of groups	4
test statistic	1.991901
p-value	0.016
number of permutations	999

Using the Weighted-Unifrac visualization for the flock data there is significance with the P-value being less than 0.05.

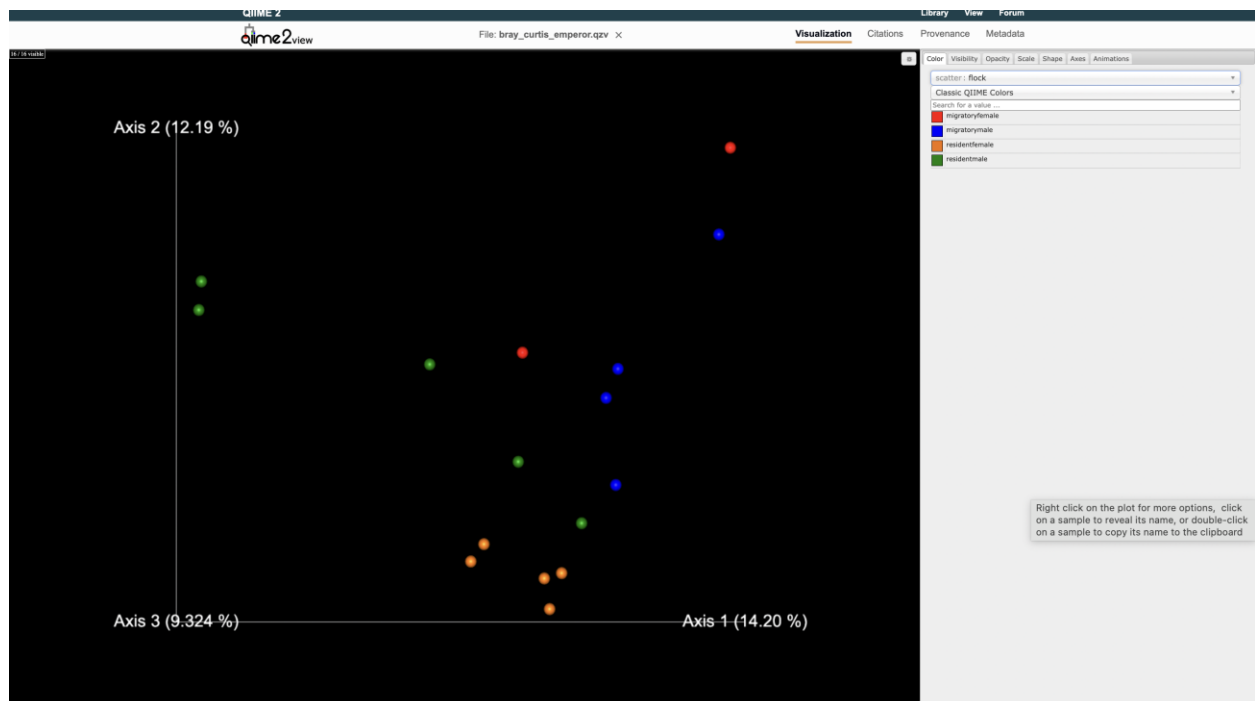
11) The **core-metrics-phylogeny** command generates a file called **bray-curtis-emperor.qzv**. Include 3 screenshots total (where the points are colored based on the metadata metrics). How do these results help you make sense of the results you got from question 10?



This graph shows the abundance of species in the populations migratory and resident. Based on the graph, the two species are not similar which will mean there is a significant difference between them. This matches the p value from the Bray-Curtis table.



This Bray-Curtis graph shows the relationship between the two sexes. Based on this graph the two sexes are somewhat similar because they are paired together on the graph. This matches the p-value from the graph which is greater than 0.05.



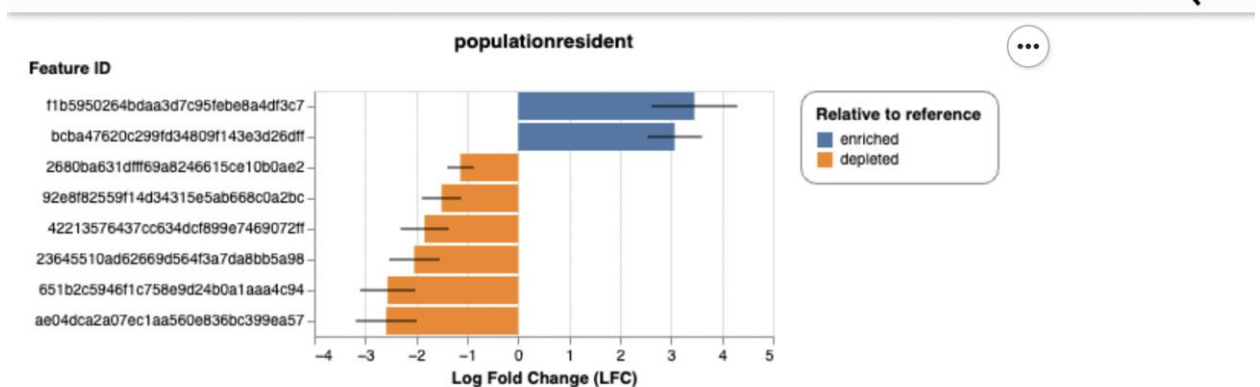
This graph shows the Bray-Curtis for the flock category. The flock category includes migratory female, migratory male, resident female, and resident male. Based on the graph, the resident female and migratory female are very different from one another; the resident male and resident

female are like one another; the migratory male and resident male are like one another; and the migratory male and migratory female are like one another. Overall, many of the flock points are far apart from one another indicating a difference between the species. This matches the low p-value which indicates a significant difference between the flocks.

**12) Using ANCOMBC, do you find any specific taxa that are differentially expressed with your three metadata categories? Please include a screenshot of any differentially expressed taxa and identify the species using the taxonomy.qzv file.**

Using ANCOMBC for the population there was a graph generated for population resident. Using the taxonomy.qzv files the species' taxonomy can be identified.

Feature ID	Taxonomy
42213576437cc634dcf899e7469072ff	s__marcusii



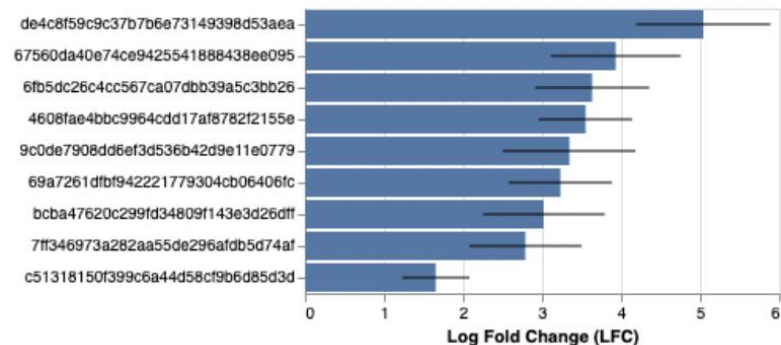
- Using ANCOMBC for the flock there were three graphs generated for flockmigratorymale, flockresidentfemale, and flock resident male.

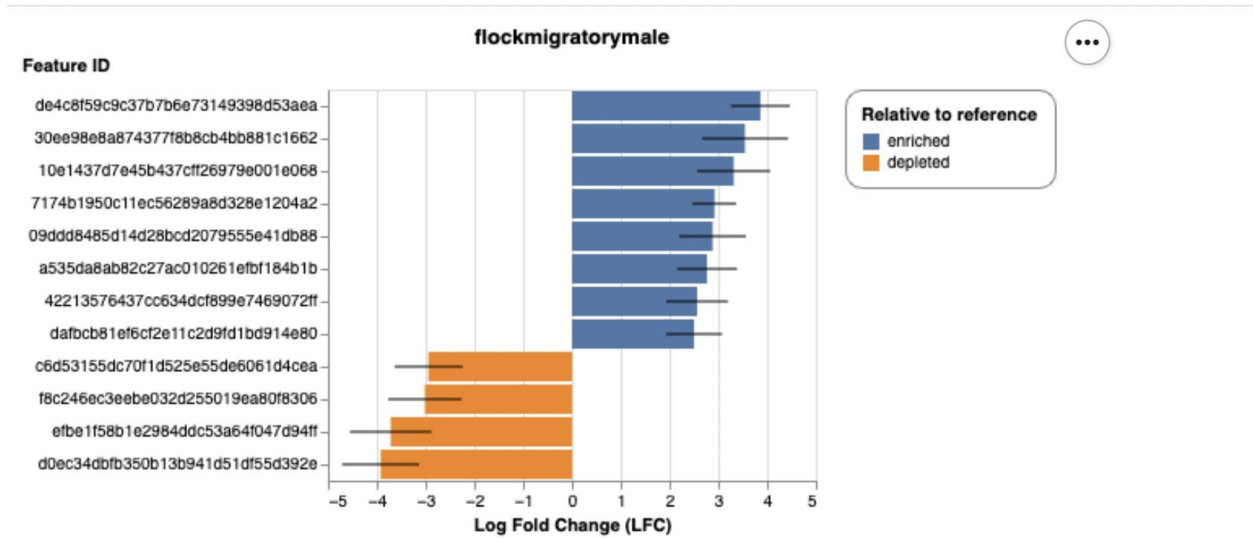
Feature ID	Taxonomy
42213576437cc634dcf899e7469072ff	s__marcusii



Feature ID

flockresidentmale





Using ANCOMBC for the sex there was a graph generated for sexmale.

Feature ID	Taxonomy
ee5cfac5c165c97e0a1475e8d9b4c297	s__pirum

