

```
In [72]: 1 # Packages
2 # from matplotlib import (pyplot as plt, lines)
3 # import seaborn as sns
4 import numpy as np
5 import os
6 import pandas as pd
7
```

```
In [63]: 1 # Globals
2 BASE_DIR = "/home/josh/PycharmProjects/eces-450/tutorial/data/algae"
3
4
5 # Given files
6 print("Assembly Dir:")
7 print('\n'.join(file for file in os.listdir(os.path.join(BASE_DIR,
8 print("\nRead Dir:")
9 print('\n'.join(file for file in os.listdir(os.path.join(BASE_DIR,
```

Assembly Dir:
simple.contig.fa

Read Dir:
CSJP002C_R1.fastq
CSJP002A_R2.fastq
CSJP002B_R2.fastq
CSJP002A_R1.fastq
CSJP002B_R1.fastq
CSJP002C_R2.fastq

```

In [73]: 1 # Read Assembly File
2 fn = os.path.join(BASE_DIR, 'assembly', 'simple.contig.fa')
3 with open(fn, "r") as fh:
4     lines = fh.readlines()
5
6 # Create list: contig_lengths
7 contig_lengths = []
8 contigs = []
9 i = 0
10 for line in lines:
11     if line[0] == '>':
12         contigs.append(line[1:]) # grab record contig id
13     elif i < 1000:
14         contig_lengths.append(len(line)) # grab record sequence length
15     if i<2:
16         print(line[0:250], end='') # print the first record
17         i+=1
18
19 # Plot the sequence lengths
20 plt.style.use('seaborn-dark-palette')
21 fig = plt.figure(figsize=(15, 5))
22 plt.hist(contig_lengths, bins=100, log=True)
23 plt.title("Length Distribution of Assembly Contigs")
24 plt.xlabel("Contig Length")
25 plt.ylabel("Count")

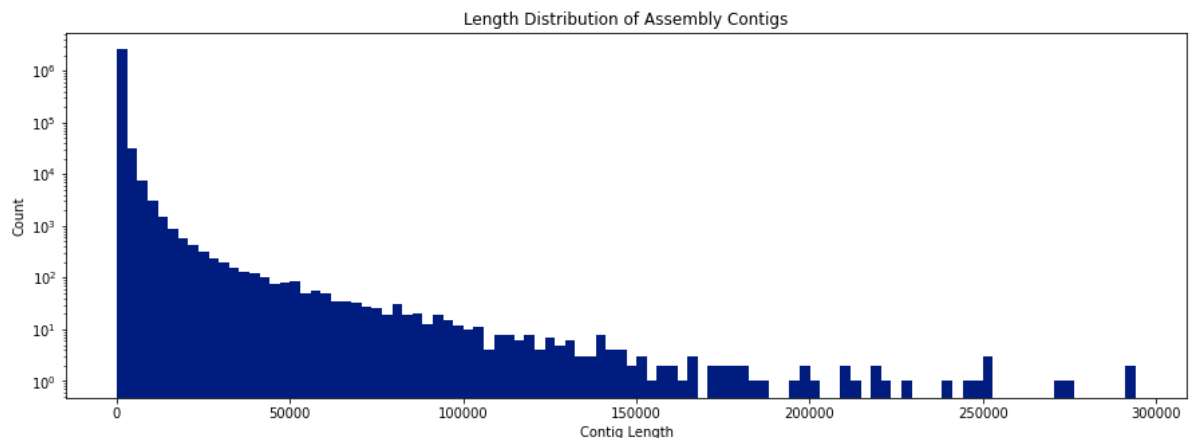
```

```

>contig-65_0
TGGCAGGCGATGTCTGGAGCGCAAACCTCCGCCGCCAGCGCGTCATCTTCTACATCTCACGTGGGTTGGAGA
ACGCCCCGACAAAGTTCCAGACCACCCCCGCCGTCGCGACTACTACAACACGCTGCTCGACCAGCTTCA
GACCGAGTTCGAGCACGTTTCATCGGCTGCTGCTGGAATCTTCGGTTCAGGCGGCGCGCGCGTGTGCCG
ATGACGGATGCCGACCACTTCCGGCACTACAAACGGTTTC

```

Out[73]: Text(0, 0.5, 'Count')



bwa index

Usage: bwa index [options]

Options: -a STR BWT construction algorithm: bwtsv, is or rb2
[auto]

-p STR prefix of the index [same as fasta name]

```

        -b INT      block size for the bwts algorithm (effective
e with -a bwts) [10000000]
        -6          index files named as .64.* instead of .*

```

Source: bwa man pages

A bit of googling and I found:

.amb is text file, to record appearance of N (or other non-ATG C) in the ref fasta.

.ann is text file, to record ref sequences, name, length, etc.

.bwt is binary, the Burrows-Wheeler transformed sequence.

.pac is binary, packaged sequence (four base pairs encode one byte).

.sa is binary, suffix array index.

Source: <http://seqanswers.com/forums/showthread.php?t=25553> (<http://seqanswers.com/forums/showthread.php?t=25553>)

```

In [74]: 1 # After running bwa index on the assembly file, several new files are
          2 print("Assembly Dir:")
          3 print('\n'.join(file for file in os.listdir(os.path.join(BASE_DIR,

```

```

Assembly Dir:
simple.contig.fa.amb
simple.contig.fa.ann
simple.contig.fa.pac
simple.contig.fa.sa
simple.contig.fa.bwt
simple.contig.fa

```

```

1 bwa mem
2
3 The BWA-MEM algorithm performs local alignment. It may produce
  multiple primary alignments for different part of a query
  sequence. This is a crucial feature for long sequences. However,
  some tools such as Picard's markDuplicates does not work with
  split alignments. One may consider to use option -M to flag
  shorter split hits as secondary.
4

```

```
In [ ]: 1 #!/bin/bash
2 ##### Create a map from the reads to the newly indexed assembly-file
3 ##### This took 10.5 hours to complete on proteus
4
5 BASE_DIR="./"
6 samples=(2A 2B 2C)
7
8 for sample in ${samples[@]}
9 do
10     echo CSJP00${sample}_R1.fastq
11     bwa mem ${BASE_DIR}assembly/simple.contig.fa ${BASE_DIR}reads/${sample}_R1.fastq > ${BASE_DIR}reads/${sample}.bam
12 done
```

```
In [4]: 1
2 print("\nMapped:")
3 print('\n'.join(file for file in os.listdir(os.path.join(BASE_DIR,
```

Mapped:
2C.bam
2A.bam
2B.bam

Usage: samtools sort [options...] [in.bam]

Description:

Sort alignments by leftmost coordinates, or by read name when -n is used. An appropriate @HD-SO sort order header tag will be added or an existing one updated if necessary.

The sorted output is written to standard output by default, or to the specified file (out.bam) when -o is used. This command will also create temporary files tmprefix.%d.bam as needed when the entire alignment data cannot fit into memory (as controlled via the -m option).

Options:

- l INT Set compression level, from 0 (uncompressed) to 9 (best)
- m INT Set maximum memory per thread; suffix K/M/G recognized [768M]
- n Sort by read name
- t TAG Sort by value of TAG. Uses position as secondary index (or read name if -n is set)
- o FILE Write final output to FILE rather than standard output
- T PREFIX Write temporary files to PREFIX.nnnn.bam
- no-PG do not add a PG line
- input-fmt-option OPT[=VAL]

```

Specify a single input file format option in the
form
    of OPTION or OPTION=VALUE
-0, --output-fmt FORMAT[,OPT[=VAL]]...
    Specify output format (SAM, BAM, CRAM)
--output-fmt-option OPT[=VAL]
    Specify a single output file format option in the
form
    of OPTION or OPTION=VALUE
--reference FILE
    Reference sequence FASTA FILE [null]
-@, --threads INT
    Number of additional threads to use [0]
--verbosity INT
    Set level of verbosity

```

```

-----
-----
samtools sort -o ./sorted/2A.sorted.bam ./mapped/2A.bam
samtools sort -o ./sorted/2B.sorted.bam ./mapped/2B.bam
samtools sort -o ./sorted/2C.sorted.bam ./mapped/2C.bam
-----
-----

```

```

In [48]: 1 # Sort the bam files for rapid processing, can also be run on proteu
2
3 print("\nSorted:")
4 print('\n'.join(file for file in os.listdir(os.path.join(BASE_DIR, 's

```

```

Sorted:
2B.sorted.bam
2C.sorted.bam
2A.sorted.bam

```

```

1 [js3973@proteusa01 Tutorial6_data]$ ls -al mapped
2 total 23896668
3 drwxrwsr-x 2 js3973 rosenclassGrp      4096 May 18 10:35 .
4 drwxrwsr-x 5 sk3389 rosenclassGrp      4096 May 18 10:32 ..
5 -rw-r--r-- 1 js3973 rosenclassGrp 5866752241 May 18 03:34 2A.bam
6 -rw-r--r-- 1 js3973 rosenclassGrp 6771359549 May 18 07:16 2B.bam
7 -rw-r--r-- 1 js3973 rosenclassGrp 5346570250 May 18 10:13 2C.bam
8 -rw-r--r-- 1 js3973 rosenclassGrp  581851433 May 18 10:37
  2C.sorted.bam
9 -rw-r--r-- 1 js3973 rosenclassGrp  361791128 May 18 10:26
  2C.sorted.bam.tmp.0000.bam
10 -rw-r--r-- 1 js3973 rosenclassGrp  361769335 May 18 10:27
  2C.sorted.bam.tmp.0001.bam
11 -rw-r--r-- 1 js3973 rosenclassGrp  363248908 May 18 10:27
  2C.sorted.bam.tmp.0002.bam

```

```

12 -rw-r--r-- 1 js3973 rosenclassGrp 364979550 May 18 10:28
    2C.sorted.bam.tmp.0003.bam
13 -rw-r--r-- 1 js3973 rosenclassGrp 362689520 May 18 10:28
    2C.sorted.bam.tmp.0004.bam
14 -rw-r--r-- 1 js3973 rosenclassGrp 360639096 May 18 10:29
    2C.sorted.bam.tmp.0005.bam
15 -rw-r--r-- 1 js3973 rosenclassGrp 363002478 May 18 10:30
    2C.sorted.bam.tmp.0006.bam
16 -rw-r--r-- 1 js3973 rosenclassGrp 363244902 May 18 10:30
    2C.sorted.bam.tmp.0007.bam
17 -rw-r--r-- 1 js3973 rosenclassGrp 355685132 May 18 10:31
    2C.sorted.bam.tmp.0008.bam
18 -rw-r--r-- 1 js3973 rosenclassGrp 357044755 May 18 10:32
    2C.sorted.bam.tmp.0009.bam
19 -rw-r--r-- 1 js3973 rosenclassGrp 359320772 May 18 10:32
    2C.sorted.bam.tmp.0010.bam
20 -rw-r--r-- 1 js3973 rosenclassGrp 358477722 May 18 10:33
    2C.sorted.bam.tmp.0011.bam
21 -rw-r--r-- 1 js3973 rosenclassGrp 355302152 May 18 10:33
    2C.sorted.bam.tmp.0012.bam
22 -rw-r--r-- 1 js3973 rosenclassGrp 358575023 May 18 10:34
    2C.sorted.bam.tmp.0013.bam
23 -rw-r--r-- 1 js3973 rosenclassGrp 360891810 May 18 10:35
    2C.sorted.bam.tmp.0014.bam
24

```

Create Depth Matrix:

First generate depth matrix from sorted bam files using `jgi_summarize_bam_contig_depths` (included with `metabat2`)

```
In [86]: 1 # Read Depth Matrix
2 fn = os.path.join(BASE_DIR, 'depth', 'depth_matrix.tab')
3 with open(fn, 'r') as fh:
4     df = pd.read_csv(fh, delimiter='\t')
5 df.head(50) # Show
```

17	contig-65_17	200413	30.7011	0.036851	0.048552	10.963000
18	contig-65_18	199480	29.6988	0.029925	0.036689	10.645100
19	contig-65_19	198319	23.7893	0.028052	0.044290	13.478700
20	contig-65_20	195195	16.9654	0.006224	0.008799	15.488800
21	contig-65_21	185758	22.6358	0.006659	0.010197	19.639300
22	contig-65_22	182600	23.3662	0.013412	0.020377	13.125900
23	contig-65_23	182193	14.7668	0.017331	0.025937	14.648100

Input to metabat2:
sorted bam files
depth matrix

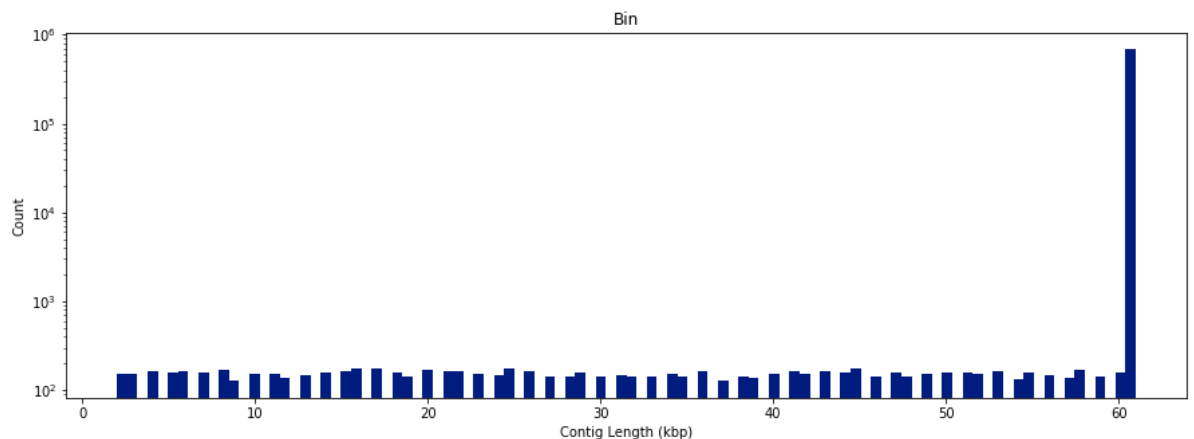
Output:
130 bins

```

In [85]: 1 # Visualize Contig lengths in largest bin
2 fn = os.path.join(BASE_DIR, 'bins', "bin.10.fa")
3 # fn = os.path.join(BASE_DIR, 'bins', "bin.22.fa")
4 with open(fn, "r") as fh:
5     lines = fh.readlines()
6
7 # Create list: contig_lengths
8 contig_lengths = []
9 for line in lines:
10     if not line[0] == '>':
11         contig_lengths.append(len(line))
12
13 # Plot the sequence lengths
14 plt.style.use('seaborn-dark-palette')
15 fig = plt.figure(figsize=(15, 5))
16 plt.hist(contig_lengths, bins=100, log=True)
17 plt.title("Bin")
18 plt.xlabel("Contig Length (kbp)")
19 plt.ylabel("Count")

```

Out[85]: Text(0, 0.5, 'Count')



```

In [80]: 1 # Binned data
2 fn = os.path.join(BASE_DIR, 'bins', "bin.10.fa")
3 with open(fn, "r") as fh:
4     lines = fh.readlines()
5
6 # Create list: contigs
7 contigs = []
8 for line in lines:
9     if line[0] == '>':
10         contigs.append(line[1:].strip())
11

```