

Advanced Microbiome Analysis 2025

Contents

I	Introduction	5
1	Welcome	7
1.1	AWS/UNIX primer	7
1.2	Course Schedule	7
1.3	Meet Your Faculty	7
II	Modules	9
2	Module 1 Lab: Introduction to metagenomics and read-based profiling	11
2.1	Overview	11
2.2	1. Initial setup and pre-processing of metagenomic reads	12
2.3	Quality Control	14
2.4	2. Generating Taxonomic Profiles	17
2.5	3. Visualisation of results in R	19
2.6	Answers	28
III	Modules	31
3	Module 2: Metagenomic Assembly and Binning	33
3.1	Table of Contents	33
3.2	Introduction	33
3.3	2.1 Initial setup	34
3.4	2.2 Assembly of raw reads with MEGAHIT	35
3.5	2.3 Make an Anvi'o contigs databases	37
3.6	2.4 Run HMMs to identify single copy genes	38
3.7	2.5 Map samples onto contigs with Bowtie2	39
3.8	2.6 Cluster contigs into bins	42
3.9	2.7 Interactive refinement of bins	45
3.10	2.8 Export the final bins into MAGs	47
3.11	2.9 Run CheckM2	49

3.12	Other things that we frequently do with MAGs	50
4	Module 3 Lab: Assigning Functions	51
4.1	Setting up the workspace	51
4.2	Overview	56
4.3	Prodigal	57
4.4	Sequence similarity	61
4.5	Machine learning	71
4.6	gapseq	82
5	Module 4 Lab: Visualization and Finding Functional Signifi-	
	cance	87
5.1	Setting up the workspace	87
5.2	Overview	89
5.3	Salmon	90
5.4	DESeq2	92
5.5	Gene set enrichment analysis (GSEA)	96
5.6	Cytoscape	106

Part I

Introduction

Chapter 1

Welcome

Welcome to CBW's Advanced Microbiome Analysis 2025 Workshop!

1.1 AWS/UNIX primer

You can find the schedule and materials for the AWS/UNIX primer [here](#).

1.2 Course Schedule

Time (PST)	Day 1	Time (PST)	Day 2
8:30	Arrivals & Check-in	8:30	Arrivals
9:00	Welcome (Coordinator)	9:00	Module 3 Lecture
9:30	Module 1 Lecture		
10:30	Break (30min)	10:00	Break (30min)
11:00	Module 1 Lab	10:30	Module 3 Lab
13:00	Lunch (1h)	12:30	Class photo + Lunch (1h)
14:00	Module 2 Lecture	13:30	Module 4 Lecture
15:00	Break (30min)	14:30	Break (30min)
15:30	Module 2 Lab	15:00	Module 4 Lab
		17:00	Survey and Closing
			Remarks
17:30	Finished	17:30	Finished

1.3 Meet Your Faculty

Coming soon!

Part II

Modules

Chapter 2

Module 1 Lab: Introduction to metagenomics and read-based profiling

This tutorial is part of the 2025 CBW Advanced Microbiome Analysis (held in Vancouver, BC, May 29-30). It is based on a version written for the 2024 CBW Advanced Microbiome Analysis workshop by Ben Fisher.

Author: Robyn Wright

2.1 Overview

The goal of this tutorial is to familiarise students with processes by which we analyse metagenomic data and classify taxa within our samples. Shotgun Metagenomic Sequencing, sometimes called MGS or WGS, is capable of capturing any DNA extracted from a given sample (however, this does not necessarily mean it captures ALL of the DNA). With MGS reads, we must consider that there could be significant host contamination, so - depending on our sample type - we must filter against host sequences in our pipeline. We will then classify the taxa in our samples using two popular approaches - an all reads approach (Kraken2 + Bracken) and a marker-gene approach (MetaPhlAn). In all steps of bioinformatics, there are many tools that can work to produce similar results, and there is never a one-size-fits-all solution. It is up to us to learn about the options that are available and choose which one is most appropriate for our application - people often have different opinions, so sometimes we need to educate ourselves, make a decision, and be able to justify this decision to others. This lab is therefore a foray into some popular tools and processes, and how to appropriately use them for our analyses.

Throughout this module, there are some questions aimed to help your understanding of some of the key concepts. You'll find the answers at the bottom of this page, but no one will be marking them. *You may wish to keep them open in another tab.*

2.2 1. Initial setup and pre-processing of metagenomic reads

2.2.1 Log into your instance

If you weren't in the beginner microbiome analysis workshop and need a refresher, check the notes in Beginner Module 1.

2.2.2 Working on the Command Line

Make a new folder in your workspace to work from:

```
cd workspace
mkdir amb_module1
cd amb_module1
```

If you weren't in the beginner workshop, feel free to take a look around the instance using the `ls`, `pwd` and `cd` commands, but make sure you come back to `~/workspace/amb_module1` before the next step.

Make a link to the files to your working directory.

```
ln -s ~/CourseData/MIC_data/AMB_data/raw_data_subsampled_subset/ raw_data
cp ~/CourseData/MIC_data/AMB_data/metagenome_metadata_subset.txt metagenome_metadata.t
```

2.2.3 About the samples

These samples are from a study that investigated gut microbiome functional capacity of chicken ceca. You can read more in the original study here: - Week-Old Chicks with High *Bacteroides* Abundance Have Increased Short-Chain Fatty Acids and Reduced Markers of Gut Inflammation

Or in the additional analysis done by John's lab here: - Metabolic modeling of microbial communities in the chicken ceca reveals a landscape of competition and co-operation

2.2.4 A Crash Course in GNU Parallel

First, activate our conda environment for this tutorial:

```
conda activate amb_module1
```

Sometimes in bioinformatics, the number of tasks you have to complete can get VERY large (e.g. when we have thousands of samples). Fortunately, there are

2.2. 1. INITIAL SETUP AND PRE-PROCESSING OF METAGENOMIC READS13

several tools that can help us with this. One such tool is GNU Parallel. This tool can simplify the way in which we approach large tasks, and as the name suggests, it can iterate through many tasks in *parallel*, i.e. at the same time. If you aren't familiar with using `top` or `htop` to look at the number of processes that you have available to you, take a look at this section of the beginner workshop.

We can use a simple command to demonstrate how to use `parallel`:

```
parallel 'echo {}' ::: a b c
```

With the command above, the program contained within the quotation marks ' ' is `echo`. This program is run 3 times, as there are 3 inputs listed after the `:::` characters. What happens if there are multiple lists of inputs? Try the following:

```
parallel 'echo {}' ::: a b c ::: 1 2 3
```

Here, we have demonstrated how `parallel` treats multiple inputs. It uses all combinations of one of each from `a b c` and `1 2 3`. But, what if we wanted to use 2 inputs that were sorted in a specific order? This is where the `--link` flag becomes particularly useful. Try the following:

```
parallel --link 'echo {}' ::: a b c ::: 1 2 3
```

In this case, the inputs are “linked”, such that only one of each is used. If the lists are different lengths, `parallel` will go back to the beginning of the shortest list and continue to use it until the longest list is completed. You do not have to run the following command, as the output is provided to demonstrate this.

```
parallel --link 'echo {}' ::: light dark ::: red blue green
```

Notice how `light` appears a second time (on the third line of the output) to satisfy the length of the second list.

Another useful feature is specifying *which* inputs we give `parallel` are to go *where*. This can be done intuitively by using multiple brackets `{ }` containing numbers corresponding to the list we are interested in. Again, you do not have to run the following command, as the output is provided to demonstrate this.

```
parallel --link 'echo {1} {3}; echo {2} {3}' ::: one red ::: two blue ::: fish
```

Finally, a handy feature is that `parallel` accepts files as inputs. This is done slightly differently than before, as we need to use four colon characters `:::` instead of three. `Parallel` will then read each line of the file and treat its contents as a list. You can also mix this with the three-colon character lists `:::` you are already familiar with. Using the following code, create a test file and use `parallel` to run the `echo` program:

```
echo -e "A\nB\nC" > test.txt
parallel --link 'echo {2} {1}' :::: test.txt ::: 1 2 3
```

Take a look inside `test.txt` with the `less` command if you like. Remember that you can use `q` to exit the file again.

And with that, you're ready to use parallel for all of your bioinformatic needs! We will continue to use it throughout this tutorial and show some additional features along the way. There is also a cheat-sheet here for quick reference.

2.3 Quality Control

Quality control is an important step in any pipeline. For this tutorial, we will use FastQC to inspect the quality of our samples.

2.3.1 Visualization with FastQC

First, make your desired output directory `fastqc_out`. Then, run FastQC as follows:

```
fastqc -t 4 raw_data/*fastq.gz -o fastqc_out
```

Run MultiQC on the FastQC results:

```
multiqc fastqc_out --filename multiqc.html
```

Remember that you can go to <http://##.uhn-hpc.ca/> (where ## is your personal number) to find all of the files that we're creating.

Copy across the `multiqc.html` file to your computer.

Click on the `.html` files to view the results for each sample. Let's look at the Per base sequence quality tab. This shows a boxplot representing the quality of the base call for each position of the read. In general, due to the inherent degradation of quality with increased sequence length, the quality scores will trend downward as the read gets longer. However, you may notice that for our samples this is not the case! This is because for the purpose of this tutorial, your raw data has already been trimmed.

More often, per base sequence quality will look like the following. The FastQC documentation provides examples of "good" and "bad" data. These examples are also shown below:

Now, you may have also noticed that most of the samples fail the "Per Base Sequence Content" module of FastQC.

This specific module plots out the proportion of each base position in a file, and raises a warning/error if there are large discrepancies in base proportions. In a given sequence, the lines for each base should run in parallel, indicating that the base calling represents proper nucleotide pairing. Additionally, the A and T lines may appear separate from the G and C lines, which is a consequence of the GC content of the sample. The relative amount of each base reflects the overall amount of the bases in the genome, and should not be imbalanced. When this is not the case, the sequence composition is **biased**. A common cause of this is the use of primers, which throws off our sequence content at the beginning of the read. Fortunately, although the module error stems from this bias, according

to the FastQC documentation for this module it will not largely impact our downstream analysis.

The other modules are explained in depth in the FastQC Module Help Documentation

2.3.2 Filtering with KneadData

KneadData is a tool which “wraps” several programs to create a pipeline that can be executed with one command. For this tutorial though, we will use KneadData to filter our reads for contaminant sequences against a human database. KneadData will: * Run Trimmomatic to remove adapter sequences * Run Bowtie2 to remove reads mapping to the reference genome and the PhiX genome (commonly used as a sequencing control)

We don’t have paired-end data here, but if we did, it would also: * Check whether both pairs of a read exist * Check whether they both map to the reference genome

Bowtie2 needs a reference genome/index file for its mapping step. There are some pre-made indexes on this page.

MAKE A CHOICE!

You can now choose to use the database that we have already made (quicker option) or you can see the steps involved in building this for yourself (slower option).

Quicker option - use pre-made database

```
ln -s ~/CourseData/MIC_data/amb_module1/bowtie2_db/ .
```

Slower option - build the database yourself

First, download the NCBI RefSeq assembly summary:

```
wget ftp://ftp.ncbi.nlm.nih.gov/genomes/refseq/assembly_summary_refseq.txt
```

Have a look in this file. You’ll see that it contains information on all of the genomes in NCBI RefSeq - we often need reference genomes for constructing various kinds of databases, so it’s useful to get familiar with what is available on NCBI. You’ll see that each assembly/accession has information on the NCBI taxonomy ID, whether it’s a complete genome, whether it’s a reference or representative genome, and a link to where we can find the associated files. Because the samples we are using are from chickens, we’re going to look for the files that are associated with the chicken *Gallus gallus*:

```
grep "Gallus gallus" assembly_summary_refseq.txt
```

The first one has “reference genome” in the first column, so that’s what we’ll use. Copy the link that starts `https://ftp.ncbi....` into a new browser window. You will see a bunch of different file types, so find the one that ends `_genomic.fna.gz`, right click, and click “Copy Link Address”. Now use the `wget` command to download this to your instance.

Now we’re going to do the same with PhiX. PhiX is added as a sequencing control and is often removed before we’re given back our samples, but it’s good to make sure. - Use `grep` to search for “phiX174” - Go to the link that you find - Download the `_genomic.fna.gz` file using `wget`

QUESTION!

Question 1: This one isn’t really a question, but you can see the code needed for these steps in the answers section if you need to!

Now we’ll combine these two files into one:

```
cat GCF_016699485.2_bGalGal1.mat.broiler.GRCg7b_genomic.fna.gz GCF_000819615.1_ViralPro
```

And finally, we’ll build our Bowtie2 database:

```
mkdir bowtie2_db
bowtie2-build chicken_Gallusgallus_phiX174.fna.gz bowtie2_db/chicken_Gallusgallus_phiX
```

This will take about ten minutes to run. If this is taking too long, feel free to stop this command using `ctrl+c` and link to the one that we already made:

```
ln -s ~/CourseData/MIC_data/amb_module1/bowtie2_db/ .
```

2.3.2.1 Back to KneadData

Now we are ready to run KneadData using parallel:

```
parallel -j 1 --eta 'kneaddata -un {1} -o kneaddata_out -db bowtie2_db/chicken_Gallusg
```

You’ll see that after the first sample has finished, you’ll get an estimate of the amount of time remaining.

While `kneaddata` is running, consider the following:

Note that if we had paired reads to start with, we’d also have some files containing reads that were unmatched between the forward and reverse files.

We’ll also run the following to get counts of how many reads we had in each file after each step:

```
kneaddata_read_count_table --input kneaddata_out --output kneaddata_read_counts.txt
```

QUESTION!

Question 2: Take a look at this file. Are there any surprises? Which of the output files in 'kneaddata_out' will you use for analysis?

Question 3: How many reads are in each sample before and after KneadData?

We'll move the other output files that we're not interested into a new folder:

```
mkdir kneaddata_out/contam_seq
mv kneaddata_out/*_contam*.fastq kneaddata_out/contam_seq
```

2.4 2. Generating Taxonomic Profiles

2.4.1 Annotation with Kraken2/Bracken

Now that we have our reads of interest, we want to understand *what* these reads are. To accomplish this, we use tools which **annotate** the reads based on different methods and databases. There are many tools which are capable of this, with varying degrees of speed and precision. For this tutorial, we will be using Kraken2 for fast exact k-mer matching against a database.

We have also investigated which parameters impact tool performance in this Microbial Genomics paper. One of the most important factors is the contents of the database, which should include as many taxa as possible to avoid the reads being assigned an incorrect taxonomic label. Generally, the bigger and more inclusive database, the better. However, due to the constraints of our AWS cloud instances, we will be using a "PlusPF 8GB" index provided by the Kraken2 developers.

First of all, we want to download and extract the Kraken database:

```
wget https://genome-idx.s3.amazonaws.com/kraken/k2_pluspf_08gb_20250402.tar.gz
mkdir k2_pluspf_08gb
tar -xvf k2_pluspf_08gb_20250402.tar.gz -C k2_pluspf_08gb
```

Then we can clean-up by removing the original file:

```
rm k2_pluspf_08gb_20250402.tar.gz
```

NOTE!

First, you must create the appropriate output directories, or Kraken2 will not write any files. Use the 'mkdir' command to make the directories to match what we're using below. Using 'parallel', we will then run Kraken2 for our concatenated reads.

Note that it's always a good idea to first try out a `--dry-run` of parallel before you run any long jobs. If you're satisfied with what it's going to be running, remove the `--dry-run` flag.

```
parallel -j 1 --eta --dry-run 'kraken2 --db k2_pluspf_08gb --output kraken2_outraw/{/.}
```

This process can take some time. While this runs, let's learn about what our command is doing! * We first specify our options for parallel, where: * the `-j 2` option specifies that we want to run two jobs concurrently; * the `--eta` option will count down the jobs as they are completed; * after the program contained in quotation marks, we specify our input files with `:::`, and use a regex to match all of the concatenated, unzipped `.fastq` files. * We then describe how we want **kraken** to run: * by first specifying the location of the database with the `--db` option; * then specifying the output directory for the raw kraken annotated reads; * notice that we use a special form of the brackets here, `{/.}`, this is a special function of **parallel** that will remove both the file path and extension when substituting the input into our **kraken** command. This is useful when files are going into different directories, and when we want to change the extension. * similarly, we also specify the output of our "report" files with the `--report` option; * the `-confidence` option allows us to filter annotations below a certain threshold (a float between 0 and 1) into the unclassified node. We are using 0 because our samples are already subset, however this should generally be higher. See our paper for more information. * and finally, we use the empty brackets `{}` for **parallel** to tell **kraken** what our desired input file is.

As Kraken runs, you should see it printing out a summary of the number of reads within each sample that it was able to classify.

With Kraken2, we have annotated the reads in our sample with taxonomy information. If we want to use this to investigate diversity metrics, we need to find the abundances of taxa in our samples. This is done with Kraken2's companion tool, Bracken (Bayesian Reestimation of Abundance with KrakEN).

Let's run Bracken on our Kraken2 outputs! **First, make the expected output directory**, then run the following:

```
parallel -j 2 --eta 'bracken -d k2_pluspf_08gb -i {} -o bracken_out/{/.}.species.bracken
```

Some notes about this command: * `-d` specifies the database we want to use. It should be the same database we used when we ran Kraken2; * `-i` is our input file(s); * `-o` is where and what we want the output files to be; * `-r` is the read length to get all classifications for, the default is 100 but our reads are 150bp; * `-l` is the taxonomic level at which we want to estimate abundances; * `-t` is the number of reads required prior to abundance estimation to perform re-estimation

Finally, let's merge our bracken outputs:

```
combine_bracken_outputs.py --files bracken_out/*.species.bracken -o merged_output.spec
```

2.4.2 Annotation with MetaPhlAn

Another tool that is commonly used for taxonomic annotation of metagenomic sequences is MetaPhlAn. This tool is different from Kraken2 in that it uses a

database of marker genes, instead of a collection of genomes, and it identifies only these marker genes within our reads, rather than trying to classify all reads. It then attempts to estimate the abundance of the taxa it identified within our whole samples, but it's important to remember that this is an **estimation**, and not the actual number of reads classified. We will use MetaPhlAn 3 for this tutorial due to the constraints of the AWS instances, but a newer version (MetaPhlAn4) utilizing a larger database is available.

Usually, we would build this ourself so that we make sure to have the most recent version. This can take quite a while, so we will create a link to it:

```
ln -s ~/CourseData/MIC_data/ref_dbs/CHOCOPhlan_201901 .
```

And then activate the environment:

```
conda activate metaphlan3
```

Now, let's make an output folder `metaphlan_out`. Then, we can run the following:

```
parallel -j 1 --eta 'metaphlan --bowtie2db CHOCOPhlan_201901 --input_type fastq -o metaphlan_out/
```

This might take a while to run, so move on to the next step while this runs and then come back to this after.

2.5 3. Visualisation of results in R

Now we'll go to RStudio server: <http://##.uhn-hpc.ca:8080> (remember to replace ## with your own instance number). If you weren't in the beginner microbiome workshop, you may be prompted to enter your username and password. It will also take a little while to log in the first time, but will be quicker afterwards.

2.5.1 Create the R Notebook

Using the menus, click **File > New File > R Notebook**, which will open an Untitled R markdown (Rmd) document. R notebooks are helpful in that you can run several lines, or **chunks** of code at the same time, and the results will appear within the document itself (*in the whitespace following the chunk*).

The default R Notebook contains a header and some information you may find helpful. Try running the chunk containing `plot(cars)` to see what happens!

You do not need to preserve most of the information in the new, untitled document. Select all of its contents by click+dragging your cursor or entering the **ctrl+a** (Windows) / **cmd+a** (Mac) shortcut, and press **backspace** or **delete** to clear the document.

The **chunks** are distinguished by the grey shading. Everything between the first ````{r}` and subsequent ````` belongs to the chunk. Anything written in the

white space surrounding the chunk is meant to be annotation. *Although you can run lines of code outside of the chunks, the chunks are useful for running multiple lines in series with one click.*

2.5.1.1 Adding new chunks

To add a new chunk into your R notebook by navigating to **Code > Insert Chunk** from the toolbar, or clicking on the little green C and selecting R.

And then save the document - it doesn't really matter what you call it, but something like `amb_module1.rmd` would be sensible.

2.5.2 Setup, Importing and Formatting Data

Now we'll start building our R markdown notebook. Paste the following into a chunk, and then click the little green "play" button on the top right of the chunk to run it.

```
library(phyloseq)
library(vegan)
library(ggplot2)
library(gridExtra)
```

```
setwd("~/workspace/amb_module1/")
```

In this chunk, we've imported the libraries/packages that we're going to use (phyloseq, vegan, ggplot2 and gridExtra) and we've told R which directory it should be working from.

Now we're going to read in the Bracken output, and do some formatting of it:

```
bracken <- read.csv("merged_output.species.bracken", sep="\t") #read in the file as a c
rownames(bracken) = bracken[,1] #set the row names to those in the first column
cols = c() #make an empty list/vector
for (colname in colnames(bracken)) { #for each of the columns in bracken
  if (grepl("bracken_num", colname, fixed = TRUE)) { #if "bracken_num" is in the column
    cols <- c(cols, colname) #add it to the columns list - we're going to use this to
  }
}
bracken = bracken[,cols] #subset the dataframe to only the columns in the cols list (i
colnames(bracken) = sapply(strsplit(colnames(bracken), "_"), `[, 1) #rename the column
```

If you're not familiar with R or Python, anything after the `#` can be used for making comments, as it won't be read by them. So you can see what I've written about each step.

It's good practice to take a look at what we're changing each time! In the "Environment" section of RStudio, you should be able to see all of these objects. Click on them to have a look. I won't always tell you to click on them, but it's a

good idea to either click on them or print them out so that you can understand what is changing each time. There are some other ways to look at this too:

```
View(bracken)
print(bracken)
bracken
```

Some of these options are more or less appropriate depending on what you're doing, but regularly printing things out is the easiest way to troubleshoot a script that isn't working. If you were wanting to print something out within a loop, you'd need to use the `print()` function.

Now we're going to get the genus name of each of the species in our dataframe:

```
all_genus = c() #make an empty list/vector
for (species in rownames(bracken)) { #for each of the rows (species) in the bracken dataframe
  this_species = species
  if (grepl("uncultured", this_species, fixed = TRUE)) { #if "uncultured" is in the name
    this_species = gsub("uncultured ", "", this_species) #rename the species by removing the "un
  }
  this_genus = strsplit(this_species, " ")[[1]][1] #get the first word of the species name (the g
  all_genus = c(all_genus, this_genus) #add this genus name to our list
}
tax_table = data.frame(genus=all_genus, species=rownames(bracken)) #make a new table with the gen
rownames(tax_table) = rownames(bracken) #and add the rownames to it
```

Next, we'll read in our metadata and reformat it:

```
metadata = read.csv("metagenome_metadata.txt", sep="\t") #read in the file as a dataframe
samples = data.frame(metadata[,2], stringsAsFactors = FALSE) #get only the second column
rownames(samples) = metadata[,1] #add the sample names
colnames(samples) = c('group') #change the column name
samples['classifier'] = 'kraken' #add a new column that shows us that these samples were classifi
```

And now we'll combine these objects to a phyloseq object:

```
BRACKEN = otu_table(bracken, taxa_are_rows = TRUE) #convert bracken to an otu_table
SAMPLE = sample_data(samples) #convert the samples dataframe to a phyloseq sample_data format
TAX = tax_table(tax_table) #convert the tax_table dataframe to a phyloseq tax_table format
taxa_names(TAX) <- rownames(bracken) #get the taxa names from the bracken row names
physeq_bracken = phyloseq(BRACKEN, TAX, SAMPLE) #make the phyloseq object
```

Print out `physeq_bracken` to see if it has all of the taxa and samples that you are expecting.

And collapse the table at the genus level:

```
physeq_bracken_genus = tax_glom(physeq_bracken) #collapse at the genus level (the default is the
taxa_names(physeq_bracken_genus) = data.frame(tax_table(physeq_bracken_genus))$ta1 #add the genus
```

Have a look at it:

```
View(physeq_bracken_genus)
```

2.5.3 Filtering and rarefying

Now we can take a look at how many reads have been classified within our samples:

```
sample_sums(physeq_bracken_genus)
print(c("Minimum sample size:", min(sample_sums(physeq_bracken_genus)), "Maximum sample
```

QUESTION!

Question 4: Roughly what proportion of the reads you have in the samples after KneadData have been classified? Is this surprising to you? Remember that you can look in the ‘kneaddata_read_counts.txt’ file to see how many reads the files all ended up with!

It’s also always good to have a look at a rarefaction curve. This tells us about whether we’ve sequenced our samples sufficiently to capture all of the diversity within them:

```
rarecurve(as.data.frame(t(otu_table(physeq_bracken_genus))), step=50, cex=0.5, label=TRUE)
```

QUESTION!

Question 5: How do these look to you? If these were your samples, would you be happy with these?

Let’s also take a look at the number of reads assigned to each taxon within our samples:

```
length(taxa_sums(physeq_bracken_genus))
print(c("Minimum reads:", min(taxa_sums(physeq_bracken_genus)), "Maximum reads:", max(taxa_sums(physeq_bracken_genus))
```

You’ll see that we have a minimum of 1 read assigned to a genus, and a maximum of >14M reads. Usually, we’ll remove some of the very low abundance taxa in our samples.

```
physeq_bracken_genus<-prune_taxa(taxa_sums(physeq_bracken_genus)>=10, physeq_bracken_genus)
```

If you take a look at this, you’ll see that we’ve already gone from >1500 genera to only 490. But we can reduce this a little further:

```
physeq_bracken_genus<-prune_taxa(taxa_sums(physeq_bracken_genus)>=100, physeq_bracken_genus)
```

Now we only have 112 genera, and with a minimum of a sum of 100 reads (average 10 per sample), this seems fairly reasonable for filtering.

Let’s take a look at the rarefaction curves again:

```
rarecurve(as.data.frame(t(otu_table(physeq_bracken_genus))), step=50, cex=0.5, label=TRUE, ylim = c(0, 1000))
```

QUESTION!

Question 6: How do you feel about the samples now?

Next we'll rarefy the samples prior to our diversity calculations - the number of classified reads varies a lot between our different samples!

```
set.seed(711)
rarefied_bracken_genus<-rarefy_even_depth(physeq_bracken_genus, rngseed = FALSE, sample.size = min(sample_sums(physeq_bracken_genus)))
```

See that the samples are all the same size now:

```
sample_sums(rarefied_bracken_genus)
```

NOTE!

We are only going to use this rarefied object for the alpha diversity analyses!

2.5.4 Alpha and beta diversity analyses

We're not going to go into explaining all of these again as we did that in the beginner workshop, but we'll take a look at the samples with a few different methods.

First, we'll look at the alpha diversity:

```
plot_richness(physeq = rarefied_bracken_genus, x="group", color = "group", measures = c("Observed", "Shannon", "Simpson"))
```

QUESTION!

Question 7: What do you notice about the alpha diversity differences between the samples?

Now we'll convert the non-rarefied phyloseq object to relative abundance:

```
relabun_bracken_genus <- transform_sample_counts(physeq_bracken_genus, function(x) x*100 / sum(x))
```

And then we'll use this for a Bray-Curtis PCoA plot:

```
ordination<-ordinate(relabun_bracken_genus, "PCoA", "bray")
plot_ordination(relabun_bracken_genus, ordination, color = "group") +
  geom_point(size=10, alpha=0.1) + geom_point(size=5) + stat_ellipse(type = "t", linetype = 2) +
  theme_minimal()
```

QUESTION!

Question 8: do the samples group how you expected them to based on their alpha diversity?

Next we'll make a bar plot. Plotting all 112 of the taxa on a barplot is a bit useless, so first we'll filter to include only the 20 most abundant genera:

```
relabun_bracken_genus_top<-prune_taxa(names(sort(taxa_sums(relabun_bracken_genus),decr  
plot_bar(relabun_bracken_genus_top, fill="ta1") + facet_wrap(c(~group), scales="free_x
```

QUESTION!

Question 9: What do you notice about the composition of these samples? What are the samples from the Non-Bacteroides group dominated by?

Often heatmaps are a little easier to look at than a stacked barplot, and it's easier to see differences between samples:

```
plot_heatmap(relabun_bracken_genus_top, method = "PCoA", distance = "bray", low = "whi
```

This function is grouping our samples by their Bray-Curtis distance from one-another, although note that this may not look quite the same as it would in the PCoA plot because we're only looking at the 20 most abundant genera.

That was the plotting inside phyloseq, but sometimes we also want a dendrogram to show how our samples cluster together, so let's also try something else:

```
heatmap(otu_table(relabun_bracken_genus_top))
```

Note that this has just clustered our taxa together, and the tree that we see for the genera is not a phylogenetic tree!

QUESTION!

Question 10: What can you notice about our samples now? Particularly the taxonomic composition of the Non-Bacteroides samples?

2.5.5 Now lets do the same with the MetaPhlAn3 output

MetaPhlAn 3 should have finished running by now! So go back to your terminal window.

You may see a warning come out, but as long as it ran, that doesn't matter!

Similar to Kraken2, MetaPhlAn will output individual files for each sample. We can use a utility script from MetaPhlAn to merge our outputs into one table.

```
merge_metaphlan_tables.py metaphlan_out/*.mpa > metaphlan_merged_out.txt
```


If we view this new combined table, we will see three key things: 1. First, the output format is different to that of Kraken2, where the full taxonomic lineages are expanded line by line. 2. Second, MetaPhlAn only outputs relative abundances for each taxonomic node, whereas Kraken2 (before re-analysis with Bracken) will output absolute numbers of reads assigned to each taxonomic node. 3. Third, the number of taxa that MetaPhlAn finds is *much* smaller than Kraken2. This is partially due to us using a low confidence threshold with Kraken2, but this discrepancy between the two tools tends to hold true at higher confidence thresholds as well. See our paper for more info about how these tools perform compared to each other.

We're going to perform some similar steps on this MetaPhlAn output as we did on the Bracken output, but we'll be able to just keep the lines of the file that are at the genus level, rather than grouping by genus. I haven't explained things much here as you can see more explanation above, so see if you can work out what each of the lines of code is doing as you go.

Read in the file and filter to only the relevant lines:

```
metaphlan <- read.csv("metaphlan_merged_out.txt", sep="\t", skip=1)
rownames(metaphlan) = metaphlan[,1]

keeping = c()
genus = c()
for (row in rownames(metaphlan)) {
  if (grepl("g__", row, fixed = TRUE)) {
    if (grepl("s__", row, fixed = TRUE)) {
      skip_this = TRUE
    } else {
      keeping = c(keeping, row)
      genus = c(genus, strsplit(row, "__")[[1]][7])
    }
  }
}

metaphlan = metaphlan[keeping,]
metaphlan = metaphlan[,c(-1,-2)]
rownames(metaphlan) = genus
colnames(metaphlan) = sapply(strsplit(colnames(metaphlan), "_"), `[, 1)

Read in the metadata again:

metadata = read.csv("metagenome_metadata.txt", sep="\t")
samples = data.frame(metadata[,2], stringsAsFactors = FALSE)
rownames(samples) = metadata[,1]
colnames(samples) = c('group')
samples['classifier'] = 'metaphlan'
```

Make a taxonomy table:

```
tax_table = data.frame(genus=rownames(metaphlan))
rownames(tax_table) = rownames(metaphlan)
```

Combine these into a phyloseq object:

```
METAPHLAN = otu_table(metaphlan, taxa_are_rows = TRUE) #convert asv_table_num to an otu
SAMPLE = sample_data(samples)
TAX = tax_table(tax_table)
taxa_names(TAX) <- rownames(metaphlan)
physeq_metaphlan_genus = phyloseq(METAPHLAN, TAX, SAMPLE)
```

We're going to skip some of the first steps looking at read depths and the number of reads assigned to different species, as we've identified marker genes using MetaPhlAn3 and not reads, and this is relative abundance and not a count.

There are only two alpha diversity metrics that we can look at without counts:

```
plot_richness(physeq = physeq_metaphlan_genus, x="group", color = "group", measures = c("richness", "shannon", "simpson"))
```

Make a PCoA plot with Bray-Curtis distance:

```
ordination<-ordinate(physeq_metaphlan_genus, "PCoA", "bray")
plot_ordination(physeq_metaphlan_genus, ordination, color = "group") +
  geom_point(size=10, alpha=0.1) + geom_point(size=5) + stat_ellipse(type = "t", linetype="dashed")
```

QUESTION!

Question 11: Are these results for alpha and beta diversity as you expect them to be from the Kraken/Bracken results?

We don't have very many taxa in our MetaPhlAn3 output, so we won't bother filtering this before plotting the stacked barplot:

```
plot_bar(physeq_metaphlan_genus, fill="ta1") + facet_wrap(c(~group), scales="free_x", ncol=1)
```

Then make the first of the heatmaps again:

```
plot_heatmap(physeq_metaphlan_genus, method = "PCoA", distance = "bray", low = "white", high = "black")
```

You'll probably notice a lot of grey spots here - that means that genus wasn't found in the sample at all.

And the second one:

```
heatmap(otu_table(physeq_metaphlan_genus))
```

QUESTION!

Question 12: What do you notice here compared with the Kraken/Bracken results?

2.5.6 Combine Kraken and MetaPhlAn output

First we'll rename both of the Bracken and MetaPhlAn phyloseq objects so that we don't have duplicated sample names in them:

```
sample_names(relabun_bracken_genus) = paste(sample_names(relabun_bracken_genus), "_bracken", sep="")
sample_names(physeq_metaphlan_genus) = paste(sample_names(physeq_metaphlan_genus), "_metaphlan", sep="")
```

Hopefully you can see how we're just pasting together each of the sample names with either “_bracken” or “_metaphlan”

Now join the two phyloseq objects:

```
physeq_combined = merge_phyloseq(relabun_bracken_genus, physeq_metaphlan_genus)
```

Make two PCoA plots, one with Jaccard distance (looking at presence/absence only) and then one with Bray-Curtis dissimilarity (that also accounts for the abundance of each genus):

```
ordination<-ordinate(physeq_combined, "PCoA", "binary")
plot1 <- plot_ordination(physeq_combined, ordination, color = "group", shape = "classifier") +
  geom_point(size=10, alpha=0.1) + geom_point(size=5) + stat_ellipse(type = "t", linetype = 2) +
ordination<-ordinate(physeq_combined, "PCoA", "bray")
plot2 <- plot_ordination(physeq_combined, ordination, color = "group", shape = "classifier") +
  geom_point(size=10, alpha=0.1) + geom_point(size=5) + stat_ellipse(type = "t", linetype = 2) +
grid.arrange(plot1, plot2, ncol=2)
```

Does this look weird? Click the little button on the right between the code and the plot to “Show in new window”.

QUESTION!

Question 13: Anything surprising here? What do you notice about the two plots?

We had included all of the taxa previously, but Bracken has a lot more genera identified than MetaPhlAn. So we'll filter the Bracken output to be a little more similar:

```
relabun_bracken_genus_top<-prune_taxa(names(sort(taxa_sums(relabun_bracken_genus),decreasing=TRUE)),
physeq_combined = merge_phyloseq(relabun_bracken_genus_top, physeq_metaphlan_genus)
```

Now our PCoA plot looks a little nicer:

```
ordination<-ordinate(physeq_combined, "PCoA", "bray")
plot_ordination(physeq_combined, ordination, color = "group", shape = "classifier") +
  geom_point(size=10, alpha=0.1) + geom_point(size=5) + stat_ellipse(type = "t", linetype = 2) +
```

Let's have a look at the taxa in the samples:

```
plot_bar(physeq_combined, fill="ta1") + facet_wrap(c(~"group", "classifier"), scales="free_x", nr=
```

Can you see all of the taxa? Try filtering to only the top 20, modifying the code that we used above.

And look at the heatmap:

```
heatmap(otu_table(physeq_combined))
```

QUESTION!

Question 14: What differences do you notice between the two classifiers?

2.6 Answers

Question 1: This one isn't really a question, but you can see the code needed for these steps in the answers section if you need to!

```
grep "phiX174" assembly_summary_refseq.txt
#go to https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/819/615/GCF_000819615.1_ViralPro
wget https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/819/615/GCF_000819615.1_ViralPro
```

Question 2: Take a look at this file. Are there any surprises? Which of the output files in `kneaddata_out` will you use for analysis? Nope! There are no samples where we seem to have lost a really large number of reads at any step - we don't lose many reads that couldn't be trimmed (difference between `raw single` and `trimmed single`), and not too many of the reads are host-associated (difference between `trimmed single` and `decontaminated chicken_Gallusgallus_phiX174 single`). We'll use the `sample_kneaddata.fastq` files for analysis, and not the `sample_kneaddata_chicken_Gallusgallus_phiX174_bowtie2_contam.fastq` files.

Question 3: How many reads are in each sample? The original number of reads per sample ranges between 931,686 and 1,017,130 (average 964,187) and after KneadData ranges between 912,712 and 998,264 (average 950,470).

Question 4: Roughly what proportion of the reads you have in the samples after KneadData have been classified? Is this surprising to you? We have between 52,405 (`min(sample_sums(physeq_bracken_genus))`) and 446,368 (`max(sample_sums(physeq_bracken_genus))`) - average 221,440 (`mean(sample_sums(physeq_bracken_genus))`). $221,440/950,470 = \sim 23\%$ of reads classified. This does seem relatively low, although it's not uncommon for only a small proportion of the reads in samples not from humans (i.e. from less well-characterised environments) to be classified. In this case, it's likely mainly due to the really small database that we've used for taxonomic classification. We've used an 8GB database due to the constraints of these AWS instances, but we recommend using the largest database that you possibly can - the ones that we typically use in our lab are 800GB-1.1TB in size.

Question 5: How do these look to you? If these were your samples, would you be happy with these? These curves don't look great - they don't saturate at all.

Question 6: How do you feel about the samples now? This is much better - and we can see that our sample rarefaction curves not being saturated is down to these very low abundance taxa. While some of them could be important, anything present at abundances this low could be due to false-positive noise, and we didn't do any filtering with Bracken.

Question 7: What do you notice about the alpha diversity differences between the samples? There are huge differences between our *Bacteroides* and Non-*Bacteroides* sample groups. The samples in the *Bacteroides* group have much lower diversity than the Non-*Bacteroides* samples.

Question 8: do the samples group how you expected them to based on their alpha diversity? Yes. All of the *Bacteroides* samples group away from the Non-*Bacteroides* samples, which isn't surprising based on how different their alpha diversity is.

Question 9: What do you notice about the composition of these samples? What are the samples from the Non-*Bacteroides* group dominated by? The *Bacteroides* samples have very low abundance of anything other than *Bacteroides*, while the Non-*Bacteroides* samples have many other genera also present. There isn't really any one genus that dominates the Non-*Bacteroides* samples, although a couple do seem to have high *Escherichia*.

Question 10: What can you notice about our samples now? Particularly the taxonomic composition of the Non-*Bacteroides* samples? There is very little of anything but *Bacteroides* in the *Bacteroides* samples. The Non-*Bacteroides* are not particularly similar to one-another, with a range of taxa present in mid-high abundances in each of them.

Question 11: Are these results for alpha and beta diversity as you expect them to be from the Kraken/Bracken results? Yes - the *Bacteroides* and Non-*Bacteroides* samples are still separating from one another and the alpha diversity is much lower in *Bacteroides* than Non-*Bacteroides* samples.

Question 12: What do you notice here compared with the Kraken/Bracken results? These samples are very sparse! Lots of the taxa are only identified in one or a few samples. This is one of the draw-backs of the marker-gene approach used by MetaPhlAn - we can typically be very confident in the taxa that it identifies as present, but it is likely to miss lots of taxa, too (low recall, high precision).

Question 13: Anything surprising here? What do you notice about the two plots? Jaccard distance only takes into account presence/absence of taxa, while Bray-Curtis also accounts for abundance. It's surprising how similar the two groups (*Bacteroides*/Non-*Bacteroides*) look with Bray-Curtis when they are so different with Jaccard, but this means that the most abundant taxa must be similar between the two different classification methods.

Question 14: What differences do you notice between the two classifiers? The Bacteroides samples look very similar between the two methods, but the Non-Bacteroides samples seem to be more different. Escherichia seems to dominate three MetaPhlAn samples but only two Kraken samples. There are lots of other differences too! See if you can come up with some explanations for these differences.

Part III

Modules

Chapter 3

Module 2: Metagenomic Assembly and Binning

This tutorial is part of the 2025 CBW Advanced Microbiome Analysis (held in Vancouver, BC, May 29-30). It is based on a version written for the 2024 CBW Advanced Microbiome Analysis workshop by Ben Fisher.

Author: Robyn Wright, updated by Ryan J. McLaughlin (2025)

3.1 Table of Contents

- Introduction
- 2.1 Initial setup
- 2.2 Assembly of raw reads with MEGAHIT
- 2.3 Make an Anvi'o contigs databases
- 2.4 Run HMMs to identify genes
- 2.5 Map samples onto contigs with Bowtie2
- 2.6 Cluster contigs into bins
- 2.7 Interactive refinement of bins
- 2.8 Export the final bins into MAGs
- 2.9 Run CheckM
- Other things that we do with MAGs
- Answers

3.2 Introduction

The main goal of this tutorial is to introduce students to the assembly of genomes from metagenomic reads (Metagenome Assembled Genomes/MAGs). There is not a one-size-fits-all pipeline for assembling MAGs. MAG assembly is incredibly

computationally intensive, and the approach that we are taking here is therefore a very “bare bones” approach designed to demonstrate the main steps involved and give you some familiarity with the methods used. At the end of this tutorial we’ve provided a few other pipelines for MAG assembly that you may wish to look into if you are looking to assemble MAGs with your own metagenome data.

Throughout this module, there are some questions aimed to help your understanding of some of the key concepts. You’ll find the answers at the bottom of this page, but no one will be marking them.

3.2.1 Anvi’o

Anvi’o is an open-source, community-driven analysis and visualization platform for microbial ‘omics. It packages together many different tools used for genomics, metagenomics, metatranscriptomics, phylogenomics, etc. and has great interactive visualisations that can be used to help this. We are just touching the surface of what Anvi’o can do today, but the website has great tutorials and learning resources for all of its capabilities - I recommend browsing through to get some inspiration!

3.3 2.1 Initial setup

First of all, we need to create the directory that we’ll be working from and change into it:

```
cd ~/workspace/
mkdir amb_module2
cd amb_module2
```

Then we’ll create a symlink to the data that we’re going to be using, and copy in the sample names and metadata:

```
ln -s ~/CourseData/MIC_data/AMB_data/reads_for_assembly/ .
ln -s ~/CourseData/MIC_data/AMB_data/fastqs_for_assembly .
cp ~/CourseData/MIC_data/AMB_data/sample_ids.txt .
ls
```

Finally, we’ll activate the conda environment that we’re going to be using:

```
mamba activate anvio-8
anvi-interactive --version
```

NOTE: This lab has several steps that are quite time consuming and resource intensive!!! We have included the commands for your knowledge, but they are will commented out and alternative linking or copying commands will be provided. Read carefully and copy-paste only after you understand what the step you are performing should be doing.

3.4 2.2 Assembly of raw reads with MEGAHIT

The first step in the assembly of MAGs is the assembly of the shorter reads from sequencing into contigs. A contig is a contiguous sequence assembled from a set of sequence fragments, and there are different bioinformatic methods for assembling our reads into contigs. For this part of the tutorial, we are continuing to use reads from the same samples that we used for taxonomic annotation, but we have sub-sampled these to contain reads from only a few species, to ensure that we have the read depth required for MAG assembly while keeping the files small enough to be able to run on these small AWS instances. These are all still short reads of approximately 150 base pairs each.

We are using a tool called **MEGAHIT** for this assembly because it allows co-assembly of the reads from multiple samples. This means that rather than the assembly being performed separately for each sample, it is performed on all of the reads from all of the samples at the same time. In this tutorial, we will co-assemble all of the samples together, but in your own analyses you should think about what makes the most sense. Are you expecting microbial taxa to overlap between different samples? Would it make sense to find similar genomes in multiple samples? If you answered “yes” to those questions, then it might be worth thinking about co-assembly. If not, then it is probably worth assembling each sample separately. You may also want to consider assembling the samples from different treatments separately.

Prior to assembling your samples you would usually run quality checks and remove potentially contaminating sequences, but seeing as we already did that during the Taxonomic annotation tutorial, we don’t need to repeat this step here.

First, we’ll make a directory for the output to go into:

```
mkdir anvio
```

Now we’ll run **MEGAHIT** on our samples. You may already know about a program called **tmux** (and if you were in the beginner workshop then you definitely should), but it’s really useful for running programs that may take a while, or just keeping track of where you are. There are several tools that are pre-installed on most Linux systems that we can use to make sure that our program carries on running even if we get disconnected from the server. **tmux** is one of the most frequently used. To activate it, just type in **tmux** and press enter. It should take a second to start up, and then load up with a similar looking command prompt to previously, but with a coloured bar at the bottom of the screen.

```
tmux
```

To get out of this window again, press **ctrl+b** at the same time, and then **d**. You should see your original command prompt and something like:

```
[detached (from session 0)]
```

We can actually use `tmux` to have multiple sessions, so to see a list of the active sessions, use:

```
tmux ls
```

We can rename the `tmux` session that we just created with this:

```
tmux rename-session -t 0 anvio
```

Note that we know it was session 0 because it said that we detached from session 0 when we exited it.

If we want to re-enter this window, we use:

```
tmux attach-session -t anvio
```

Now, we can run MEGAHIT inside this `tmux` session, but we will need to activate the conda environment again inside here:

```
mamba activate anvio-8
```

NOTE: I commented out the below because you DON'T want to run it today, keep reading below!

Run MEGAHIT:

```
#SE=$( ls reads_for_assembly/*.fastq.gz | tr '\n' ',' | sed 's/,,$//' )
#megahit -r ${SE} \
#      --min-contig-len 1000 \
#      --num-cpu-threads 4 \
#      --presets meta-large \
#      --memory 0.8 \
#      -o anvio/megahit_out \
#      --verbose
```

The arguments here are:

- `-r` - The single-end reads
- `--min-contig-len` - The minimum length in base pairs for contigs that we want to use. NOTE: 1000 is about the minimum length that you will ever want to use, but sometimes people might increase this to e.g. 2500 or 5000
- `--num-cpu-threads` - The number of threads to use
- `--presets` - This is for a set of parameters needed within MEGAHIT, and meta-large is suggested for large and complex metagenomes
- `--memory` - The amount of available memory that we want to allow MEGAHIT to use - this means it will use up to 80% of the memory available. Keeping this below 100% just means that we would still be able to use the Amazon instance for other things, and could be important if you're sharing a server with other people that might also need to be carrying out some work!
- `-o` - The output folder name

- `--verbose` - MEGAHIT will print out what it is doing

MEGAHIT can take quite a long time to run, so it's up to you if you want to wait for it to run. If not, press `ctrl+c` to stop the run and you can copy across the output that we already ran:

```
mkdir -p anvio/megahit_out
cp ~/CourseData/MIC_data/AMB_data/final.contigs.fa anvio/megahit_out/
```

If you ran it yourself, a step that we'll often do is removing the intermediate contigs to save space:

```
rm -r anvio/megahit_out/*/intermediate_contigs
```

The main output at this point is a fasta file containing the contigs `anvio/megahit_out/final.contigs.fa`. You can inspect the file with the `less` command if you like (remember to press `q` to exit this view).

Let's take a look at some basic stats for the assembly using a software tool called `seqkit`, but we'll add a little code after it so it's terminal friendly to view:

```
seqkit stat -T -a anvio/megahit_out/final.contigs.fa | awk -F'\t' 'NR==1 { for (i=1; i<=NF; i++)
```

Question 1: How many contigs are there in the `anvio/megahit_out/final.contigs.fa` file?

Question 2: What is the total number of basepairs?

Question 3: What is the N50?

3.5 2.3 Make an Anvi'o contigs databases

First of all, we'll run a script to reformat our final contigs file from MEGAHIT. This ensures that they're in the right format for reading into Anvi'o in the next step:

```
anvi-script-reformat-fasta anvio/megahit_out/final.contigs.fa \
    --simplify-names \
    --min-len 2500 \
    -o anvio/megahit_out/final.contigs.fixed.fa
```

Take a quick look at both `anvio/megahit_out/final.contigs.fa` and `anvio/megahit_out/final.contigs.fixed.fa` using the `less` or `head` commands to see what was changed! You can read more about the `anvi-script-reformat-fasta` command on this page. You should be able to see that we're giving this command a few options (aside from obviously giving it the `final.contigs.fa` file from the previous step:

- `--simplify-names` - this is simply telling the program to simplify the names of the contigs in the file - you should have seen that in the original `anvio/megahit_out/final.contigs.fa` file, they had a description as

well as a name, and this included information that `Anvi'o` doesn't need and may confuse it.

- `--min-len` - although we already used a minimum contig length of 1000 in the previous step, we're further filtering here to ensure that we include only the best contigs, as well as reducing the computational time for the next steps.
- `-o` - the name of the output file. Now that we've prepared the file, we can read this into `Anvi'o`:

```
mkdir anvio/anvio_databases
anvi-gen-contigs-database -f anvio/megahit_out/final.contigs.fixed.fa \
                        -o anvio/anvio_databases/CONTIGS.db \
                        -n CHKN -T 4
```

You can read more about what `anvi-gen-contigs-database` is doing here, but the key parts of this are:

- `-f` - the input file of fixed contigs
- `-o` - what the output should be saved as
- `-n` - the name we're giving this contigs database

Throughout the next steps, we're going to be adding information to this contigs database.

3.6 2.4 Run HMMs to identify single copy genes

The first thing that we're going to add to the contigs database is information on the genes within the contigs:

```
anvi-run-hmms -c anvio/anvio_databases/CONTIGS.db \
              --num-threads 4
```

Because the only arguments we've given this are the contigs database (`-c`) and the number of threads to use (`--num-threads`), it will run the default Hidden Markov Models (HMMs) within `Anvi'o`. You can read more about HMMs in this paper if you're interested, but in short, they are "statistical models that can be describe the evolution of observable events that depend on internal factors, which are not directly observable". HMMs are often used in computational biology for predicting the function of a gene or protein; we can build HMMs using multiple sequence alignments of genes or proteins known to carry out the same function, and then we can run these HMMs against our protein (amino acid) or gene (DNA) sequences to find other proteins/genes that are likely to carry out this function. You can read about the default HMMs in `Anvi'o` here, but they include `Bacteria_71` (71 single-copy core genes for the bacterial domain), `Archaea_76` (76 single-copy core genes for the archaeal domain), and `Protista_83` (83 single-copy core genes for the protists, within the eukaryotic domain).

If we had our own set of HMMs, or HMMs for specific genes of interest, we could also use these here.

As an added bonus, we can also infer taxonomy from these SCGs using:

```
anvi-run-scg-taxonomy -c anvio/anvio_databases/CONTIGS.db -T 4 -P 4
```

This next step exports the sequences for all of the genes that we've identified using the HMMs to a file called `anvio/anvio_databases/gene_calls.fa`:

```
anvi-get-sequences-for-gene-calls -c anvio/anvio_databases/CONTIGS.db \
    -o anvio/anvio_databases/gene_calls.fa
```

Question 4: How many genes were identified?

3.7 2.5 Map samples onto contigs with Bowtie2

The next steps are going to generate abundance profiles of our contigs within the samples - these will help us to bin the contigs together into groups that have similar abundance profiles and are therefore likely to come from the same genome.

The first step here is to build a `bowtie2` database/index of our fixed contigs:

```
bowtie2-build anvio/megahit_out/final.contigs.fixed.fa \
    anvio/megahit_out/final.contigs.fixed
```

You can see that here we just have two positional arguments - the file name for the fasta file that we want to use, and the prefix to use for the `bowtie2` index. If you look at the files in `anvio/megahit_out/` now, you should see six files with this prefix and the file extension `.bt2`. These files are what `bowtie2` will use in the next steps for mapping the samples to the contigs index.

For the next step, we've already made a file that contains the names of all of the samples that we'll be using (`sample_ids.txt`). We can use this to loop through each of the samples using a for loop, performing the same actions on the files for each sample.

NOTE: The loop is commented out and replaced with a soft-link.

```
# We add a soft-link here instead of running bowtie because it takes a long time to run :)
ln -s ~/CourseData/MIC_data/amb_module2/anvio/bam_files anvio/bam_files

#mkdir anvio/bam_files # because soft-link instead for today

#for SAMPLE in `awk '{print $1}' sample_ids.txt`
#do
#    # 1. do the bowtie mapping to get the SAM file:
#    bowtie2 --threads 4 \
```

```

#           -x anvio/megahit_out/final.contigs.fixed \
#           -U fastqs_for_assembly/${SAMPLE}.fastq.gz \
#           --no-unal \
#           -S anvio/bam_files/${SAMPLE}.sam
#
# # 2. convert the resulting SAM file to a BAM file:
# samtools view -F 4 \
#           -bS anvio/bam_files/${SAMPLE}.sam \
#           > anvio/bam_files/${SAMPLE}-RAW.bam
#
# # 3. sort and index the BAM file:
# anvi-init-bam anvio/bam_files/${SAMPLE}-RAW.bam \
#           -o anvio/bam_files/${SAMPLE}.bam
#
# # 4. remove the intermediate BAM file that is no longer needed:
# rm anvio/bam_files/${SAMPLE}-RAW.bam
#done

```

Let's confirm the soft-link worked from above:

```
ls anvio/bam_files
```

What do you see? There should be some `.bam` files, among other things.

Hopefully you either remember about for loops from Beginner Module 4 or you already have some experience with them! Otherwise, I'll explain them briefly here.

For loops are an incredibly useful thing that we do in programming - to my knowledge, they exist in every programming language, and they allow us to repeat a section of code many times, with a different input each time. To explain what we are doing above, we can go through this step-by-step. The code is essentially saying that for each of the rows in the `sample_ids.txt` file, we want to: (1) run `bowtie2`, (2) convert the output SAM file from `bowtie2` to a BAM file, (3) sort and index the BAM file, and (4) remove the intermediate BAM file. You can see that we give `${SAMPLE}` in several places, as these parts will change for each of the sample names.

In each of the steps within the for loop, we are giving several options:

1. We are running `bowtie2` using the contigs index that we created `-x`, reads for this sample, `-U`, and the `--no-unal` options means it won't output unaligned reads. We'll then get `-S` as the output SAM file. SAM (Sequence Alignment MAP) files are used to store information about sequences aligned to a reference and they have 11 mandatory fields corresponding to various metrics about how well the sequence aligns to the reference. You can read more about them [here](#).
2. `samtools` is being used to convert the SAM file to a BAM file. You can see that this uses the view module of `samtools`, which would

by default print the output to the terminal, but as we've added `> anvio/bam_files/${SAMPLE}-RAW.bam`, this tells it to save the output to a file. A BAM (Binary Alignment MAP) file is a compressed, binary version of a SAM file (see more [here](#)).

3. The running of `anvi-init-bam` is fairly straight-forward, we are just giving it the raw BAM file and it is outputting a sorted, re-indexed version.
4. Then we just remove the intermediate BAM file.

3.7.1 Add coverage and detection statistics to the Anvi'o profile

Now that we've created information on the coverage and detection of our contigs within our samples (how abundant they are and whether they appear in every sample or not), we can make Anvi'o profiles with this information:

NOTE: We're going to copy over the profiles this time, not a soft-link because we need to manipulate them.

```
mkdir anvio/anvio_databases/profiles
cp -r ~/CourseData/MIC_data/amb_module2/anvio/anvio_databases/profiles/* anvio/anvio_databases/profiles/

# also need to overwrite the CONTIG.DB so they match
cp ~/CourseData/MIC_data/amb_module2/anvio/anvio_databases/CONTIGS.db anvio/anvio_databases/CONTIGS.db

#for SAMPLE in `awk '{print $1}' sample_ids.txt`
#do
#
#    anvio-profile -c anvio/anvio_databases/CONTIGS.db \
#                -i anvio/bam_files/${SAMPLE}.bam \
#                --num-threads 4 \
#                -o anvio/anvio_databases/profiles/${SAMPLE}
#done
```

We're again using a for loop to carry this out on each of our samples. The `anvi-profile` program quantifies coverage per nucleotide position within the contigs, and averages them per contig. It also calculates single-nucleotide, single-codon, and single-amino acid variants, as well as structural variants such as insertion and deletions and stores these data into appropriate tables.

3.7.2 Merge sample profiles

The command above created contig profiles for each sample individually, so now we need to merge these profiles into one single profile that we can use for clustering the contigs into bins. We'll do this with the `anvi-merge` program:

```
anvi-merge -c anvio/anvio_databases/CONTIGS.db \
           -o anvio/anvio_databases/merged_profiles \
           anvio/anvio_databases/profiles/*/PROFILE.db
```

And then it's always a good idea to have a look at some basic stats about the contigs before we go further:

```
anvi-display-contigs-stats --report-as-text \
                          --output-file contigs_stats.txt \
                          anvio/anvio_databases/CONTIGS.db
```

If you take a look at the `contigs_stats.txt` file, you'll see (you can also see these explanations here:

- **Total Length** - the total number of nucleotides in your contigs
- **Num Contigs** - the number of contigs in your database
- **Num Contigs > X kb** - the number of contigs that are longer than X
- **Longest Contig** - the longest contig in your databases (in nucleotides)
- **Shortest Contig** - the shortest contig in your databases (in nucleotides), hopefully this is longer than the 2500 that we set earlier!
- **Num Genes (prodigal)** - the number of genes that Prodigal predicts are in your contigs
- **L50, L75, L90** - if you ordered the contigs in your database from longest to shortest, these stats describe the number of contigs you would need to go through before you had looked at a certain percent of a genome. For example, L50 describes the number of contigs you would have to go through before you reached 50 percent of the entire dataset
- **N50, N75, N90** - if you ordered the contigs in your database from longest to shortest, these stats describe the length of the contig you would be looking when you had looked at a certain percent of a genome. For example, N50 describes the length of contig you would be on when you reached 50 percent of the entire genome length
- The number of HMM hits in your database
- The number of genomes that Anvi'o predicts are in your samples, based on how many hits the single-copy core genes got in your database

As long as we're satisfied with all of this, then we can carry on to the clustering.

Question 5: How many contigs are there? Is this the same as what we started with? Why or why not? **Question 6:** What are the longest and shortest contigs? What do you think of this?

3.8 2.6 Cluster contigs into bins

Next we will bin - or cluster - the contigs to create genome "bins", or MAGs. To do this, we typically use information on the coverage, or abundance, of contigs within samples (that we generated in step two) and the binning algorithms will usually identify patterns in nucleotide composition or k-mer frequencies in order to group together contigs that they think are likely to have originated from the same genome. We are using CONCOCT for this, but as with every step of this pipeline, there are multiple tools that are capable of carrying this out

and may be better or worse in different circumstances. Drivers of choice include: CONCOCT, MAXBIN2, and METABAT2, but be mindful that each driver may have different flags that need to be set so the command below may vary.

Now we'll run the clustering:

```
DRIVER=CONCOCT # or MAXBIN2 or METABAT2
anvi-cluster-contigs -c anvio/anvio_databases/CONTIGS.db \
                    -p anvio/anvio_databases/merged_profiles/PROFILE.db \
                    -C "merged_${DRIVER}_2500" \
                    --driver ${DRIVER} \
                    --length-threshold 2500 \
                    --num-threads 4 \
                    --just-do-it
```

Here you can see that we're giving a few options:

- **-c** - the contigs database
- **-p** - the profile database
- **-C** - the name for the resulting collection of bins to be given
- **--driver** - the clustering algorithm to be used. We could also have used metabat2, maxbin2, dastool, or binsanity. If you're running this pipeline on your own data then you will probably want to think about which of these algorithms you should use; dastool has the advantage that it can take the input of multiple binning algorithms to calculate an optimised set of bins. It's also worth noting that no binning algorithm is perfect, and you will likely need to check the bins and refine them manually after any of these are run
- **--num-threads** - the number of threads to use
- **--just-do-it** - ignore any warnings (like CONCOCT being implemented experimentally into Anvi'o) and just run it anyway

Next, we'll summarise the bins that CONCOCT has given us:

```
mkdir anvio/${DRIVER}_summary_2500
anvi-summarize -c anvio/anvio_databases/CONTIGS.db \
               -p anvio/anvio_databases/merged_profiles/PROFILE.db \
               -C "merged_${DRIVER}_2500" \
               -o anvio/${DRIVER}_summary_2500/summary/
```

We can also now use the previously annotated contigs to infer taxonomy based on the SCGs from section 2.5:

```
anvi-estimate-scg-taxonomy -p anvio/anvio_databases/merged_profiles/PROFILE.db \
                           -c anvio/anvio_databases/CONTIGS.db \
                           -C merged_${DRIVER}_2500 \
                           -o anvio/${DRIVER}_summary_2500/summary/bins_taxonomy.txt \
                           --simplify-taxonomy-information --compute-scg-coverages
```

And we'll take a look at the bin summaries in the terminal quickly just to check that everything looks OK:

```
less -S anvio/${DRIVER}_summary_2500/summary/bins_summary.txt
less -S anvio/${DRIVER}_summary_2500/summary/bins_across_samples/abundance.txt
less -S anvio/${DRIVER}_summary_2500/summary/bins_across_samples/mean_coverage.txt
less -S anvio/${DRIVER}_summary_2500/summary/bins_taxonomy.txt
```

Notice we used the `-S` flag for less which chops long lines, making big tables easier to view. Remember that you can exit the less view by pressing `q`!

We can also look at this in a more interactive way in our workspace: <http://##.uhn-hpc.ca/> (remember to replace the `##` with your number!) Go to `amb_module2/anvio/CONCOCT_summary_2500/`. Right click on `summary` > open in new tab. You can scroll through and explore what this says about our bins so far, and then go back to the other tab with `amb_module2/anvio/CONCOCT_summary_2500/`. Add `/summary/bins_summary.txt` to the URL bar and copy and paste the resulting page into a new Excel (or whatever you usually use for viewing spreadsheets) document - it will be useful to refer back to. If you're in Excel, you can easily get this into columns by going to the "Data" tab > click on "Text to columns" > check "Delimited" > Next > Check "Space" > Finish.

Looking through the bins (the rows), you should see that there are a number of columns giving some stats on each of the bins (clusters of contigs), including:

- **total_length** - the total number of nucleotides in this bin
- **num_contigs** - the total number of contigs in this bin
- **N50** - N50 is a metric widely used to assess the contiguity of an assembly, and it is defined as the sequence length of the shortest contig that covers at least 50% of the total assembly length
- **GC_content** - GC content (%) is the percentage of bases in the contigs that make up this bin that are either guanine (G) or cytosine (C)
- **percent_completion** and **percent_redundancy** - **completeness** and **redundancy** (also called contamination by some tools) are both important when we're assessing the quality of the MAGs that we've assembled. Both scores are calculated based on the presence of the ubiquitous, single-copy marker genes - ones that all bacteria are known to possess - and the completion is a prediction of how complete the genome is likely to be, so whether it possesses a copy of all of the marker genes that are expected. The contamination/redundancy is a measure of whether those marker genes that are expected to be present in only a single copy are duplicated. Typically speaking, a MAG that has >50% completeness and <10% redundancy/contamination is considered to be reasonable, although >90% completeness is desirable. There is a good explanation on completeness and redundancy here. You'll see that we have quite a lot of bins that are 0% completion (also 0% redundancy), and typically these are made up of a single contig. Now, we'll just be focusing on the

bins that have completion >50% and we'll be trying to refine them so that their redundancy is <10%.

BONUS: Try running the MAXBIN2 and METABAT2 drivers through this process then compare all the bins from all the methods and see what you get!

Question 7: How many bins are there? **Question 8:** How many bins >50% completion are there? **Question 9:** What is the redundancy in these bins?

3.9 2.7 Interactive refinement of bins

Now we're going to be doing the interactive refinement, and first of all we'll have a look at the bin that's 100% complete and <10% redundant - for me, this is called Bin_13, but this naming is NOT consistent, so if that's a different bin for you then don't worry! Additionally, be sure you are using the correct `${DRIVER}` variable for the bin you're trying to view below, for me it is CONCOCT.

To do this, we'll run the Anvi'o `anvi-refine` command:

```
DRIVER=CONCOCT
anvi-refine -c anvio/anvio_databases/CONTIGS.db \
            -p anvio/anvio_databases/merged_profiles/PROFILE.db \
            -C "merged_${DRIVER}_2500" \
            -b Bin_13 \
            --server-only \
            -P 8081
```

You'll see that we're telling Anvi'o the contigs database, profile and collection name, as well as the name of the bin we want to look at and:

```
--server-only
-P 8081
```

Both of these parts are to do with Anvi'o being run on the Amazon instances rather than on our local computers. The `--server-only` part is telling it that we will want to create an SSH tunnel to the server, and then the `-P` port is telling it which port to use. This could be one of many ports, just like we are using port 8080 for accessing RStudio.

Now the instructions are different depending on whether we're on Mac/Linux or Windows.

3.9.0.1 Mac

Open a second terminal window or tab and run:

```
ssh -L 8081:localhost:8081 -i BMB.pem ubuntu@##.uhn-hpc.ca
```

Where `##` is your number. It should just look like you logged into the server in a new window. Now go to `http://localhost:8081/` in your browser. This

should have an Anvi'o page loaded up.

3.9.0.2 Windows

Open a new Putty window.

Click on “Session” and under “Saved Sessions”, click on the “Amazon node” and then click “Load”.

Now go to “Connection” > “SSH” > “Tunnels”. In “Source port” type in 8081. In “Destination” type in 90.uhn-hpc.ca:8081.

Click “Add” and then click “Open”. You should see a new Putty window open. Now go to <http://localhost:8081/> in your browser. This should have an Anvi'o page loaded up.

3.9.0.3 Back to both

Anvi'o opens in the **main** tab (look at the top left).

Now press the “Draw” button. You should see something that looks like a phylogenetic tree get drawn. Each branch of this is for one of the contigs that makes up the bin, and you can see information about their abundance in different samples in the rings.

Because there are so many samples (the track lines around the center tree), the tree in the middle is quite small. You can fix this by setting the radius to be larger for the tree. To do this find the **Show Additional Settings** button under the **Display** section (right under the top tabs). If you click that button it will open a menu of more advanced settings. The **Dendrogram** section is where you can modify certain aspects of the center tree. Setting the **Radius** so a larger number will expand the tree size and make it easier to view and manipulate. Try setting this value to 8000 and then click **Draw** again. Notice how the visual changes.

Now click on the “Bins” tab.

If we select the whole tree (click on it), we can see that this is 100% complete and 5.63% redundant. If we un-select (right click) some of the tree branches, this makes the completion and the length of the genome go down slightly, but it doesn't affect the redundancy.

Now that we've seen how this looks for a bin with 100% completion and 5.63% redundancy, we can leave this and try another.

Go back to the first Terminal window (where we ran **anvi-refine**) and press **ctrl+c** to stop.

Now run this with another bin. I want to find one that has higher **redundancy** but is also the same taxa as the previous bin, so I review my tables again and

find that in there is a bin in the MAXBIN2 set that is 97.18% completion and 30.98% redundancy, bin MAXBIN__002 for me but yours will vary.

```
DRIVER=MAXBIN2
anvi-refine -c anvio/anvio_databases/CONTIGS.db \
            -p anvio/anvio_databases/merged_profiles/PROFILE.db \
            -C "merged_${DRIVER}_2500" \
            -b MAXBIN__002 \
            --server-only \
            -P 8081
```

Refresh your Anvi'o browser window and you should see a different page now.

Again, click “Draw” and then go to the “Bins” tab. We ideally want to get the redundancy below 10% while keeping the completion above 50%.

Start by selecting the whole tree, and then begin by deselecting (right clicking) branches that appear to have a different abundance profile than others. For example, there is one that seems to only be present in sample CSM7KOMH - deselecting this one only reduces the completion a little (to 56.3%), but reduces the redundancy to 15.5%.

You can deselect large areas to determine whether they’ll have an impact on the redundancy, and then go through in more detail and deselect/reselect individual parts of this. Make sure you use the zoom +/- and zoom in and out to see the tree in more detail, as well as checking the abundance/GC content profiles of the parts that you’ve deselected.

Sometimes you might want to add another bin to keep some of the contigs that you don’t think belong in the first one (this particularly might be the case if you had really high completion initially).

When you are happy with what you’ve selected, click on “Store refined bins in the database”. You can now go back to the first console window and press ctrl+c.

The other two bins that are >50% completion are already <10% redundancy, so it’s not necessary to do anything to these, but you can have a look at them in the same way as we did with these two if you like.

3.10 2.8 Export the final bins into MAGs

Now that we’ve refined our bins into MAGs for our specified **driver**, we can rename them:

```
DRIVER=MAXBIN2
anvi-rename-bins -c anvio/anvio_databases/CONTIGS.db \
                 -p anvio/anvio_databases/merged_profiles/PROFILE.db \
                 --collection-to-read merged_${DRIVER}_2500 \
```

```

--collection-to-write FINAL_${DRIVER} \
--call-MAGs \
--min-completion-for-MAG 50 \
--max-redundancy-for-MAG 10 \
--prefix CHKN \
--report-file ${DRIVER}_renaming_bins.txt

```

You can see that here we're defining a new collection called `FINAL_${DRIVER}` (MAXBIN2 for me), and we're saying that to rename these bins as MAGs, they should be >50% completion and <10% redundancy. You can look at the file that's created with the renamed bins if you like, to check that this did what you expected: `${DRIVER}_renaming_bins.txt`.

Now we'll summarise this new collection that we've made:

```

anvi-summarize -c anvio/anvio_databases/CONTIGS.db \
                -p anvio/anvio_databases/merged_profiles/PROFILE.db \
                -C "FINAL_${DRIVER}" \
                -o anvio/FINAL_${DRIVER}_summary/

```

And we can have a look at the bins that we've created. Take a look at the FASTA files for each bin in the folder `anvio/FINAL_${DRIVER}_summary/bin_by_bin/` using the `ls` command. You should see that some of these are called *MAG* while some are called *Bin* - this is because in the command above, we're only calling things a MAG if they're >50% completion and <10% redundancy.

```
ls anvio/FINAL_${DRIVER}_summary/bin_by_bin/*/*.fa
```

Have a look at the summary file `anvio/FINAL_${DRIVER}_summary/bins_summary.txt`. You can see that most of them are the same as before, but the refined MAG is now present in place of the original bin.

Now have a look in one of those MAG folders (`anvio/FINAL_summary/bin_by_bin/`) - you'll see a lot of statistics, as well as a `*-contigs.fa`. This is a fasta file of the contigs used for each MAG, and we can use it as their genome for further analyses.

Now we'll create a copy of these in a new folder:

```

mkdir MAG_fasta
cp anvio/FINAL_*_summary/bin_by_bin/*/*contigs.fa MAG_fasta/
ls MAG_fasta/

```

You can also move all the bins/MAGs for all drivers if you would like to use them for downstream analysis!!!

Now that we have our refined MAGs, we can do anything that we like with them!

3.11 2.9 Run CheckM2

Next, we're going to run `checkm2`. `checkm2` can give us information on the quality of genomes.

First, we'll activate the conda environment:

```
mamba activate checkm2
```

Now we'll use the `predict` submodule of `checkm2`. This will take a little while to run, so you could do it in a `tmux` session if you like:

```
checkm2 predict -i MAG_fasta/ --output_directory MAGs_checkm2_output --allmodels -x .fa -t 4 --re
```

Using these models, `checkm2` predicts completeness and contamination scores, and assigns a taxonomic domain where possible. The results are saved in a summary file called `quality_report.tsv` inside the specified output directory, `MAGs_checkm2_output/`. This report includes genome names along with quality metrics such as `completeness`, `contamination (redundancy)`, `N50`, `genome size`, `GC content`, and `number of contigs`. By specifying `--remove_intermediates`, the pipeline deletes temporary processing files after prediction to reduce disk usage. The `-t 4` flag ensures the analysis runs using four CPU threads, speeding up processing for multiple genomes. This approach provides a fast, scalable, and model-driven alternative to traditional marker gene-based quality checks.

In this command, we're telling `checkm2` to run its predict quality assessment submodule on all of the `.fa` files within the `MAG_fasta` folder. It begins by scanning the `MAG_fasta/` directory for input files with the `.fa` extension, each representing a MAG. For each genome, `checkm2` extracts features such as `GC content` and `k-mer composition`, which are then passed into a set of pre-trained deep learning models. These models—general-purpose, bacteria-specific, and archaea-specific — are all applied, and the best-performing model is selected for each MAG.

Once it's finished running, you can take a look at the output from it, `MAGs_checkm2_output/quality_report.tsv` (either with `less -S` or in the file viewer).

NOTE: This is an estimate! Not a fact! Compare this to the purity and redundancy values from Anvi'o. There are many ways to assess your MAGs, you need to understand which are best for your purpose.

Question 10: What are the purity, redundancy, completeness, and contamination of your favorite MAG from today?

3.12 Other things that we frequently do with MAGs

While I hope that you now have a reasonable idea of the steps that go into assembling MAGs, how to assess their quality, and some of the ways in which we might annotate their functions, this has really only scratched the surface of what they might be used for.

If you still have time in the workshop and want to have a look, I've added some papers that I think are nice uses of MAGs as well as some other things that can be done with MAGs that you might want to explore.

3.12.0.1 Papers:

Bacterial ecology and evolution converge on seasonal and decadal scales

Nitrogen-fixing populations of Planctomycetes and Proteobacteria are abundant in surface ocean metagenomes

Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life

A genomic catalog of Earth's microbiomes

3.12.0.2 Other tools:

GTDB-toolkit - this can be used for assigning taxonomy to MAGs and has a much larger reference database than CheckM, which is designed for generating quality metrics of MAGs

Other Anvi'o workflows and tutorials

Programs available within Anvi'o

BlastKOALA and GhostKOALA

IslandViewer

METABOLIC

Chapter 4

Module 3 Lab: Assigning Functions

In this lab, we will predict and briefly examine the functions encoded in the assembled metagenome from the previous labs. Optional context and questions are provided as collapsible sections like so:

Click me

Optional context is provided to make this workshop as inclusive as possible. Feel free to skip any topics that you are already familiar with. Questions and answers will also use this collapsible format.

All instructions required to completed this lab are provided below, so you are free to work asynchronously at your own pace. Active discussions with both peers and instructors are encouraged. Say hi to the person next to you!

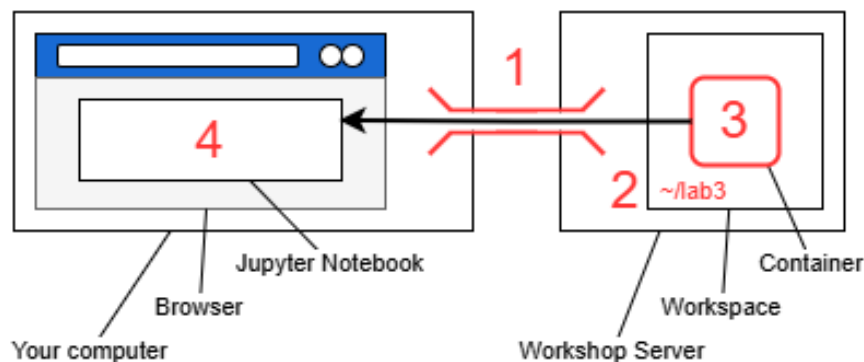
4.1 Setting up the workspace

Before we begin, this section will guide you through the set up of your workspace and compute environment. The software tools used in this lab all have dependencies that must be installed separately. Constructing the environment for each tool is non-trivial and error-prone. The difficulty of replicating compute environments for bioinformatics software is one of the key contributing factors to the reproducibility crisis [1].

We will instead be using prepackaged compute environments called “containers”, implemented by Apptainer (formerly Singularity) [2]. One of these containers will enable the use of jupyter notebooks along with a collection of utilities.

1. Mangul S, Martin LS, Eskin E, Blekhman R. Improving the usability and archival stability of bioinformatics software. *Genome Biol.* 2019;20(1):47. <https://doi.org/10.1186/s13059-019-1649-8>
2. Kurtzer GM, Sochat V, Bauer MW. Singularity: Scientific containers for mobility of compute. *PLOS ONE.* 2017;12(5):e0177459. <https://doi.org/10.1371/journal.pone.0177459>

4.1.0.1 Steps for setup:



1. Open an additional port via SSH for the new jupyter notebook
2. Create a new workspace directory for this lab
3. Unpack and start the container
4. Open the jupyter notebook in your browser

4.1.1 SSH port forwarding

We will use SSH port forwarding to open the additional port.

4.1.1.1 Terminal (Unix/MacOS)

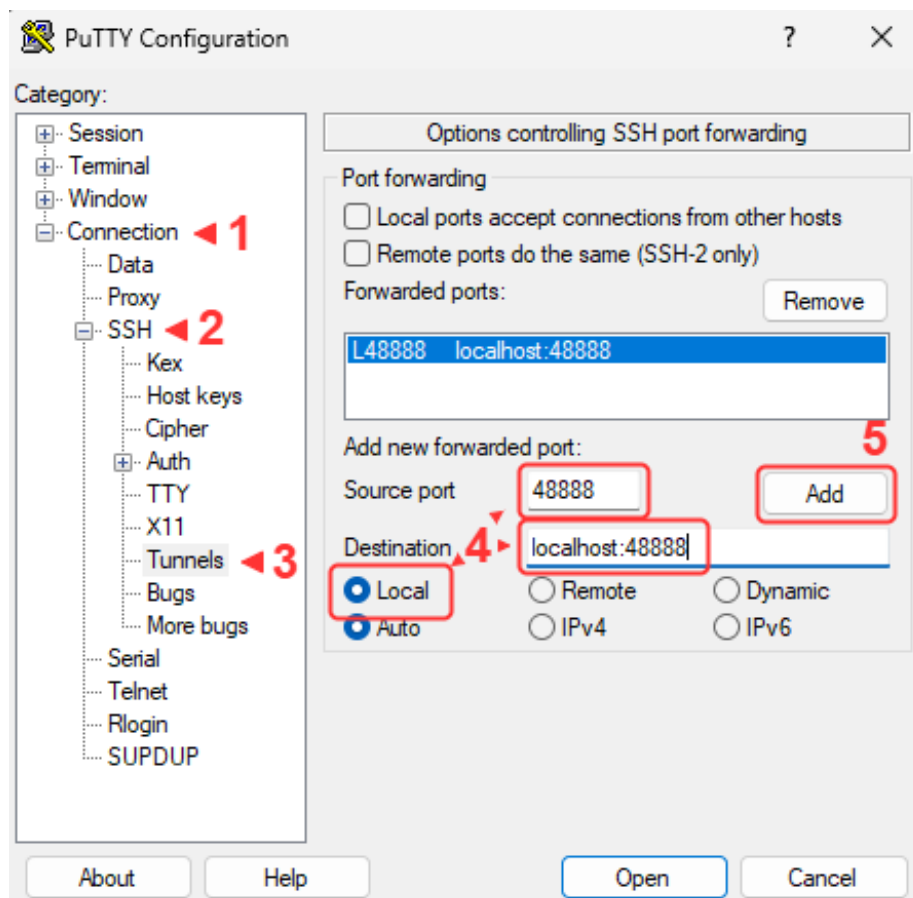
Add `-L 48888:localhost:48888` to the SSH command you use to connect to the server. For example:

```
1 ssh -L 48888:localhost:48888 ubuntu@${your_number}.uhn-hpc.ca -i ${path_to_your_key}
```

This means port 48888 on the remote server will be forwarded to your machine at `localhost:48888`.

4.1.1.2 Putty (Windows)

Modern versions of Windows should have `ssh` available through Powershell. Putty is available as a backup.



4.1.2 Workspace

Create a new workspace in your home directory.

```
cd ~
mkdir amb_lab3
cd amb_lab3
```

Link the resource folder for easy access.

```
ln -s ~/CourseData/MIC_data/amb_module3/ ./lib
```

4.1.3 Start container

Most containers for this lab were obtained from the biocontainers project, which provides a large collection of bioinformatics software in a standard format. For

each tool, we will provide URIs for the the specific container image used and where to get future versions. The `pull` command of `apptainer` can be used to download images (docs: https://apptainer.org/docs/user/main/cli/apptainer_pull.html).

Here is an example which downloads `prodigal` as the container image `example.sif` from `docker://quay.io/biocontainers/prodigal:2.6.3--h7b50bb2_10`. The image URI takes the form of `<protocol>://<address>:<tag>`. Tags are used to identify versions <https://quay.io/repository/biocontainers/prodigal?tab=tags>.

```
1 apptainer pull example.sif docker://quay.io/biocontainers/prodigal:2.6.3--h7b50bb2_10
```

We can now run the `prodigal` executable inside the container using `exec` (execute) like so.

```
apptainer exec ./example.sif prodigal -h
```

All containers for this lab were pre-downloaded in the resource folder now linked as `./lib`, including a main container with Jupyter and miscellaneous utilities. The main container provides an automated setup script named `unpack`. Let's run it.

```
apptainer exec ./lib/amb_lab3.sif unpack
```

The following will dedicate the current terminal to running the jupyter notebook on port 48888. You may want to connect another terminal to maintain terminal access.

```
./start_lab
```

4.1.4 Jupyter

JupyterLab should now be accessible at <http://localhost:48888/lab>. Create a new python notebook and import dependencies.

```
1 import pandas as pd
2 from pathlib import Path
3 from local.ipc import Shell
4
5 LIB = Path("./lib")
6 CPUS = 4 # your workshop instance has been allotted 4 cores
```

- `pandas` enables parsing and reading of tables
- `pathlib` enables manipulation of filesystem paths

- Pre-downloaded container images are in LIB
- Pre-computed outputs for this lab are available at LIB/"outputs".

Since containers can not be run inside other containers, the provided `Shell` function grants access to the terminal outside of the main container.

```
1 Shell(f"""
2 pwd -P
3 echo {LIB}
4 """)
```

Python comments `# ...`

Comments are free text not meant to be interpreted as python code and are indicated by a preceding `#`. They are useful for leaving notes.

```
1 # <helpful information>
```

Python variables `CPUS = 4`

Useful values can be saved to variables using `=` for later use.

```
1 x = 1 # the "value" 1 is assigned to the "variable" x
```

Python functions `Shell(...)`

Functions are saved procedures that can be repeatedly executed through “calls”. Some functions receive additional information as “arguments” or “return” information back when called. Called functions execute their procedure and are “evaluated” to their return value. Functions without a return evaluate to the dedicated null value `None`.

```
1 # function definition syntax
2 def function_name(): # no arguments
3     pass # do nothing, return nothing
4
5 # define the function "add"
6 def add(a, b):
7     return a + b
8
9 # call function "add"
10 x = add(1, 2) # add evaluates to 1 + 2, then to 3, then is assigned to x
11 x # x is 3
```

`Shell` accepts text and executes it as a bash command.

Python `import` and dot notation `local.ipc`

```
1 from local.ipc import Shell # access "ipc" from "local"
```

`local` is a folder containing the file `ipc.py`, which in turn defines the function `Shell`. `local` is a python “module” and `ipc.py` is a python “script”. `import` makes python components available from other files, typically called “libraries”. Dot notation (the period between `local` and `ipc`) retrieves the named components of complex objects. An alternative syntax to the above is:

```
1 import local.ipc.Shell # retrieve Shell from ipc in local
```

imported components can be renamed using the `as` keyword:

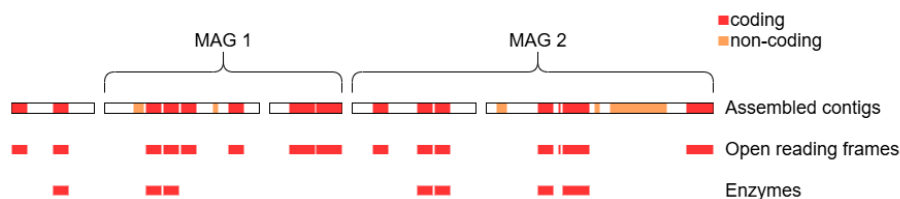
```
1 import local.ipc.Shell as Terminal # Shell is now "Terminal"
```

Python strings `f"..."`

Text values in python (strings), such as the contents of a bash command, must be declared in quotes like `"text"` or `'text'`. Triple quotes allows multi-line strings. The `f` prefix enables the evaluation of python variables in the string. For example, `f"echo {LIB}"` evaluates to `echo ./lib`. This is a convenient way to simplify commands by saving repeated sections to variables.

4.2 Overview

This lab will guide you through finding and identifying the functional units (represented by the red and orange features below) of a metagenomic assembly, with a focus on proteins and enzymes.



We will begin by exploring techniques related to sequence similarity, which seeks to identify known genes through comparison with reference databases.

- Prodigal
- Diamond
- Bakta
- Resistance Gene Identifier (RGI)

Then we will transition to machine-learning tools, which move away from the “lookup table” approach of sequence similarity.

- Kofamscan
- deepFRI
- deepEC
- proteinBERT

Finally, we will use gapseq to construct toy metabolic models using techniques to predict the existence of enzymes.

- gapseq

4.3 Prodigal

The container image for prodigal:

<https://quay.io/repository/biocontainers/prodigal?tab=tags>

```
1 # aptainer pull ./lib/prodigal.sif docker://quay.io/biocontainers/prodigal:2.6.3--h7b50bb2_10
2 prodigal_sif = LIB/"prodigal.sif"
```

The goal of prodigal is to predict the protein coding features (open reading frames; ORFs) in a given nucleotide sequence. Recall that this corresponds to the red features in the overview diagram, neglecting the orange features for now. It does this by optimizing predictions based on a set of tuned rules (heuristics) which take into account gene length and overlap, start and stop codons, ribosome binding sites, and GC content [1].

1. Hyatt D, Chen GL, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic gene recognition and translation initiation site identification. BMC Bioinformatics. 2010;11(1):119. <https://doi.org/10.1186/1471-2105-11-119>

Based on this, we would expect prodigal to:

- accept a nucleotide sequence as input
- output an itemized list of predicted open reading frames

The input sequence will be the metagenomic assembly since we would otherwise neglect some contigs that were not assigned to MAGs, thereby causing information loss. Let's have a look at prodigal's help message to get an idea of specific inputs and outputs.

```
1 Shell(f"aptainer exec {prodigal_sif} prodigal -h")
```

This looks like the input:

- -i: Specify FASTA/Genbank input file (default reads from stdin).

These look like they are useful for the output:

- `-a`: Write protein translations to the selected file.
- `-d`: Write nucleotide sequences of genes to the selected file. These will be helpful next lab for aligning transcripts to genes.
- `-f`: Select output format (gbk, gff, or sco). Default is gbk. We will choose gff since the genbank output of prodigal is not in a standard formatted.
- `-o`: Specify output file (default writes to stdout).

After a bit of investigation, it also looks like we should specify that the input is a metagenomic assembly.

- `-p`: Select procedure (single or meta). Default is single.

Have a go at compiling the final command.

```

1 out_dir = "./prodigal"
2 Shell(f"""
3 mkdir -p {out_dir}/
4 apptainer exec {prodigal_sif} \
5     prodigal \
6         -i {LIB}/inputs/assembly.fa \
7         ...
8 """)

```

The purpose of \

Starting from `apptainer exec`, the command must be written in a single line. The backslash `\` is used to ignore the newline character and continue the command on the next line. This lets us write break up a long command into multiple lines for readability. —

Suggested answer

```

1 out_dir = Path("./prodigal")
2 prodigal_aa = out_dir/"orfs.faa"
3 prodigal_gff = out_dir/"orfs.gff"
4 Shell(f"""
5 mkdir -p {out_dir}/
6 apptainer exec {prodigal_sif} \
7     prodigal \
8         -p meta \
9         -i {LIB}/inputs/assembly.fa \
10        -a {prodigal_aa} \
11        -d {prodigal_aa.with_suffix('.fna')} \
12        -f gff \
13        -o {prodigal_gff} \
14 """)

```

Let's take a peek at the `gff`.

```
1 gff = pd.read_csv(out_dir/"orfs.gff", sep="\t", header=None, comment="#")
2 print(gff.shape)
3 gff.head(2)
```

From left to right, the columns are:

0. contig id
1. name of the program that generated this feature
2. feature type, where CDS stands for coding sequence
3. start position (first nucleotide is 1)
4. end position
5. raw prodigal score
6. strand, where + is forward and - is reverse
7. codon frame, being one of 0, 1 or 2
8. a semicolon-separated list of tag-value pairs, providing additional information about each feature

Selecting entries in python dataframes and matrices

Entries (rows, columns, or cells) in pandas dataframes and numpy matrices can be selected in multiple ways.

Example matrix:

```
1 mat = np.array([1, 2, 3, 4, 5])
```

Select by single index, where negative indexes count from the end:

```
1 mat[0]      # 1
2 mat[-1]     # 5
```

a slice in the form of `start:end`, where `start` is inclusive and `end` is exclusive:

```
1 mat[1:4]    # [2, 3, 4]
2 mat[:2]     # [1, 2]
3 mat[2:]     # [3, 4, 5]
4 mat[-1:]    # [5]
```

a list of indexes:

```
1 mat[[0, 2, 4]] # [1, 3, 5]
```

a boolean mask:

```
1 mat[[False, False, False, True, True]] # [4, 5]
```

Let's investigate how many ORFs were not on MAGs. To do this, we will first create a blacklist for mag contig IDs.

```
1 mag_ids = set()
2 # loop through the paths to each of the MAG files
3 for mag_file in [
4     LIB/"inputs/HMP2_MAG_00001-contigs.fa",
5     LIB/"inputs/HMP2_MAG_00002-contigs.fa",
6 ]:
7     # we can use SeqIO to iterate through each sequence in a fasta file
8     for contig in SeqIO.parse(mag_file, "fasta"):
9         mag_ids.add(contig.id) # add the contig id of each contig to mag_ids
10 len(mag_ids)
```

Python for loops

Loops enable repeated execution of a block of code for each element in a collection.

```
1 collection = [1, 2, 3] # a "list" collection
2 for element in collection:
3     print(element) # do something with each element
4
5 # result:
6 # 1
7 # 2
8 # 3
```

Then filter the gff based on the MAG blacklist.

```
1 _filter = ~gff[0].isin(mag_ids) # column 0 is the contig id
2 non_mag_orfs = gff[_filter]
3 non_mag_orfs.shape
```

How many ORFs were not from MAGs?

1060 of 7176, or about 15%

4.4 Sequence similarity

In this section, we will use tools to find proteins with known sequences that are similar to that of the predicted ORFs. By assuming that similar sequences equate to similar functions, we can assign functions through reference databases of known proteins. While not perfect [1, 2], this is a common first-pass approach to assign functions to genes in metagenomic assemblies.

1. Omelchenko MV, Galperin MY, Wolf YI, Koonin EV. Non-homologous isofunctional enzymes: A systematic analysis of alternative solutions in enzyme evolution. *Biol Direct.* 2010;5(1):31. <https://doi.org/10.1186/1745-6150-5-31>
2. Seffernick JL, de Souza ML, Sadowsky MJ, Wackett LP. Melamine deaminase and atrazine chlorohydrolase: 98 percent identical but functionally different. *J Bacteriol.* 2001;183(8):2405–10. <https://doi.org/10.1128/JB.183.8.2405-2410.2001>

4.4.1 Diamond

<https://quay.io/repository/biocontainers/diamond?tab=tags>

```
1 # docker://quay.io/biocontainers/diamond:2.1.11--h5ca1c30_2
2 diamond_sif = LIB/"diamond.sif"
```

Diamond is highly optimized for finding similar protein sequences between a list of queries and a list of subjects [1]. It has a similar command line interface (CLI) to the classic BLAST, requiring that we first build a database from the list of subjects, before searching against it with a list of queries.

1. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Methods.* 2015;12(1):59–60. <https://doi.org/10.1038/nmeth.3176>

We expect to run diamond in two steps:

1. Build a reference database
 - inputs: a list of protein sequences
 - outputs: a database file
2. Search the database
 - inputs: a list of protein sequences, the compiled database
 - outputs: a list of hits

```
1 Shell(f"apptainer exec {diamond_sif} diamond --help")
```

We will build a database for transporters using the transporter classification database (TCDB) [1]. TCDB has been predownloaded from these URLs.

- protein sequences: <https://www.tcdb.org/public/tcdb>

- sequences are grouped into families: <https://www.tcdb.org/cgi-bin/projectv/public/families.py>
 - transporter substrates: <https://www.tcdb.org/cgi-bin/substrates/getSubstrates.py>
1. Saier MH Jr, Tran CV, Barabote RD. TCDB: the Transporter Classification Database for membrane transport protein analyses and information. *Nucleic Acids Res.* 2006;34(suppl_1):D181–6. <https://doi.org/10.1093/nar/gkj001>

```
1 Shell(f"apptainer exec {diamond_sif} diamond makedb")
```

The important arguments appear to be:

- `--db` database file
- `--in` input reference file in FASTA format/input DAA files for merge-daa

Have a go at compiling the command to build the diamond database. Running a your command is a good way to check for correctness and resulting error messages can be informative.

```
1 tcdb_fasta = LIB/"transporter_classification_db/tcdb.faa"
2 tcdb_db = Path("./diamond/tcdb.dmnd")
3 Shell(f"""
4 apptainer exec {diamond_sif} diamond makedb \
5     ...
6 """)
```

Possible solution

```
1 tcdb_fasta = LIB/"transporter_classification_db/tcdb.faa"
2 tcdb_db = Path("./diamond/tcdb.dmnd")
3 Shell(f"""
4 mkdir -p {tcdb_db.parent} # parent folder of given path
5 apptainer exec {diamond_sif} diamond makedb \
6     --threads {CPUS} \
7     --in {tcdb_fasta} \
8     --db {tcdb_db}
9 """)
```

Now we can search our predicted ORFs for transporters catalogued in TCDB.

```
1 Shell(f"apptainer exec {diamond_sif} diamond blastp")
```

That produced a lot of output, but these are the essential arguments:

- `--query` input query file
- `--db` database file
- `--out` output file

These may be frequently useful:

- `--fast` enable fast mode
- `--sensitive` enable sensitive mode
- `--ultra-sensitive` enable ultra sensitive mode
- `--threads` number of CPU threads
- `--outfmt` output format
- `--evalue` maximum e-value to report alignments
(default=0.001)

I find that this works well for general cases.

```

1 columns = "qseqid sseqid score pident evalue"
2 diamond_out = Path("./diamond/orfs.tsv")
3 Shell(f"""
4 apptainer exec {diamond_sif} \
5     diamond blastp \
6         --query ./prodigal/orfs.faa \
7         --db {tcdb_db} \
8         --out {diamond_out} \
9         --outfmt 6 {columns} \
10        --threads {CPUS}
11 """)

```

In addition to e-value, we can calculate a blast score ratio (BSR) to assign confidence to the hits. BSR for a hit is defined as the ratio of its score to the score of the query aligned to itself (the best possible score). A BSR of less than 0.4 was found to be a good threshold [1].

1. Rost B. Twilight zone of protein sequence alignments. Protein Eng. 1999;12(2):85–94. <https://doi.org/10.1093/protein/12.2.85>

Perform a self blast. We can skip making a database since the number of sequences is small.

```

1 diamond_self = Path("./diamond/self.raw")
2 Shell(f"""
3 apptainer exec {diamond_sif} \
4     diamond blastp \
5         --query {prodigal_aa} \
6         --db {prodigal_aa} \
7         --out {diamond_self} \
8         --outfmt 6 {columns} \
9         --max-target-seqs 1 \

```

```

10         --faster \
11         --threads {CPUS}
12     """)

```

“string”.split() and “string”.join()

Strings can be split into lists using the `split` method.

```

1  "a b c".split(" ") # ["a", "b", "c"]
2  "a/b/c".split("/") # ["a", "b", "c"]
3  "c_001_3".split("_") # ["c", "001", "3"]

```

Lists can be joined into strings using the `join` method.

```

1  " ".join(["a", "b", "c"]) # "a b c"
2  "_".join(["c", "001"]) # "c_001"

```

We now remove hits with a BSR < 0.4. First, load in the self blast results.

```

1  # get self scores
2  df_self = pd.read_csv(diamond_self, sep="\t", header=None, names=columns.split(" "))
3  df_self = df_self[df_self["qseqid"] == df_self["sseqid"]] # use boolean mask to select
4  print(df_self.shape)
5  df_self.head(2)

```

As the code gets more complex, it may be helpful to examine individual components to increase understanding. For example, have a look at:

```

1  columns.split(" ")
2  # and
3  df_self["qseqid"] == df_self["sseqid"]
4  # or even
5  df_self["qseqid"]

```

Calculate BSR

```

1  df_raw = pd.read_csv(diamond_raw, sep="\t", header=None, names=columns.split(" "))
2  # merge in self scores by matching the qseqid
3  # "left" refers to a "left join" where
4  # each row in the left dataframe (df_raw) is found a match in the right (df_self)
5  # duplicate columns from the right df are suffixed with "_self", such as "score_self"
6  df_hits = df_raw.merge(df_self, on="qseqid", suffixes=("", "_self"), how="left")
7  # calculate bsr and place in new column

```



```

8 df_hits["bsr"] = df_hits["score"] / df_hits["score_self"] # divide each value in score by the co
9 df_hits.head(2)

```

Filter out hits with BSR < 0.4

```

1 df_hits = df_hits[df_hits["bsr"] >= 0.4] # filter out rows with bsr < 0.4
2 df_hits = df_hits[columns.split(" ") + ["bsr"]] # keep only the columns we want
3 df_hits.head(2)

```

Take the best hit

```

1 # sort by bsr and group by qseqid, keeping the first row of each group
2 # this will keep the best hit for each query sequence
3 # groupby("qseqid") will set the index to qseqid, so we reset the index
4 df_hits = df_hits.sort_values("bsr", ascending=False).groupby("qseqid").first().reset_index()
5 df_hits = df_hits.sort_values("qseqid") # re-sort by qseqid
6 print(df_hits.shape)
7 df_hits.head(2)

```

Let's plot the results using a simple diverging bar chart for each MAG. First, let's create whitelists for the contigs in each MAG.

```

1 MAG1 = set()
2 for e in SeqIO.parse(LIB/"inputs/HMP2_MAG_00001-contigs.fa", "fasta"):
3     MAG1.add(e.id)
4 MAG2 = set()
5 for e in SeqIO.parse(LIB/"inputs/HMP2_MAG_00002-contigs.fa", "fasta"):
6     MAG2.add(e.id)
7 len(MAG1), len(MAG2)

```

Python set

Sets enable easy checking of membership using `in`. For example:

```

1 X = {1, 2, 3}
2 1 in X # True
3 0 in X # False
4 empty_set = set()

```

We can use the whitelists with additional parsing to get the family and MAG for each hit.

```

1 _fams = []
2 _mags = []
3 df_fam = pd.DataFrame(df_hits)
4 for _, row in df_fam.iterrows():
5     orf = row["qseqid"]
6     parts = orf.split("_")
7     contig = "_".join(parts[:-1])
8     if contig in MAG1:
9         mag = 1
10    elif contig in MAG2:
11        mag = 2
12    else:
13        mag = 0 # not in either MAG
14    _mags.append(mag)
15
16    description = row["sseqid"]
17    tcdb_id = description.split("|")[-1]
18    parts = tcdb_id.split(".")
19    family = ".".join(parts[:3])
20    _fams.append(family)
21
22 df_fam["family"] = _fams
23 df_fam["mag"] = _mags
24 df_fam = df_fam[df_fam["mag"] != 0] # remove no-MAG ORFs
25 print(df_fam.shape)
26 df_fam.head(2)

```

Python dict

Dictionaries enable the lookup of values by their corresponding key.

```

1 d = {"a": 1, "b": 2, "c": 3}
2 d["a"] # 1

```

By using `enumerate` to provide the index of each element in an iterable (set, list, dataframe, etc.), dictionaries that map from values to their index can be created.

```

1 values = ["a", "b", "c"]
2 value2i = {value: i for i, value in enumerate(values)}
3 value2i["a"] # 0

```

To visualize the results, we will first organize the data into a matrix of counts

for each transporter family and MAG.

```

1 mat = np.zeros(shape=(2, df_fam["family"].nunique())) # MAGs x families
2 mag2i = {mag: i for i, mag in enumerate([1, 2])}
3 fam2j = {fam: i for i, fam in enumerate(df_fam["family"].unique())}
4 for _, row in df_fam.groupby(["family", "mag"])["qseqid"].count().reset_index().iterrows():
5     i = mag2i[row["mag"]]
6     j = fam2j[row["family"]]
7     mat[i, j] = row["qseqid"]
8 order = np.argsort(-mat.sum(axis=0)) # produces a sorted list of indexes

```

Some utility functions are provided to make simple figures with Plotly. Let's use them to compare the abundance of transporter families in each MAG with a diverging bar plot.

```

1 from local.figures.template import BaseFigure, ApplyTemplate, go
2
3 family_labels = np.array([tcdb_families[k] for k in fam2j.keys()])
4 TOP_K = 10
5 fig = BaseFigure()
6 for mag in [1, 2]:
7     i = mag2i[mag]
8     sign = -1 if mag == 1 else 1 # -1 for left bar
9     fig.add_trace(
10         go.Bar(
11             x = sign*mat[i, order[:TOP_K]], # length of each bar
12             y = family_labels[order[:TOP_K]], # name of each row
13             orientation='h',
14             name = f"MAG{mag}",
15         )
16     )
17     # break # this will only draw the bars of one MAG. Uncomment "break" to see what each trace looks like
18
19 fig = ApplyTemplate(
20     fig,
21     layout=dict(
22         legend_orientation='h',
23         bargmode="relative",
24         width=1000, height=400,
25     ),
26 )
27 fig.show()

```

How to change the figure to show the top 3 instead?

```
1 TOP_K = 3
```

Can you put MAG1 on the right and MAG2 on the left?

```
1 sign = -1 if mag == 2 else 1
```

4.4.2 Bakta

<https://quay.io/repository/biocontainers/bakta?tab=tags>

```
1 # docker://quay.io/biocontainers/bakta:1.9.4--pyhdfd78af_0
2 bakta_sif = LIB/"bakta.sif"
```

Bakta is an automated workflow for annotating both coding and non-coding features using a variety of tools [1]. Its core protocol can be reduced to ORF prediction by Prodigal, followed by sequence similarity search by Diamond, as we have just done. Many additional tools and consolidated reference databases are used to search for non-coding features, provide information on ORFs with no good hits in any database, and perform QC.

1. Schwengers O, Jelonek L, Dieckmann MA, Beyvers S, Blom J, Goemann A. Bakta: rapid and standardized annotation of bacterial genomes via alignment-free sequence identification. *Microb Genomics*. 2021;7(11):000685. <https://doi.org/10.1099/mgen.0.000685>

To run Bakta, we expect to first set up reference databases before it will: - accept contigs, for which it will first run Prodigal for us - or accept a list of regions + contigs, for which Prodigal will be skipped - and produce many files containing useful annotations

Why is a protein sequence fasta insufficient to full utilize Bakta?

Non-coding features would not be represented.

The database was pre-installed, but here is some code to guide you if you decide to revisit in the future.

```
1 Shell(f"apptainer exec {bakta_sif} bakta_db list")

1 # bakta_db = LIB/"bakta_db/db"          # 2 hours, ~70 GB
2 bakta_db = LIB/"bakta_db/db-light"      # 30 mins, ~3.5 GB
```

```

3  if not bakta_db.exists():
4      if "light" in str(bakta_db):
5          _type = "light"
6      else:
7          _type = "full"
8      Shell(f"apptainer exec {bakta_sif} bakta_db download --output {bakta_db} --type {_type}")

```

There may be other annotation tools that we wish to run in addition to Bakta, which depend on predicted ORFs. Waiting for Bakta to finish before we receive the predicted ORFs may be inconvenient, especially for large metagenomic datasets. Running Prodigal ourselves and passing the ORFs to Bakta frees us to run additional tools in parallel. Bakta does not accept ORFs that run off the edge of contigs, so we will need to filter these out.

```

1  gff = pd.read_csv(prodigal_gff, sep="\t", header=None, comment="#")
2  complete_orfs = [] # a boolean mask
3  for _, row in gff.iterrows():
4      meta = tuple(row)[-1]
5      if "partial=00" in meta:
6          complete_orfs.append(True)
7      else:
8          complete_orfs.append(False)
9
10 gff_complete = "./prodigal/orfs.complete.gff"
11 gff[complete_orfs].to_csv(gff_complete, sep="\t", header=False, index=False)

```

Running Bakta, even on this reduced dataset, may take **15 minutes or more**. Precomputed outputs are available at {LIB}/outputs/bakta.

```

1  Shell(f"""
2  apptainer exec {bakta_sif} \
3      bakta --meta --threads {CPUS} --force --skip-plot \
4      --db {bakta_db} \
5      --keep-contig-headers \
6      --output ./bakta/ \
7      --regions {gff_complete} \
8      {LIB}/inputs/assembly.fa
9  """)

```

4.4.3 Resistance Gene Identifier (RGI)

<https://quay.io/repository/biocontainers/rgi?tab=tags>

```

1 # docker://quay.io/biocontainers/rgi:6.0.4--pyh05cac1d_0
2 rgi_sif = LIB/"rgi.sif"

```

Sometimes, specific functions are of interest, such as antibiotic resistance. RGI is a tool for searching the Comprehensive Antibiotic Resistance Database (CARD) [1, 2] using tuned protocols for Diamond.

1. McArthur AG, Waglehner N, Nizam F, Yan A, Azad MA, Baylay AJ, et al. The Comprehensive Antibiotic Resistance Database. *Antimicrob Agents Chemother.* 2013;57(7):3348–57. <https://doi.org/10.1128/aac.00419-13>
2. Alcock BP, Huynh W, Chalil R, Smith KW, Raphenya AR, Wlodarski MA, et al. CARD 2023: expanded curation, support for machine learning, and resistome prediction at the Comprehensive Antibiotic Resistance Database. *Nucleic Acids Res.* 2023;51(D1):D690–9. <https://doi.org/10.1093/nar/gkac920>

```

1 Shell(f"apptainer exec {rgi_sif} rgi --help")

```

```

1 Shell(f"apptainer exec {rgi_sif} rgi main --help")

```

Errors are a common occurrence for bioinformatics tools. The following is a command which appears to be correct, but will nonetheless fail. Read the error messages and use the hints to fix the issue. Additional warnings will appear that do not crash the program. These can be ignored. RGI also provides no indication of progress or successful completion, so checking the output folder will be necessary.

```

1 rgi_out = Path("./rgi/chicken")
2 Shell(f"""
3 apptainer exec {rgi_sif} \
4     rgi main \
5         --num_threads {CPUS} \
6         --clean \
7         -i {prodigal_aa} \
8         -o {rgi_out}
9 """)

```

`FileNotFoundError` hint

Is the path to the missing file related to any of the paths we gave it?

`FileNotFoundError` explanation

We indicated an output path within a new folder `rgi_out = Path("./rgi/chicken")`. It appears that RGI can't make this folder for us, so we will have to create

it ourselves with `mkdir -p {rgi_out.parent}`. `parent` refers to the parent folder.

invalid fasta hint

Nucleotide?

invalid fasta explanation

We can use the argument `-t protein` to indicate that the input is a protein sequence.

solution

```

1 rgi_out = Path("./rgi/chicken")
2 Shell(f"""
3 mkdir -p {rgi_out.parent}
4 apptainer exec {rgi_sif} \
5     rgi main \
6     --num_threads {CPUS} \
7     --clean \
8     -t protein \
9     -i {prodigal_aa} \
10    -o {rgi_out}
11 """)
12 rgi_out = rgi_out.with_suffix(".txt") # RGI adds the .txt suffix
13 # or
14 # rgi_out = Path("./rgi/chicken.txt")

```

Let's take a quick peak at the results.

```

1 Shell(f"""
2 _df = pd.read_csv(rgi_out, sep="\t", header=0)
3 _df
4 """)

```

4.5 Machine learning

Alignment-based approaches that assign functional annotations based on sequence similarity rely on databases of known sequences to transfer their known functions to the ORFs of a given sample. Machine learning may be able to discover more general patterns that can extend annotations to previously unseen sequences.

4.5.1 Kofamscan

https://quay.io/repository/hallamlab/external_kofamscan?tab=tags

```
1 # docker://quay.io/hallamlab/external_kofamscan:1.3.0
2 kofamscan_sif = LIB/"kofamscan.sif"
```

Kofamscan uses a machine learning technique called hidden markov models to learn the amino acid sequence patterns associated with various functions, organized in the KEGG database [1, 2]. Kofamscan sits somewhere in the middle between alignment-based and machine learning-based annotation methods.

1. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, et al. KofamKOALA: KEGG Ortholog assignment based on profile HMM and adaptive score threshold. *Bioinformatics*. 2020;36(7):2251–2. <https://doi.org/10.1093/bioinformatics/btz859>
2. Kanehisa M, Goto S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Res*. 2000;28(1):27–30. <https://doi.org/10.1093/nar/28.1.27>

```
1 Shell(f"apptainer exec {kofamscan_sif} kofamscan --help") # expect crash
```

It looks like we can't even get a help message with a naive call to the tool. Let's check the documentation for some clues. Specifically, there is short section marked **## Usage** that will be helpful.

https://www.genome.jp/ftp/tools/kofam_scan/README.md

Ah, so the executable is `exec_annotation`.

```
1 Shell(f"apptainer exec {kofamscan_sif} exec_annotation --help")
```

Kofamscan needs pre-trained HMM models, which have been downloaded from these links.

- ftp://ftp.genome.jp/pub/db/kofam/ko_list.gz
- <ftp://ftp.genome.jp/pub/db/kofam/profiles.tar.gz>

Kofamscan will take a long time to run. You may want to try running it on the hypothetical ORFs from the bakta output at `LIB/"outputs/bakta.full/assembly.hypotheticals.fasta"` instead. Pre-computed outputs are available at `LIB/"outputs/kofamscan/kos.out"`

```
1 kofamscan_raw = Path("./kofamscan/kos.out")
2 Shell(f"""
3 mkdir -p ./kofamscan
4 apptainer exec {kofamscan_sif} \
5     exec_annotation \
6     --cpu={CPUS} --format=detail --no-report-unannotated \
```



```

7         --profile={LIB}/kofamscan_db/profiles/prokaryote.hal \
8         --ko-list={LIB}/kofamscan_db/ko_list \
9         -o {kofamscan_raw} \
10        {prodigal_aa}
11    """)

```

Kofamscan attempts to assign KEGG orthology (KO) numbers to each ORF, linking them to rich KEGG databases which provide hierarchical organization and metabolic context. We can use the provided utility to parse the results into a partial model of KEGG for local use.

```

1  from local.models.kegg_orthology import ParseKofamScanResults
2
3  model = ParseKofamScanResults(
4      kofamscan_raw,
5      LIB/"kofamscan_db/api_kegg.db",
6      LIB/"kofamscan_db/brite.json",
7  )

```

```

1  model.Summary()

```

```

1  df = model.df
2  df = df[df.score > df.hmm_threshold] # retrieve only assignments that pass the model's adaptive cutoff
3  print(df.shape)
4  df.head(2)

```

Let's take a peek at a KEGG database entry.

```

1  model.GetComponentModel("K06442")._raw

```

We can also look at where the assigned function sits in the BRITE functional hierarchy.

```

1  for lineage in model.GetLineage("K06442"):
2      for k in lineage:
3          print(model.GetNodeInfo(k))

```

GetCategories provides a shorthand to retrieve 2 useful levels from the BRITE hierarchy.

```

1  for l2, l1 in model.GetCategories("K06442"):
2      for k in [l2, l1]:
3          print(model.GetNodeInfo(k))

```

Just like what we did with transporters, let's visualize the abundance of functions aggregated by BRITE categories for each MAG. First, we will compile a count of each category by MAG.

```

1 rows = []
2 for _, row in df.iterrows():
3     ko = row["ko"]
4     orf = row["orf"]
5     contig = "_".join(orf.split("_")[:-1])
6     if contig in MAG1:
7         mag = 1
8     elif contig in MAG2:
9         mag = 2
10    else:
11        mag = 0 # not in either MAG
12    for category, _ in model.GetCategories(ko):
13        rows.append((mag, category))
14 df_temp = pd.DataFrame(rows, columns=["mag", "category"])

```

Then, we can format it as a matrix for plotting.

```

1 # create some index lookup dictionaries for the matrix
2 c2i = {c: i for i, c in enumerate(df_temp["category"].unique())} # category to index
3 mag2i = {mag: i for i, mag in enumerate([1, 2])} # mag to index
4 mat = np.zeros(shape=(len(mag2i), len(c2i))) # MAG x categories
5 # iterate through the counts for each mag
6 # assigning them to the correct cell using the index lookup dictionaries
7 for _, row in df_temp.iterrows():
8     mag = row["mag"]
9     if mag == 0: continue
10    i = mag2i[mag]
11    j = c2i[row["category"]]
12    mat[i, j] += 1
13
14 order = np.argsort(-mat.sum(axis=0)) # get the order of categories by size (largest ->
15 mat.sum(axis=1)

```

Plot the top 10 largest categories.

```

1 from local.figures.template import BaseFigure, ApplyTemplate, go
2
3 category_labels = np.array([model.GetNodeInfo(k)[1] for k in c2i.keys()])
4 K = 10
5 fig = BaseFigure()
6 for mag in [1, 2]:

```

```

7     i = mag2i[mag]
8     sign = -1 if mag == 1 else 1
9     fig.add_trace(
10         go.Bar(
11             x = sign*mat[i, order[:K]],
12             y = category_labels[order[:K]],
13             orientation='h',
14             name = f"MAG{mag}",
15         )
16     )
17
18 fig = ApplyTemplate(
19     fig,
20     layout=dict(
21         legend_orientation='h',
22         barmode="relative",
23         width=1000, height=400,
24     ),
25 )
26 fig.show()

```

4.5.2 deepEC

https://quay.io/repository/hallamlab/external_deepec?tab=tags

```

1 # docker://quay.io/hallamlab/external_deepec:0.4.1
2 deepec_sif = LIB/"deepec.sif"

```

The DeepEC model uses convolutional neural networks, a deep learning architecture that is not far from HMM approach of Kofamscan. It too scans the protein sequence and attempts to predict function based on the “topology” or motif structure of the amino acid sequences. However, the model is able to make additional abstract decisions after the results of the initial scan in the form of additional model layers, further differentiating it from the sequence similarity approach [1]. As the name suggests, deepEC predicts enzyme commission (EC) numbers, which can be linked into the KEGG database, among others. It is straightforward to run.

1. Ryu JY, Kim HU, Lee SY. Deep learning enables high-quality and high-throughput prediction of enzyme commission numbers. *Proc Natl Acad Sci.* 2019;116(28):13996–4001. <https://doi.org/10.1073/pnas.1821905116>

```

1 Shell(f"""
2 aptainer exec {deepec_sif} \

```

```

3     deepec \
4         -p {CPUS} \
5         -i {prodigal_aa} \
6         -o ./deepec
7     """
)

1 dfec = pd.read_csv("./deepec/DeepEC_Result.txt", sep="\t")
2 print(dfec.shape)
3 dfec.head(2)

```

How many ORFs did deepEC assign EC numbers to?

87 out of 7176

4.5.3 deepFRI

https://quay.io/repository/hallamlab/external_deepfri?tab=tags

```

1 # docker://quay.io/hallamlab/external_deepfri:1.0.1
2 deepfri_sif = LIB/"deepfri.sif"

```

So deep learning models are able to consider the protein sequence at a higher-level of abstraction, but what is it abstracting to? We often don't know the mental model of protein function that these approaches have constructed for themselves, but deepFRI provides an example that can ease us into the deep learning world. DeepFRI predicts function by first examining the potential structure of a protein sequence by predicting the pairwise distances between all amino acids in the form of a “contact map” [1]. The final output are Gene ontology (GO) numbers, which are an alternative to KEGG more common in medical settings.

1. Gligorijević V, Renfrew PD, Kosciółek T, Leman JK, Berenberg D, Vatanen T, et al. Structure-based protein function prediction using graph convolutional networks. *Nat Commun.* 2021;12(1):3168. <https://doi.org/10.1038/s41467-021-23303-9>

DeepFRI is quite temperamental. We must remove the * character indicating stop codons from the prodigal output.

```

1 prodigal_aa_no_stop = Path("./prodigal/orfs.no_stop.faa")
2 with open(prodigal_aa_no_stop, "w") as f:
3     for e in SeqIO.parse(prodigal_aa, "fasta"):
4         if e.seq[-1] == "*":

```

```

5         e.seq = e.seq[:-1]
6         SeqIO.write(e, f, "fasta")

```

We must also create and run it within a dedicated output folder, while using absolute paths. As it runs, it will dump files into the current directory despite us telling it otherwise. Try with the reduced set of hypotheticals (`{LIB}/"outputs/bakta.full/assembly.hypotheticals.faa"`) or use the pre-computed results at (`{LIB}/"outputs/deepfri/"`).

```

1 prodigal_aa_no_stop = Path("./prodigal/orfs.no_stop.faa")
2 deepfri_out = Path("./deepfri")
3 Shell(f"""
4 mkdir -p {deepfri_out}
5 cp {prodigal_aa_no_stop} {deepfri_out}/
6 cd {deepfri_out}
7 aptainer exec -B ./:/ws -W /ws ../{deepfri_sif} \
8     deepfri \
9         --ont mf bp ec \
10         --fasta_fn /ws/{prodigal_aa_no_stop.name} \
11         --output_fn_prefix /ws/chicken
12 """)

```

EC:

```

1 df = pd.read_csv(LIB/"outputs/deepfri/chicken_EC_predictions.csv", comment="#")
2 print(df.shape)
3 df.head(2)

```

GO molecular function

```

1 df = pd.read_csv(LIB/"outputs/deepfri/chicken_MF_predictions.csv", comment="#")
2 print(df.shape)
3 df.head(2)

```

GO biological process

```

1 df = pd.read_csv(LIB/"outputs/deepfri/chicken_BP_predictions.csv", comment="#")
2 print(df.shape)
3 df.head(2)

```

4.5.4 proteinBERT

https://quay.io/repository/hallamlab/external_proteinbert?tab=tags

```

1 # docker://quay.io/hallamlab/external_proteinbert:2024.03.28
2 proteinbert_sif = LIB/"proteinbert.sif"

```

ProteinBERT “reads” protein sequences like reading text. The core of its architecture is an attention mechanism that is part of what’s called a “transformer”, often used in natural language processing [1, 2]. Residue context is considered and pairwise importance is predicted to enable the model to “focus” on what it thinks are important motifs. The pairwise attention could be spatial distance, as in deepFRI, but is not strictly limited to this interpretation, hence increasing the abstraction capabilities of the model. The core transformer architecture is used in current large language models like ChatGPT. The only difference is that ChatGPT is much, much larger.

1. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. Curran Associates Inc.; 2017. (NIPS’17).
2. Brandes N, Ofer D, Peleg Y, Rappoport N, Linial M. ProteinBERT: a universal deep-learning model of protein sequence and function. *Bioinformatics*. 2022;38(8):2102–10. <https://doi.org/10.1093/bioinformatics/btac020>

Since the model did not have friendly interface, a CLI was constructed courtesy of the Hallam Lab.

```

1 aptainer exec {proteinbert_sif} \
2     pbert run \
3         --threads {CPUS} \
4         -i {prodigal_aa} \
5         -o protein_bert

```

The output of proteinBERT is an array of numbers for each ORF, representing their predicted function. These numbers can be used as coordinates in proteinBERT’s mental model. We say each ORF is embedded in proteinBERT’s functional space. The euclidean distance between each ORF’s embedding represents predicted functional similarity.

```

1 df_index = pd.read_csv("./protein_bert/orfs.csv")
2 print(df_index.shape)
3 df_index.head(2)

```

ProteinBERT has 6 transformer layers, each of which reads then paraphrases the aminoacid sequence for the next. This information is stored in a 512-dimensional vector for each ORF, so the total embedding matrix is $N_{\text{orfs}} \times (6 \times 512)$.

```

1 mats = []
2 for batch in sorted(df_index.batch.unique()):
3     mats.append(np.load(f"./protein_bert/orfs.{batch}.embedding.npy"))

```

```

4 mat = np.vstack(mats)
5 mat.shape

```

Examining the standard deviation of the embeddings for each layers reveals that the last layer has the most to say, so we will use it to reduce needed compute.

```

1 fig = BaseFigure()
2 fig.add_trace(
3     go.Scatter(
4         x=np.arange(mat.shape[0]),
5         y=mat.std(axis=0),
6         mode="markers",
7         marker=dict(size=2),
8     )
9 )
10 fig = ApplyTemplate(fig, default_xaxis=dict(dtick=512))
11 fig.show()

```

```

1 mat_last = mat[:, -512:]
2 mat_last.shape

```

Let's load the kofamscan results to provide KEGG landmarks for these embeddings.

```

1 LIB = Path("./lib")
2 kofamscan_raw = Path("./kofamscan/kos.out")

```

```

1 from local.models.kegg_orthology import ParseKofamScanResults
2
3 ko_model = ParseKofamScanResults(
4     kofamscan_raw,
5     LIB/"kofamscan_db/api_kegg.db",
6     LIB/"kofamscan_db/brite.json",
7 )

```

```

1 df_kofam = ko_model.df
2 df_kofam = df_kofam[df_kofam["score"]>df_kofam["hmm_threshold"]]
3 print(df_kofam.shape)
4 df_kofam.head(2)

```

```

1 categories = {}
2 for _, row in df_kofam.iterrows():

```

```

3     ko = row["ko"]
4     orf = row["orf"]
5     for brite2, brite1 in ko_model.GetCategories(ko):
6         categories[brite2] = categories.get(brite2, set()) | {orf}
7
8 categories = sorted(categories.items(), key=lambda x: len(x[1]), reverse=True)
9 len(categories)

```

This assigns each ORF a category based on the BRITE hierarchy.

```

1 orf2cat = {}
2 for brite2, orfs in categories:
3     for orf in orfs:
4         if orf in orf2cat: continue
5         orf2cat[orf] = brite2
6 cat2i = {cat: i+1 for i, (cat, _) in enumerate(categories)}
7
8 y_categories = []
9 for _, row in df_index.iterrows():
10    orf = row["id"]
11    brite2 = orf2cat.get(orf, "None")
12    i = cat2i.get(brite2, 0)
13    y_categories.append(i)

```

We can use UMAP to reduce the 512 dimensional embeddings to 2 dimensions for visualization and use the BRITE categories as a guide.

```

1 from umap import UMAP
2
3 model = UMAP(
4     n_components=2, # reduce to x, y
5     spread=5,      # spread out the points
6     min_dist=0.1,  # reduce the maximum density of points
7     metric="cosine",
8 )
9
10 emb = model.fit_transform(mat_last, y=y_categories)

```

Let's have a look.

```

1 from local.figures.template import BaseFigure, ApplyTemplate, go
2 from local.figures.colors import Palettes, COLORS
3
4 fig = BaseFigure()

```



```

5
6 K = 10
7 seen = set()
8 todo = []
9 for brite2, _ in categories[:K]:
10     _filter = [orf2cat.get(orf) == brite2 for orf in df_index["id"]]
11     seen |= set(orfs)
12     todo.append((brite2, _filter))
13 _palette = Palettes.PLOTLY
14 for i, (name, _filter) in enumerate(todo):
15     fig.add_trace(
16         go.Scatter(
17             name=ko_model.GetNodeInfo(name)[1],
18             x=emb[_filter, 0],
19             y=emb[_filter, 1],
20             text=df_index[_filter]["id"],
21             mode="markers",
22             marker=dict(
23                 size=10,
24                 color=COLORS.TRANSPARENT,
25                 line=dict(
26                     width=1,
27                     color=_palette[i%len(_palette)].color_value,
28                 ),
29             ),
30         )
31     )
32
33 # plot rest as gray dots
34 _filter = ~df_index["id"].isin(seen)
35 fig.add_trace(
36     go.Scatter(
37         name="other",
38         x=emb[_filter, 0],
39         y=emb[_filter, 1],
40         text=df_index[_filter]["id"],
41         mode="markers",
42         marker=dict(
43             size=5, opacity=0.1,
44             color=COLORS.GRAY,
45         ),
46     )
47 )
48 no_axis = dict(linecolor=COLORS.TRANSPARENT, ticks=None, showticklabels=False)

```

```

49 fig = ApplyTemplate(
50     fig,
51     default_xaxis=no_axis, default_yaxis=no_axis,
52     layout=dict(
53         width=1000,
54         height=700,
55     ),
56 )
57 fig.show()

```

What are the gray dots?

These are ORFs not in the top K categories, including ORFs of unknown function.

4.6 gapseq

<https://quay.io/repository/biocontainers/gapseq?tab=tags>

```

1 # docker://quay.io/biocontainers/gapseq:1.4.0--h9ee0642_1
2 gapseqt_sif = LIB/"gapseq.sif"

```

So far, we have used various tools to predict the locations of genes and their functions in our metagenomic assembly. A portion of these functions were enzymatic reactions which convert between chemical compounds called metabolites. Others facilitate the transport of metabolites across membranes. Imagine reactions and transporters to be the edges of a network, chained together by common metabolites. The combined network would model the metabolic capacity encoded in the examined genetic sequences, drawing out possible pathways by which compounds are chemically transformed. Models from single genomes can indicate viability by possessing all necessary pathways to produce the constituents of biomass (nucleotides, amino acids, cofactors, lipids, etc.). Prediction errors that cause reactions to be missing create gaps in pathways, removing metabolic capacity, and potentially rendering the model non-viable. **gapseq** is a tool for filling these gaps from reference metabolic models by weighing candidate reactions against sequence evidence [1]. We will use gapseq to generate a model for each of the 2 MAGs by assuming viability on M9 minimal media.

1. Zimmermann J, Kaleta C, Waschina S. gapseq: informed prediction of bacterial metabolic pathways and reconstruction of accurate metabolic models. *Genome Biol.* 2021;22(1):81. <https://doi.org/10.1186/s13059-021-02295-1>

Can we run gapseq on the entire metagenomic assembly?

This would be a crude model of the community's metabolic capacity and can be used to detect possible pathways distributed across different members. A major source of error would be the lack of compartmentalization. The fact that a metabolite would have to be transported out of the cytosol of one cell and into the cytosol of another cell is not considered.

First, we separate the ORFs by MAG.

```

1 mag1_orfs, mag2_orfs = [], []
2 for e in SeqIO.parse(prodigal_aa, "fasta"):
3     contig = "_".join(e.id.split("_")[:-1])
4     if contig in MAG1:
5         mag1_orfs.append(e)
6     elif contig in MAG2:
7         mag2_orfs.append(e)
8
9 mag1_aa, mag2_aa = "./prodigal/mag1.faa", "./prodigal/mag2.faa"
10 with open(mag1_aa, "w") as f:
11     SeqIO.write(mag1_orfs, f, "fasta")
12 with open(mag2_aa, "w") as f:
13     SeqIO.write(mag2_orfs, f, "fasta")

```

We will have to repeat the gapseq protocol for each MAG.

```

1 gapseq_mag = "mag1"
2 # gapseq_mag = "mag2" # repeat with mag2
3 aa_fasta = f"./prodigal/{gapseq_mag}.faa"
4 out_dir = f"./gapseq/{gapseq_mag}"

```

Candidate reactions are searched for by gapseq *one pathway at a time*. This makes multiple (thousands) of redundant calls to the alignment tool which makes the process incredibly slow. This is because gapseq was intended to be an aid to manual curation and not a fully automated tool. Pre-computed outputs are available at {LIB}/outputs/gapseq. A quick trial with only a single pathway (glycolysis) is possible.

- -p keywords such as pathways or subsystems (for example amino,nucl,cofactor,carbo,polyamine)
- -l Select the pathway database (MetaCyc, KEGG, SEED, all; default: metacyc,custom)
- -K Number of threads for sequence alignments.
- -O Force offline mode
- -f Path to directory, where output files will be saved

```

1 gapseq_p = "glycolysis"
2 # gapseq_p = "all" # takes about 4 hours!!
3 Shell(f"""
4 mkdir -p {out_dir}
5 apptainer exec {gapseq_sif} \
6     gapseq find -p {gapseq_p} -l KEGG -K {CPU} -O -f {out_dir} {aa_fasta} \
7     > {out_dir}/find.out 2> {out_dir}/find.err
8 """)

```

> and 2>

> redirects the standard output and 2> redirects the standard error to a files for logging.

Transporters are faster, but still takes several minutes. Please use pre-computed results.

```

1 Shell(f"""
2 mkdir -p {out_dir}
3 apptainer exec {gapseq_sif} \
4     gapseq find-transport -K {CPU} -f {out_dir} {aa_fasta} \
5     > {out_dir}/findtr.out 2> {out_dir}/findtr.err
6 """)

```

Construct a draft model using candidate reactions, transporters, and pathways.

```

1 Shell(f"""
2 apptainer exec {gapseq_sif} \
3     gapseq draft \
4     -r {out_dir}/{gapseq_mag}-{gapseq_p}-Reactions.tbl \
5     -p {out_dir}/{gapseq_mag}-{gapseq_p}-Pathways.tbl \
6     -t {out_dir}/{gapseq_mag}-Transporter.tbl \
7     -c {aa_fasta} \
8     -f {out_dir} \
9     > {out_dir}/draft.out 2> {out_dir}/draft.err
10 """)

```

Add additional reactions to enable growth on M9 minimal media.

```

1 Shell(f"""
2 apptainer exec {gapseq_sif} \
3     gapseq fill --quick.gf \
4     -m {out_dir}/{gapseq_mag}-draft.RDS \
5     -c {out_dir}/{gapseq_mag}-rxnWeights.RDS \

```

```

6         -g {out_dir}/{gapseq_mag}-rxnXgenes.RDS \
7         -n {LIB}/gapseq/m9.csv \
8         --output.dir {out_dir} \
9         > {out_dir}/fill.out 2> {out_dir}/fill.err
10    """)

```

We can take a peek at the observed completion of each pathway. The model will be examined further in the next lab.

```

1    Shell(f"""
2    mag = 1
3    dfp = pd.read_csv(f"./gapseq/mag{mag}/mag{mag}-all-Pathways.tbl", sep="\t", comment="#")
4    dfp.sort_values("Completeness", ascending=False, inplace=True)
5    dfp
6    """)

```

The following code will generate a link to the KEGG pathway viewer with found reactions coloured. Pick a few and have a look.

```

1    from urllib.parse import quote
2    # https://www.genome.jp/kegg/webapp/color_url.html
3    pathway = "map00970"
4    color = "#00cc96"
5
6    reactions = dfp[dfp["ID"] == pathway]["ReactionsFound"].tolist()
7    reactions = [row.strip().split(" ") for row in reactions]
8    reactions = [f"{r} {color}" for row in reactions for r in row]
9    q = "\n".join(reactions)
10   url = f"https://www.kegg.jp/kegg-bin/show_pathway?map={pathway}&multi_query={quote(q)}"
11   url

```

This concludes the lab!

Chapter 5

Module 4 Lab: Visualization and Finding Functional Significance

In this lab, we will work through typical transcriptomic analyses to find significantly expressed genes and gene sets. Metagenomic reads will be used in place of transcriptomic data, so we will be examining abundance instead of expression, but the analyses shown will remain unchanged. For the sake of consistency, the abundance data will be referred to “expression” for the rest of the lab.

5.1 Setting up the workspace

The following procedure is near identical to that of lab3. Please refer to the setup section of lab3 for justifications and context. We will need to perform the following steps:

1. Open an additional port via SSH for the new jupyter notebook
2. Create a new workspace directory for this lab
3. Unpack and start the container
4. Open the jupyter notebook in your browser

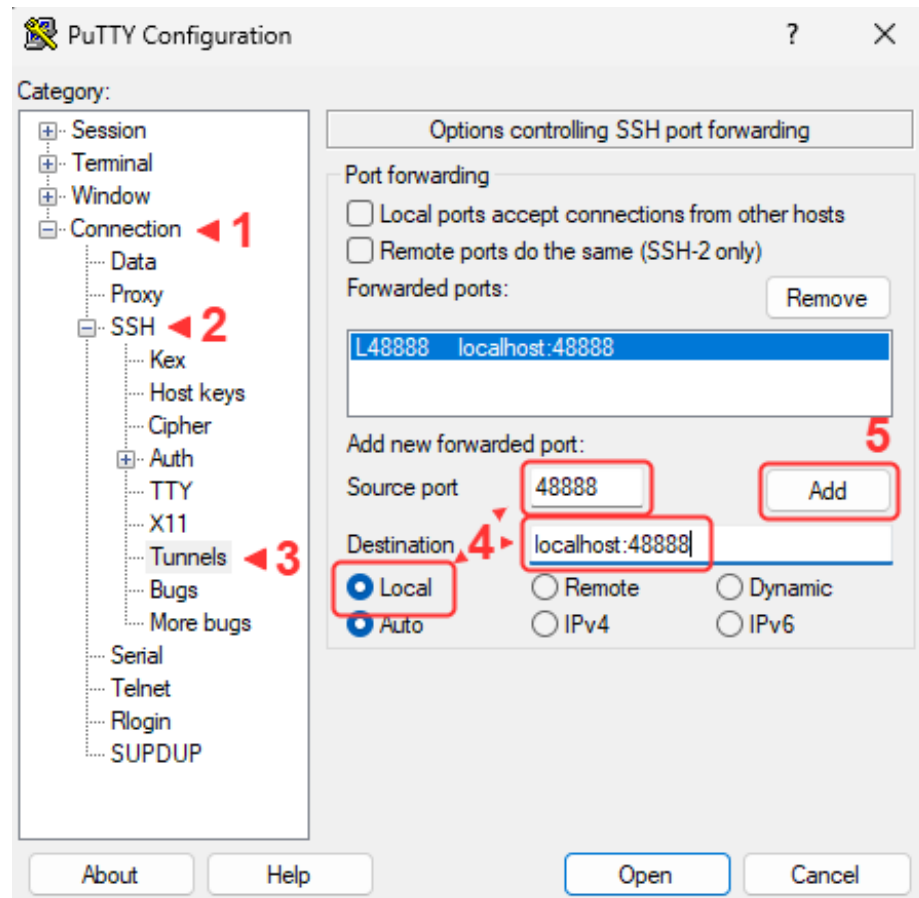
5.1.1 SSH port forwarding

5.1.1.1 Terminal (Unix/MacOS)

Add `-L 48888:localhost:48888` to the SSH command you use to connect to the server. For example:

```
1 ssh -L 48888:localhost:48888 ubuntu@${your_number}.uhn-hpc.ca -i ${path_to_your_key}
```

5.1.1.2 Putty (Windows)



5.1.2 Workspace

Create a new workspace in your home directory.

```
cd ~
mkdir amb_lab4
cd amb_lab4
```

Link the resource folder for easy access.


```
ln -s ~/CourseData/MIC_data/amb_module4/ ./lib
```

5.1.3 Start container

Unpack the main container for this lab.

```
apptainer exec ./lib/amb_lab4.sif unpack
```

The following will dedicate the current terminal to running the jupyter notebook on port 48888. You may want to connect another terminal to maintain terminal access.

```
./start_lab
```

5.1.4 Jupyter

JupyterLab should now be accessible at <http://localhost:48888/lab>. Create a new python notebook and import dependencies.

```
1 import os
2 import pandas as pd
3 from pathlib import Path
4
5 LIB = Path("./lib")
6 CPUS = 4 # your workshop instance has been allotted 4 cores
```

All software tools are packaged in the main container, so we can use `os.system(...)` instead of `Shell`.

```
1 os.system(f"""
2 pwd -P
3 echo {LIB}
4 """)
```

5.2 Overview

To understand the objective of this lab, let's first have a look at the metadata.

```
1 df_meta = pd.read_csv(LIB/"inputs/metagenome_metadata.txt", sep="\t")
2 print(df_meta.shape)
3 df_meta.head(2)
```

How many samples and how many categories?

We have 33 samples divided into “Bacteroides” and “Non-Bacteroides” groups.

We also have a set of predicted genes.

```

1 from Bio import SeqIO
2
3 count = 0
4 for i, sequence in enumerate(SeqIO.parse(GENES, "fasta")):
5     count += 1
6 count

```

The objective of this lab is to identify and visualize the pathways whose gene expression patterns correlate with the Bacteroides vs Non-Bacteroides grouping. To do this we will:

1. Use salmon to obtain a count matrix representing the expression level of each gene in each sample.
2. Use deseq2 to identify significantly expressed genes and generate a null distribution for gene set enrichment analysis.
3. Perform Gene set enrichment analysis (GSEA) to identify significantly enriched pathways.
4. Visualize the metabolic network in cytoscape with expression information.

5.3 Salmon

The alignment of transcriptomic reads to genes is often ambiguous. A single read can map to multiple genes and the length of genes compared to reads means that multiple reads will map to different regions of the same gene. Salmon uses probabilistic and statistical models to quantify the expression of each gene while also addressing potential biases [1]. We will use salmon to obtain a count matrix for deseq2.

1. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. Nat Methods. 2017;14(4):417–9. <https://doi.org/10.1038/nmeth.4197>

Let’s have a look at the help message.

```

1 os.system(f"salmon --help")

1 os.system(f"salmon index --help")

1 os.system(f"salmon quant --help")

```

We will need to index the reference genes before performing quantification.

```

1 salmon_dir = Path("./salmon")
2 GENES = LIB/"inputs/lab3_orfs.fna"
3 salmon_index = salmon_dir/"index"
4 os.system(f"salmon index -t {GENES} -i {salmon_index}")

```

Quantification will take about 6 minutes. `tqdm` is used to display a progress bar. Precomputed outputs are at `./lib/outputs/salmon/`.

```

1 from tqdm import tqdm
2
3 reads_dir = LIB/"inputs/reads"
4 samples = list(reads_dir.iterdir())
5 for read_file in tqdm(samples, total=len(samples)):
6     srr = read_file.name.split("_")[0]
7     out = salmon_dir/f"{srr}"
8     os.system(f"""
9         mkdir -p {out}
10        salmon quant \
11            -l A \
12            -p {CPUS} \
13            -i {salmon_index} \
14            -r {read_file} \
15            -o {out} \
16            >{out}/salmon.log 2>{out}/salmon.err
17        """)

```

Salmon requires information about the reads' library type, such as stranded, unstranded, facing, paired-end, etc. `-l A` will tell Salmon to infer the type from the first 1000 reads. More information is [here](#).

Let's have a look at one of the results. The file we want is `quant.sf`.

```

1 for srr in salmon_dir.iterdir():
2     if not srr.name.startswith("SRR"): continue
3     df = pd.read_csv(srr/"quant.sf", sep="\t")
4     break
5 df

```

The TPM (transcripts per million) column is useful on its own, but `deseq2` prefers raw counts, so we will use the `NumReads` column. Let's create the count matrix of genes by samples. First, create an index of genes.

```

1 gene2i = {}
2 for i, sequence in enumerate(SeqIO.parse(GENES, "fasta")):
3     gene2i[sequence.id] = i

```

```
4 len(gene2i)
```

Now, we will load the counts for each sample into the matrix by matching the gene name using the index.

What are the expected dimensions of the count matrix?

Genes x Samples, so 3298 x 33.

```
1 mtx = np.zeros(shape=(len(gene2i), len(samples)), dtype=np.int32) # genes x samples
2 for i, srr in tqdm(enumerate(df_meta["sample_name"]), total=len(df_meta)):
3     df = pd.read_csv(salmon_dir/f"{srr}/quant.sf", sep="\t")
4
5     gene_indexes = [gene2i[row["Name"]] for _, row in df.iterrows()]
6     counts = df["NumReads"]
7     mtx[gene_indexes, i] = counts
8 mtx.shape
```

Finally, we will label the rows and columns of the matrix with gene IDs and sample names before saving it to a file.

```
1 gene_ids = list(gene2i.keys())
2 sample_ids = list(df_meta["sample_name"])
3
4 df = pd.DataFrame(mtx, columns=sample_ids)
5 df["gene_id"] = gene_ids
6 df = df.set_index("gene_id")
7 print(df.shape)
8 df.to_csv(salmon_dir/"counts.mtx")
9 df.head(2)
```

We are now ready for `deseq2`.

5.4 DESeq2

DESeq2 is a statistical package for R that is used to analyze count data from RNA-Seq experiments [1]. It uses a negative binomial distribution to model count data, taking into account large biological diversity and technical noise (overdispersion). We will use the statistical models and tools provided by DESeq2 to calculate fold changes and estimate significance. The main purpose of DESeq2 in this lab is to generate the needed data for GSEA, which includes estimating the correlation between gene expression and the experimental groups: Bacteroides and Non-Bacteroides. We will also generate a null distribution for

GSEA by permuting the sample labels.

1. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. Nat Methods. 2017;14(4):417–9. <https://doi.org/10.1038/nmeth.4197>

Transition to an R notebook. <http://localhost:48888/lab>

```
1 library(DESeq2)
2 library(readr)
```

load the count matrix

```
1 cts <- as.matrix(read.csv(
2   "salmon/counts.mtx",
3   row.names = "gene_id"
4 ))
```

load the sample metadata

```
1 coldata <- read.csv(
2   "lib/inputs/metagenome_metadata.txt",
3   sep = "\t",
4 )
5 coldata$grouping
```

Create a Deseq dataset object.

```
1 dds <- DESeqDataSetFromMatrix(
2   countData = cts,
3   colData = coldata,
4   design = ~ grouping
5 )
6 dds
```

This is how we would run DESeq2, but don't run it just yet.

```
1 dds <- DESeq(dds)
```

We can cache the results to an “R data serialized” (RDS) to prevent having to re-run this analysis when restarting the notebook.

```
1 deseq_save <- "deseq2/dds.rds"
2 if (!file.exists(deseq_save)) {
3   dds <- DESeq(dds)
4   saveRDS(dds, deseq_save)
```

```

5 } else {
6   dds <- readRDS(deseq_save)
7 }

```

Extract the results and save to a dataframe. This will contain log2 fold changes and p-values for each gene.

```

1 res <- results(dds, contrast = c("grouping", "Bacteroides", "Non-Bacteroides"))
2 results_table <- as.data.frame(res)
3 write.csv(results_table, "./deseq2/results.csv", row.names = TRUE)

```

Next, we will compute the correlation between gene expression and the experimental groups and generate a null distribution for GSEA by permuting the sample labels. This effectively randomizes the assignment of samples to each experimental condition, but preserves any relationships between genes within a single sample. When the number of samples is low, genes can be randomly assigned to gene sets instead, but this fails to take into account gene-gene relationships.

What is the minimum number of samples to simulate a null distribution with a precision of at least 0.05?

4 samples ($4! = 24$, $1/24$ is about 0.04) This will be explored further in the next section on GSEA.

We will use the suggested variance stabilized transform (VST) to obtain normalized counts. VST removes the dependence of variance on the mean, which limits the high variance of the logarithm of counts when the mean is low. (more info [here](#)).

```

1 vst <- vst(dds, blind = FALSE)
2 vst_counts <- assay(vst)

```

The human-readable class labels need to be converted to more machine-friendly ordinal labels.

```

1 coldata$grouping
2 class_labels_ordinal <- 1 - (as.numeric(factor(coldata$grouping)) - 1)
3 class_labels_ordinal

```

We will use the point biserial correlation, a special case of the Pearson correlation designed to compare a continuous variable with a binary variable. Alternatively the biserial correlation achieves nearly the same, but assumes that the binary variable is the result of thresholding a continuous variable. For example, it

would be appropriate to use the point biserial correlation for a coin toss, while the biserial correlation would be more appropriate for pass/fails from percentage grades. Practically, the biserial correlation also requires the relative proportions of each class. While there likely is a genetic spectrum of bacteria spanning from Bacteroides to Non-Bacteroides, we don't have any information on this spectrum or the genetic threshold at which an organism stops being a Bacteroides.

```

1 point_biserial <- function(expression, classes) {
2   if (sd(expression) == 0) return(NA)
3   return(cor(expression, classes, method = "pearson"))
4 }

```

Estimate the observed correlation.

```

1 gene_correlations <- apply(vst_counts, 1, point_biserial, classes = class_labels_ordinal)
2 gene_correlations <- as.data.frame(gene_correlations)
3 write.csv(gene_correlations, "./deseq2/gene_correlations.csv", row.names = TRUE)

```

To generate the null distribution, we will be taking many samples (200). The following libraries will enable us to parallelize the computation.

```

1 library(pbapply)
2 library(parallel)

```

Since this is computationally expensive, we will use the caching technique shown earlier.

```

1 null_save <- "./deseq2/null_distr.csv"
2
3 if (!file.exists(null_save)) {
4
5   tryCatch(                                # if anything goes wrong, we must run the "finally" block
6     {
7       cl <- makeCluster(4) # make a cluster with 4 workers
8       # make these variables available to each worker
9       clusterExport(cl, c("vst_counts", "point_biserial", "class_labels_ordinal"))
10
11      n_permutations <- 200
12      permutation_results <- pbsapply(
13        1:n_permutations, # iterate 200 times
14        function(i) {     # and run this function each time
15          # Shuffle labels
16          n <- length(class_labels_ordinal)
17          shuffled_labels <- sample(class_labels_ordinal, size = n, replace = FALSE)
18

```

```

19         # Calculate correlations for all genes
20         return(apply(vst_counts, 1, point_biserial, classes = shuffled_labels))
21     },
22     cl = cl # use workers to compute in parallel
23 )
24 },
25 finally = {          # no matter what happens
26     stopCluster(cl)   # remember to stop the cluster
27 }
28 )
29
30 # cache results as a dataframe
31 null_distr <- data.frame(
32     sample = permutation_results,
33     stringsAsFactors = FALSE
34 )
35 write.csv(
36     null_distr,
37     null_save,
38     row.names = TRUE
39 )
40 } else {
41     null_distr <- read.csv(null_save)
42 }
43
44 dim(null_distr)

```

5.5 Gene set enrichment analysis (GSEA)

Back to python.

GSEA ranks genes by their correlation with experimental groups, then examines the overrepresentation of gene sets in highly ranked genes by calculating an enrichment score (ES) [1]. The observed ES compared to the null distribution of enrichment scores is then used to estimate the significance of the observed ES. To apply GSEA, we need:

- N samples where M continuous features are measured (eg. N transcriptomic samples for M genes, N metagenomes of M members, N pictures of M pixels, etc.)
- Groupings (overlap permitted) of these M features to test

Despite GSEA being widely applicable, implementations are overly specific to human transcriptomics, rendering them difficult to adapt to even adjacent domains. Fortunately, the technique itself is rather straightforward to implement.

1. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. Nat Methods. 2017;14(4):417–9. <https://doi.org/10.1038/nmeth.4197>

We will use KEGG pathways as our gene sets. These were predicted using Kofamscan from the previous lab.

```

1 from local.models.kegg_orthology import ParseKofamScanResults
2
3 kegg_model = ParseKofamScanResults(
4     LIB/"inputs/kofamscan/kos.out",
5     LIB/"kofamscan_db/api_kegg.db",
6     LIB/"kofamscan_db/brite.json",
7 )

```

We will find the pathways for each gene by tracing: gene -> function -> reaction -> pathway.

```

1 function2genes = {}
2 for _, r in kegg_model.df.iterrows():
3     ko = r["ko"]
4     orf = r["orf"]
5     function2genes[ko] = function2genes.get(ko, set()) | {orf}
6
7 reactions2functions = {}
8 for ko in function2genes:
9     meta = kegg_model.Functions[ko]
10    for r in meta.reactions:
11        reactions2functions[r] = reactions2functions.get(r, set()) | {ko}
12
13 pathway2reactions = {}
14 for r in reactions2functions:
15     meta = kegg_model.Reactions[r]
16     for p in meta.pathways:
17         pathway2reactions[p] = pathway2reactions.get(p, set()) | {r}
18
19 gene_sets = {} # collect the genes for each pathway
20 for p, reactions in pathway2reactions.items():
21     genes = set()
22     for r in reactions:
23         functions = reactions2functions[r]
24         for f in functions:
25             genes |= function2genes[f]
26     gene_sets[p] = genes

```

Load correlations from DESeq2.

```

1 df_null = pd.read_csv(deseq_dir/"null_distr.csv", index_col=0)
2 print(df_null.shape)
3 df_null.head(2)

1 df_corr = pd.read_csv(deseq_dir/"gene_correlations.csv", index_col=0)
2 print(df_corr.shape)
3 df_corr.head(2)

```

Next we will calculate the ES and p-value for an example pathway `rn00983`.

```

1 gene_set = gene_sets["rn00983"]
2 kegg_model.Pathways["rn00983"].name

```

First, create a matrix of permutations x genes, only taking those that are in gene sets.

```

1 enzymes = {x for g in gene_sets.values() for x in g}
2 def _only_enzymes(df: pd.DataFrame):
3     return df.loc[df.index.isin(enzymes)]
4
5 corr = _only_enzymes(df_corr).to_numpy()
6 permuted = _only_enzymes(df_null).to_numpy()
7 mat = np.hstack([corr, permuted]) # true correlations is the first permutation
8 mat.shape

```

Rank order the genes in each permutation by their correlation.

```

1 mat_index = np.argsort(-mat.T, stable=True)
2 mat_sorted = np.zeros_like(mat)
3 for i, order in enumerate(mat_index):
4     mat_sorted[:, i] = mat[:, i][order]
5 mat_sorted.shape

```

To calculate the ES, a running sum is calculated by scanning down the ranked list of genes, increasing the score for genes that are in the gene set and decreasing it otherwise. Genes in the gene set are called “hits” and those not in the gene set are called “misses”. The formulation for the running sum of hits (P_{hit}) at position i for gene set S is given in the paper. Don’t look at all of it at once, we will break it down.

First, note that the summations only consider genes in the gene set S .

$$P_{\text{hit}}(S, i) = \sum_{\substack{g_j \in S \\ j \leq i}} \frac{|r_j|^p}{N_R}, \quad \text{where } N_R = \sum_{g_j \in S} |r_j|^p$$

Let's create a mask to select only genes in S .

```

1 labels = _only_enzymes(df_corr).index.to_numpy()
2 S = np.array([k in gene_set for k in self.labels], dtype=int) # g in S
3 S = S[mat_index.T] # same order as mat

```

Next, have a look at N_R . r_j is simply the correlation values of each gene. p is a hyperparameter set to 1. Let's calculate N_R .

```

1 p = 1
2 corr = np.pow(np.abs(self.mat_sorted), p)
3 N_R = np.sum(corr*S, axis=0)

```

We now need to combine the repeated $|r_j|^p$ term in the numerator with N_R in the denominator and cumulatively sum up to each gene. `np.cumsum()` does this for us.

```

1 P_hit = corr/N_R
2 P_hit = P_hit*S # g_j in S
3 P_hit = np.cumsum(P_hit, axis=0)

```

Have a look.

```

1 from local.figures.template import BaseFigure, ApplyTemplate, go
2 fig = BaseFigure()
3 fig.add_trace(go.Scatter(
4     x = np.arange(len(P_hit[:, 0])),
5     y = P_hit[:, 0],
6 ))
7 fig = ApplyTemplate(fig)
8 fig.show()

```

For P_{miss} , we cumulatively sum the fraction of misses.

$$P_{\text{miss}}(S, i) = \sum_{\substack{g_j \notin S \\ j \leq i}} \frac{1}{(N - N_H)}.$$

```

1 P_miss = 1/(len(labels) - len(gene_set))
2 _mask = 1-S.T # invert hits for misses
3 P_miss = np.ones_like(labels)*_mask*P_miss
4 P_miss = np.cumsum(P_miss.T, axis=0)

```

That's it, have a look.

```

1 fig = BaseFigure()
2 fig.add_trace(go.Scatter(
3     x = np.arange(len(P_miss[:, 0])),
4     y = -P_miss[:, 0],
5 ))
6 fig = ApplyTemplate(fig)
7 fig.show()

```

The running statistic is the difference between P_{hit} and P_{miss} . Finally, the ES is the maximum magnitude of the running statistic.

```

1 S_distr = (P_hit - P_miss).T # running statistic
2
3 ES_pos = S_distr.max(axis=1)
4 ES_neg = S_distr.min(axis=1)
5 ES_candidates = np.vstack([ES_pos, ES_neg]).T # candidate is either the positive
6 to_take = np.argmax(np.abs(ES_candidates), axis=1) # we must select for each permutat
7 ES = (ES_neg*to_take) + (ES_pos*(1-to_take))

```

Have a look.

```

1 fig = BaseFigure()
2 fig.add_trace(go.Scatter(
3     x = np.arange(len(S_distr[0])),
4     y = S_distr[0],
5 ))
6 fig = ApplyTemplate(fig)
7 fig.show()

```

We can now split the score and null distribution and estimate the p-value.

```

1 score, null = ES[0], ES[1:]
2 pvalue = np.sum(null >= score) / null.shape[0]
3 if pvalue == 0:
4     pvalue = 1 / null.shape[0]
5 pvalue

```

Let's clean things up a bit. We will return a `dataclass`, which is just a convenient way to group many variables together. The procedure above is also attached to the class object `Gsea`, but otherwise unchanged.

```

1 from dataclasses import dataclass
2
3 @dataclass
4 class GseaResult:
5     set_name: str
6     set_size: int
7     score: float
8     pvalue: float
9     hits: np.ndarray
10    running_sum_statistic: np.ndarray
11
12 class Gsea:
13     def __init__(self, labels: np.ndarray, corr: np.ndarray, permuted: np.ndarray) -> None:
14         self.labels = labels
15         mat = np.hstack([corr, permuted])
16         self.mat_index = np.argsort(-mat.T, stable=True)
17         self.mat_sorted = np.zeros_like(mat)
18         for i, order in enumerate(self.mat_index):
19             self.mat_sorted[:, i] = mat[:, i][order]
20
21     def run(self, name: str, gene_set: set, p=1):
22         S = np.array([k in gene_set for k in self.labels], dtype=int)
23         if S.sum() == 0: return
24         S = S[self.mat_index.T]
25
26         corr = np.pow(np.abs(self.mat_sorted), p)
27         N_R = np.sum(corr*S, axis=0)
28
29         P_hit = corr/N_R
30         P_hit = P_hit*S # g_j in S
31         P_hit = np.cumsum(P_hit, axis=0)
32
33         P_miss = 1/(len(self.labels) - len(gene_set))

```

```

34     _mask = 1-S.T # invert hits for misses
35     P_miss = np.ones_like(self.labels)*_mask*P_miss
36     P_miss = np.cumsum(P_miss.T, axis=0)
37
38     S_distr = (P_hit - P_miss).T
39
40     ES_pos = S_distr.max(axis=1)
41     ES_neg = S_distr.min(axis=1)
42     ES_candidates = np.vstack([ES_pos, ES_neg]).T
43     to_take = np.argmax(np.abs(ES_candidates), axis=1)
44     ES = (ES_neg*to_take) + (ES_pos*(1-to_take))
45     score, null = ES[0], ES[1:]
46     pvalue = np.sum(null >= score) / null.shape[0]
47     if pvalue == 0:
48         pvalue = 1 / null.shape[0]
49
50     return GseaResult(
51         set_name=name,
52         set_size=len(gene_set),
53         score=score,
54         pvalue=pvalue,
55         hits=S.T[0],
56         running_sum_statistic=S_distr[0],
57     )
58
59 def _only_enzymes(df: pd.DataFrame):
60     return df.loc[df.index.isin(enzymes)]
61 _df_corr_e = _only_enzymes(df_corr)
62 gsea = Gsea(
63     labels=_df_corr_e.index.to_numpy(),
64     corr=_df_corr_e.to_numpy(),
65     permuted=_only_enzymes(df_null).to_numpy(),
66 )

```

We can now run GSEA for all pathways.

```

1 results: list[GseaResult] = []
2 for k, s in tqdm(gene_sets.items(), total=len(gene_sets)):
3     res = gsea.run(k, s)
4     results.append(res)
5 results_map = {r.set_name: r for r in results}

```

And print out some relevant results.

```

1 results = sorted(results, key=lambda x: x.pvalue)
2 for r in results:
3     if r.pvalue > 0.05: continue
4     meta = kegg_model.Pathways[r.set_name]
5     print(f"{r.pvalue:.2} {r.score:.2} {r.hits.sum()} {r.set_name}:{meta.name}")

```

Let's save the results to a table.

```

1 rows = []
2 for res in results:
3     rows.append((
4         res.set_name,
5         res.pvalue,
6         res.score,
7     ))
8
9 df_gsea = pd.DataFrame(rows, columns=["set_name", "pvalue", "score"])
10 Path("./gsea").mkdir(exist_ok=True, parents=True)
11 df_gsea.to_csv("./gsea/results.csv", index=False)

```

The nice GSEA plot is composed of 3 elements. First, we indicate the rank of genes within the selected gene set.

```

1 res = results_map["rn00983"]
2
3 fig = BaseFigure()
4 _f = res.hits.astype(bool)
5 fig.add_trace(
6     go.Scatter(
7         x = np.arange(gsea.labels.shape[0])[_f],
8         y = np.zeros_like(gsea.labels)[_f],
9         mode = "markers",
10        marker=dict(
11            size=25,
12            color=COLORS.TRANSPARENT,
13            symbol="line-ns",
14            line=dict(width=1),
15        ),
16        showlegend=False,
17    ),
18 )
19 fig = ApplyTemplate(fig)
20 fig.show()

```

Next, we draw out the magnitudes of the correlations.

```

1 fig = BaseFigure()
2 fig.add_trace(
3     go.Bar(
4         x = np.arange(gsea.labels.shape[0]),
5         y = gsea.mat_sorted.T[0],
6         showlegend=False,
7         marker=dict(
8             color=COLORS.BLACK,
9             line=dict(width=0),
10        ),
11    ),
12 )
13 fig = ApplyTemplate(fig)
14 fig.show()
```

And draw out the running sum statistic for the ES.

```

1 fig = BaseFigure()
2 fig.add_trace(
3     go.Scatter(
4         x = np.arange(gsea.labels.shape[0]),
5         y = res.running_sum_statistic,
6         mode = "lines",
7         showlegend=False,
8         line=dict(
9             color = COLORS.BLACK,
10        ),
11    ),
12 )
13 fig = ApplyTemplate(fig)
14 fig.show()
```

Finally, we stack the 3 together to replicate the figures used in the GSEA paper.

```

1 def draw_gsea(res: GseaResult, gsea: Gsea):
2     k = res.set_name
3     fig = BaseFigure(
4         shape=(1, 3), # 3 subplots, vertically stacked
5         row_heights=[0.2, 0.2, 0.6], # spaced like so
6         vertical_spacing=0.05,
7     )
8
9     _f = res.hits.astype(bool)
```



```

10     fig.add_trace(
11         go.Scatter(
12             x = np.arange(gsea.labels.shape[0])[_f],
13             y = np.zeros_like(gsea.labels)[_f],
14             mode = "markers",
15             marker=dict(
16                 size=25,
17                 color=COLORS.TRANSPARENT,
18                 symbol="line-ns",
19                 line=dict(width=1),
20             ),
21             showlegend=False,
22         ),
23         col=1, row=1,
24     )
25
26     fig.add_trace(
27         go.Bar(
28             x = np.arange(gsea.labels.shape[0]),
29             y = gsea.mat_sorted.T[0],
30             showlegend=False,
31             marker=dict(
32                 color=COLORS.BLACK,
33                 line=dict(width=0),
34             ),
35         ),
36         col=1, row=2,
37     )
38
39     fig.add_trace(
40         go.Scatter(
41             x = np.arange(gsea.labels.shape[0]),
42             y = res.running_sum_statistic,
43             mode = "lines",
44             showlegend=False,
45             line=dict(
46                 color = COLORS.BLACK,
47             ),
48         ),
49         col=1, row=3,
50     )
51
52     hidden = dict(ticks=None, linecolor=COLORS.TRANSPARENT, showticklabels=False)
53     fig = ApplyTemplate(

```

```

54         fig,
55         default_xaxis=hidden,
56         default_yaxis=hidden,
57         layout=dict(
58             width=600, height=400,
59             margin=dict(l=15, r=15, t=50, b=15),
60             title=dict(text=f"(p={res.pvalue}) {k}:{kegg_model.Pathways[k].name}", font
61         ),
62     )
63     fig.show()
64
65 draw_gsea(results_map["rn00983"], gsea)

```

If you have time, use the `draw_gsea` function to explore other pathways.

5.6 Cytoscape

Cytoscape is a powerful visualization environment for biological networks [1, 2]. We will combine the previously generated gapseq metabolic model with the gene expression data from this lab for visualization.

1. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Res.* 2003;13(11):2498–504. <https://doi.org/10.1101/gr.1239303>
2. Cline MS, Smoot M, Cerami E, Kuchinsky A, Landys N, Workman C, et al. Integration of biological networks and gene expression data using Cytoscape. *Nat Protoc.* 2007;2(10):2366–82. <https://doi.org/10.1038/nprot.2007.324>

The gapseq model is saved in systems biology markup language (SBML) format, which is a standard for representing biochemical networks. We will load the model and convert it to a python dictionary for easy usage.

```

1 import xmltodict
2
3 MAG = "mag1"
4 # MAG = "mag2" # Uncomment to use the second MAG
5 sbml = LIB/f"inputs/gapseq/{MAG}/{MAG}.xml"
6 with open(sbml, "r") as f:
7     gs_model = xmltodict.parse(f.read())

```

It will be helpful to examine the raw file itself for orientation. Much of the information we need is burried in a matryoshka doll-like structure of nested objects. Our first objective is to extract the metabolic network consisting of

reactions connected to compounds. We will use both reactions and compounds as nodes in the network, with edges connecting reactions to compounds. As we traverse the network, we will alternate between reactions and compounds.

```

1  # open the matryoshka data structure and find the information we need
2  def find(d, k):
3      if k in d: # found it
4          return [d[k]]
5      candidates = []
6      for v in d.values(): # we didn't find it, so go deeper
7          if isinstance(v, dict):
8              candidates += find(v, k)
9          elif isinstance(v, list):
10             for x in v:
11                 candidates += find(x, k)
12     return candidates
13
14 # retrieve list of reactions
15 raw_reactions = gs_model["sbml"]["model"]["listOfReactions"]["reaction"]
16 # store some useful information
17 names = {} # node id to reaction or compound names
18 r2genes = {} # reactions to genes
19 c2r = {} # compound to reactions
20 r2c = {} # reactions to compounds
21
22 # collect reactions and their associated genes and compounds
23 for r in raw_reactions:
24     rxn_id = r["@id"]
25     names[rxn_id] = r.get("@name", rxn_id)
26     for gene_id in find(r, "@fbc:geneProduct"):
27         r2genes[rxn_id] = r2genes.get(rxn_id, []) + [gene_id]
28     compounds = find(r, "@species")
29     for cpd_id in compounds:
30         c2r[cpd_id] = c2r.get(cpd_id, []) + [rxn_id]
31         r2c[rxn_id] = r2c.get(rxn_id, []) + [cpd_id]
32
33 # collect compound names
34 raw_compounds = gs_model["sbml"]["model"]["listOfSpecies"]["species"]
35 for c in raw_compounds:
36     cpd_id = c["@id"]
37     names[cpd_id] = c.get("@name", cpd_id)
38
39 len(r2c), len(c2r)

```

Certain compounds, such as H⁺, ATP, NADH, are ubiquitous and reactions connected by these compounds may not be very related. We can generate weights

for edges to filter out low-information connections using the “promiscuity” of each compound.

```

1 compound_promiscuity = {k: len(v) for k, v in c2r.items()}
2
3 names["M_cpd00011_c0"], compound_promiscuity["M_cpd00011_c0"]

```

The positions of network nodes on the 2D canvas of Cytoscape is called a graph projection. It would be useful if similar or biochemically adjacent reactions were positioned close to each other in this projection. We can achieve this by using UMAP with pairwise distances calculated from the shared compounds between 2 reactions, weighted by the promiscuity of each compound. In other words, we will use weighted jaccard similarity to obtain a graph projection.

```

1 import numpy as np
2 from tqdm import tqdm
3
4 # helper function to get connected reactions that share compounds
5 def adjacent_reactions(rxn):
6     found_reactions = []
7     for c in r2c[rxn]:
8         for r in c2r[c]:
9             if r == rxn: continue
10            found_reactions.append(r)
11    return found_reactions
12
13 # we scale the promiscuity values to be between 1 and 0 using a sigmoid function
14 # the more promiscuous a compound is, the less weight it has in the similarity calculation
15 _promiscuity_values = np.array([compound_promiscuity[c] for c in c2r.keys()])
16 b = 5
17 _scaled_promiscuity = 0.5*np.e/(1+np.exp((-_cproms-b-1)/b)) # sigmoid, dist=0.5 at 8, n
18 compound_scaled_weights = {c: w for c, w in zip(c2r.keys(), _scaled_promiscuity)}
19
20 # calculate the weighted jaccard similarity between reactions
21 # jaccard = shared / total
22 reaction_similarity = {}
23 for r1 in tqdm(r2c, total=len(r2c)):
24     for r2 in adjacent_reactions(r1):
25         k = (r1, r2) if r1 < r2 else (r2, r1)
26         if k in reaction_similarity: continue
27         c1 = set(r2c[r1])
28         c2 = set(r2c[r2])
29         cintersect = c1&c2
30         cunion = c1|c2
31         total_weight = sum(compound_scaled_weights[c] for c in cunion)

```

```

32     intersect_weight = sum(compound_scaled_weights[c] for c in cintersect)
33     weighted_jaccard_sim = intersect_weight/total_weight
34     if weighted_jaccard_sim < 1e-6: continue
35     reaction_similarity[k] = weighted_jaccard_sim
36
37 min(reaction_similarity.values()), max(reaction_similarity.values()), len(reaction_similarity)

```

Many reactions connected by promiscuous compounds are given near 0 weights.

```

1 from local.figures.template import BaseFigure, ApplyTemplate, go
2 fig = BaseFigure()
3 fig.add_trace(
4     go.Histogram(
5         x=list(reaction_similarity.values()),
6         nbinsx=100,
7     )
8 )
9 fig = ApplyTemplate(fig, layout=dict(height=200))
10 fig.show()

```

Run UMAP to obtain the network projection.

```

1 from umap import UMAP
2
3 pdist = np.ones((len(r2c), len(r2c)))*1e6
4 r2i = {r: i for i, r in enumerate(r2c)}
5 for (r1, r2), w in edges.items():
6     i, j = r2i[r1], r2i[r2]
7     w = 1/w-1
8     pdist[i, j] = w
9     pdist[j, i] = w
10
11 model = UMAP(
12     n_components=2, metric='precomputed',
13 )
14 rpos = model.fit_transform(pdist)
15 rpos.shape

```

We can get the positions of compounds in the projection by averaging the positions of associated reactions. A jitter is added to prevent compounds from overlapping reactions.

```

1 def polar2cartesian(rho, phi):
2     x = rho * np.cos(phi)

```

```

3     y = rho * np.sin(phi)
4     return np.array([x, y])
5
6     cpos = np.zeros(shape=(len(c2r), 2))
7     c2i = {c: i for i, c in enumerate(c2r)}
8     for c in c2r:
9         _f = list({r2i[r] for r in c2r[c]}) # mask of associated reactions
10        _pos = np.mean(rpos[_f], axis=0)
11        if rpos[_f].std(axis=0).max() < 1e-1: # if the reactions are too close, add some jitter
12            jitter = polar2cartesian(1/50, 2*np.pi*np.random.rand())
13            _pos += jitter
14        i = c2i[c]
15        cpos[i] = _pos
16    cpos.shape

```

For network edges, we will use the inverse of the compound promiscuity.

```

1    eweights = []
2    compound_has_edge = set()
3    for r in r2c:
4        for c in r2c[r]:
5            eweights.append((r, c, 1/compound_promiscuity[c]))
6            compound_has_edge.add(c)
7    e2i = {e: i for i, e in enumerate(eweights)}

```

```

1    fig = BaseFigure()
2    fig.add_trace(
3        go.Histogram(
4            x=[w for _, _, w in eweights],
5        )
6    )
7    fig = ApplyTemplate(fig, layout=dict(height=200))
8    fig.show()

```

For each reaction node, we will include additional information, starting with their parent pathway and GSEA results.

```

1    df_gsea = pd.read_csv("./gsea/results.csv")
2    print(df_gsea.shape)
3    df_gsea.head(2)

```

We will create a lookup table for when we are adding this information to the nodes. Since a reaction may belong to multiple pathways, we will also create a ranking and take the most relevant as indicated by the GSEA results.

```

1 df_gsea = df_gsea.sort_values(["pvalue", "score"], ascending=[True, False])
2 gsea_lookup = {r["set_name"]:r for _, r in df_gsea.iterrows()}
3 pathway2rank = {p:i for i, p in enumerate(df_gsea["set_name"])}
4 df_gsea.head(2)

```

We can also add the differential abundance results from DESeq2.

```

1 df_ds2 = pd.read_csv("./deseq2/results.csv", index_col=0)
2 df_ds2.index.name = "gene"
3 df_ds2 = df_ds2.reset_index()
4 print(df_ds2.shape)
5 df_ds2.head(2)

```

Again, we create a lookup table for when this information is needed. Since each reaction may be the result of multiple genes, we will also create a ranking to determine the most relevant gene as indicated by the DESeq2 results.

```

1 df_ds2 = df_ds2.sort_values(["pvalue", "log2FoldChange"], ascending=[True, False])
2 ds2_lookup = {r["gene"]:r for _, r in df_ds2.iterrows()}
3 gene2rank = {g:i for i, g in enumerate(df_ds2["gene"])}
4 df_ds2.head(2)

```

We will use the metabolic model to link reactions to pathways. Since the names are not present in the SBML, we will use the local KEGG model.

```

1 groups = gs_model["sbml"]["model"]["groups:listOfGroups"]["groups:group"]
2 r2pathways = {}
3 pathway_names = {}
4 for g in groups:
5     g_id = g["@groups:name"]
6     pathway_id = g_id.replace("map", "rn") # convert the index prefix
7     meta = kegg_model.Pathways.get(pathway_id)
8     name = meta.name if meta else g_id
9     pathway_names[pathway_id] = name
10    for r in find(g, "@groups:idRef"):
11        r2pathways[r] = r2pathways.get(r, [])+[pathway_id]
12
13 len(r2pathways)

```

Let's start to construct the cytoscape save file, starting with the nodes.

```

1 i = 0 # a counter to assign unique ids to nodes. increment for each node
2 nodes = [] # we will store the nodes here
3 n2i = {} # this lookup will convert

```

```

4  for _n, _pos, _type in [      # add reactions and compounds
5      (r2c, rpos, "reaction"),
6      (c2r, cpos, "compound"),
7  ]:
8      for k, (x, y) in zip(_n, _pos):
9          n2i[k] = i
10
11         if _type == "reaction":
12             pathways = r2pathways.get(k, [])
13             pathways = sorted(pathways, key=lambda p: pathway2rank.get(p, 1e6))
14             relavent_pathway = pathways[0] if pathways else None
15             pathway_name = pathway_names.get(relavent_pathway, relavent_pathway)
16
17             gsea_meta = gsea_lookup.get(relavent_pathway, None)
18             gsea_score = gsea_meta["score"] if gsea_meta is not None else None
19             gsea_pvalue = gsea_meta["pvalue"] if gsea_meta is not None else None
20
21             genes = r2genes.get(k, [])
22             genes = [g[2:] for g in genes]
23             genes = sorted(genes, key=lambda g: gene2rank.get(g, 1e6))
24             relavent_gene = genes[0] if genes else None
25
26             ds_meta = ds2_lookup.get(relavent_gene, None)
27             ds_lfc = ds_meta["log2FoldChange"] if ds_meta is not None else None
28             ds_pvalue = ds_meta["pvalue"] if ds_meta is not None else None
29
30         else: # for compounds, we can leave these fields blank
31             relavent_pathway = None
32             pathway_name = None
33             gsea_score = None
34             gsea_pvalue = None
35             relavent_gene = None
36             ds_lfc = None
37             ds_pvalue = None
38
39         nodes.append(dict(
40             id=i, x=float(x)*1000, y=float(y)*1000, # numeric id and positiion
41             v=dict( # node attributes
42                 name=names.get(k, k),
43                 type=_type,
44                 gene=relavent_gene,
45                 gene_lfc=ds_lfc,
46                 gene_pvalue=ds_pvalue,
47                 pathway=relavent_pathway,

```



```

48         pathway_name=pathway_name,
49         pathway_score=gsea_score,
50         pathway_pvalue=gsea_pvalue,
51     )
52 ))
53 i += 1

```

Create a collection of edges, indicating the links for nodes. The only attribute added is the edge weight.

```

1 edges = []
2 for j, (r, c, w) in enumerate(eweights):
3     edges.append(dict(
4         id=j, s=n2i[r], t=n2i[c],
5         v=dict(w=float(w)),
6     ))

```

Additional parameters are obtained from creating a simple network with Cytoscape and saving it. The only modification is registering the node and edge attributes we created above.

```

1 META = [
2     {
3         "CXVersion": "2.0",
4         "hasFragments": False
5     },
6     {
7         "networkAttributes": [
8             {
9                 "name": "gen",
10                "description": "gen from script"
11            }
12        ]
13    },
14    {
15        "attributeDeclarations": [
16            {
17                "networkAttributes": {
18                    "name": {
19                        "d": "string"
20                    },
21                    "description": {
22                        "d": "string"
23                    }
24                }
25            }
26        ]
27    }
28 ]

```

```

24         },
25         "nodes": {k: {"d": t} for k, t in [
26             ("name", "string"),
27             ("type", "string"),
28             ("gene", "string"),
29             ("gene_lfc", "double"),
30             ("gene_pvalue", "double"),
31             ("pathway", "string"),
32             ("pathway_name", "string"),
33             ("pathway_score", "double"),
34             ("pathway_pvalue", "double"),
35         ]},
36         "edges": {
37             "w": {
38                 "d": "double"
39             }
40         }
41     },
42 ]
43 },
44 ]
45
46 VISUAL = [
47     {
48         "visualProperties": [
49             {
50                 "default": {
51                     "network": {
52                         "NETWORK_BACKGROUND_COLOR": "#FFFFFF"
53                     },
54                     "edge": {
55                         "EDGE_LABEL": "",
56                         "EDGE_LABEL_COLOR": "#000000",
57                         "EDGE_LABEL_FONT_FACE": {
58                             "FONT_FAMILY": "serif",
59                             "FONT_STYLE": "normal",
60                             "FONT_WEIGHT": "normal"
61                         },
62                         "EDGE_LABEL_FONT_SIZE": 12,
63                         "EDGE_LABEL_OPACITY": 1,
64                         "EDGE_LABEL_ROTATION": 0,
65                         "EDGE_LABEL_MAX_WIDTH": 100,
66                         "EDGE_LINE_STYLE": "solid",
67                         "EDGE_OPACITY": 1,

```

```

68         "EDGE_SELECTED_PAINT": "red",
69         "EDGE_SOURCE_ARROW_COLOR": "#000000",
70         "EDGE_SOURCE_ARROW_SHAPE": "none",
71         "EDGE_LINE_COLOR": "#000000",
72         "EDGE_TARGET_ARROW_COLOR": "#000000",
73         "EDGE_TARGET_ARROW_SHAPE": "none",
74         "EDGE_VISIBILITY": "element",
75         "EDGE_WIDTH": 1,
76         "EDGE_Z_LOCATION": 0
77     },
78     "node": {
79         "NODE_BORDER_COLOR": "#000000",
80         "NODE_BORDER_STYLE": "solid",
81         "NODE_BORDER_OPACITY": 1,
82         "NODE_BORDER_WIDTH": 1,
83         "NODE_BACKGROUND_COLOR": "#FFFFFF",
84         "NODE_HEIGHT": 40,
85         "NODE_LABEL": "",
86         "NODE_LABEL_COLOR": "#000000",
87         "NODE_LABEL_FONT_FACE": {
88             "FONT_FAMILY": "serif",
89             "FONT_STYLE": "normal",
90             "FONT_WEIGHT": "normal"
91         },
92         "NODE_LABEL_FONT_SIZE": 12,
93         "NODE_LABEL_OPACITY": 1,
94         "NODE_LABEL_POSITION": {
95             "HORIZONTAL_ALIGN": "center",
96             "VERTICAL_ALIGN": "center",
97             "HORIZONTAL_ANCHOR": "center",
98             "VERTICAL_ANCHOR": "center",
99             "JUSTIFICATION": "center",
100            "MARGIN_X": 0,
101            "MARGIN_Y": 0
102        },
103        "NODE_LABEL_ROTATION": 0,
104        "NODE_LABEL_MAX_WIDTH": 100,
105        "NODE_BACKGROUND_OPACITY": 1,
106        "NODE_SELECTED_PAINT": "yellow",
107        "NODE_SHAPE": "ellipse",
108        "NODE_VISIBILITY": "element",
109        "NODE_WIDTH": 40,
110        "NODE_Z_LOCATION": 0
111    }

```

```

112         },
113         "nodeMapping": {},
114         "edgeMapping": {}
115     }
116 ]
117 },
118 {
119     "nodeBypasses": []
120 },
121 {
122     "edgeBypasses": []
123 },
124 {
125     "visualEditorProperties": [
126         {
127             "properties": {
128                 "nodeSizeLocked": False,
129                 "arrowColorMatchesEdge": False
130             }
131         }
132     ]
133 },
134 ]

```

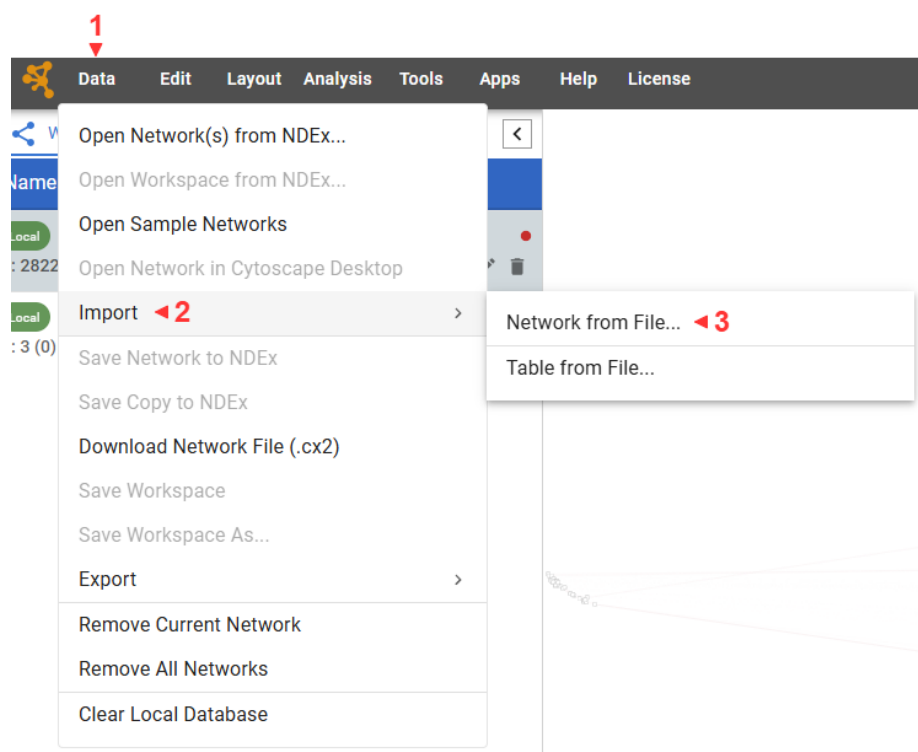
We can now create the Cytoscape save file and load it into Cytoscape.

<https://web.cytoscape.org>

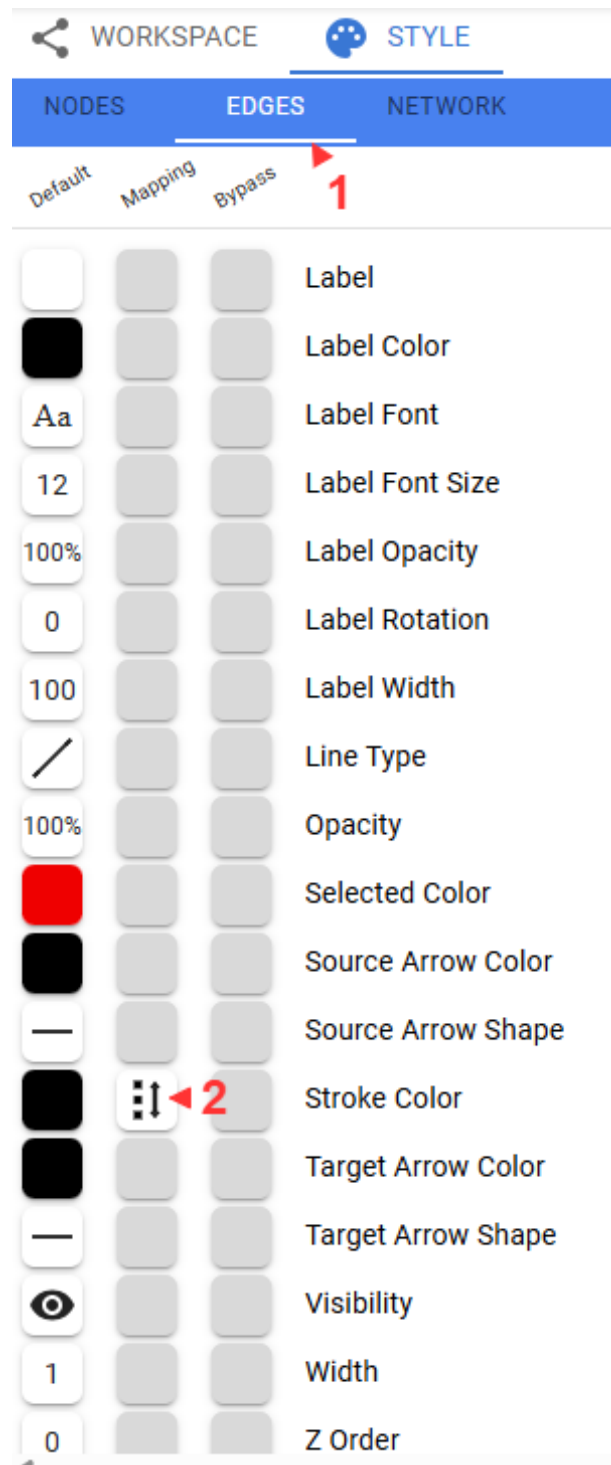
```

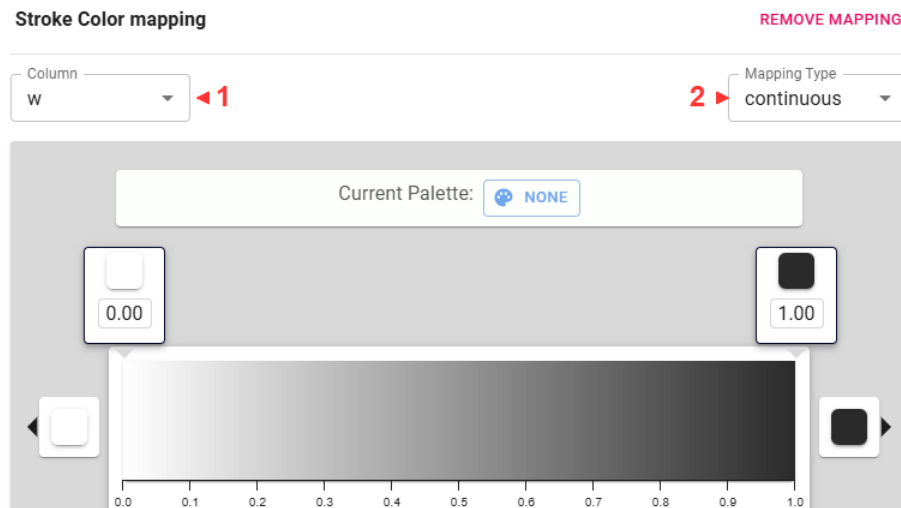
1 import json
2
3 Path("./cytoscape").mkdir(exist_ok=True, parents=True)
4 with open(f"./cytoscape/{MAG}.cx2", "w") as f:
5     json.dump(HEADER+[dict(nodes=nodes), dict(edges=edges)]+VISUAL, f, indent=4)

```



Let's use the edge weights to fade the less important edges.





Some useful navigational tips:

- Scroll to zoom.
- Click and drag to pan.
- Nodes can be dragged around
- Multiple nodes can be selected by pressing shift, clicking, and dragging a box.
- Node properties are shown in a table at the bottom

Play with mapping other values to different visual aspects of the network. Save your network with **Data -> Download Network File (.cx2)**

This concludes the lab!