Evaluating Metagenomic Assemblies Using Gene Metrics Derived from Protein Coding Sequences

- ⁴ Felix Andrade May¹
- ⁵ ¹Department of Computer Science, Aberystwyth University
- 6 Corresponding author:
- 7 Felix Andrade May¹
- 8 Email address: f.andrademay@gmail.com

ABSTRACT

The High-C binning metagenomic assembly technique sometimes produces collapsed repeats and the 10 loss of complex regions of DNA that do not assemble well. This creates less gene dense Metagenome-11 Assembled Genomes (MAGs), which are harder to accurately classify (Stewart et al., 2018). A classifi-12 cation accuracy of 100% makes the difficult task of identifying genetic samples from within the rumen 13 easier and circumvents the need to culture bacteria before sequencing. By developing new methods 14 to analyse the ability of metagenomic assembly techniques it becomes easier to remove poor MAGs 15 from classification training data and improve the classification accuracy. This paper develops and tests a 16 pipeline for calculating gene metrics, primarily gene density, on consumer grade hardware that could be 17 deployed in crowdsourced citizen science projects. Using four MAGs produced by Stewart et al. (2018), 18 the pipeline can be used to evaluate the High-C binning method of metagenomic assembly used to 19 assemble the four MAGs. The proposed method was able to calculate gene counts and densities in line 20 with expected values for bacterial genomes. When tested on the RefSeq Saccharomyces Cerevisiae 21 genome Engel et al. (2022) it was able to accurately calculate the gene density to within $\pm 40.5 \ gMBP^{-1}$ 22 (genes per Mega Base Pairs). This pipeline confirms that the High-C binning assembly technique is 23 capable of producing accurate MAGs. The ability of the pipeline can be easily improved by selecting a 24 harsher threshold for the e-value of alignments produced by BLAST+ search alignments. 25

26 INTRODUCTION

Biological environments such as the rumen of cattle contain large numbers of organisms that are difficult to culture in a laboratory environment. As a result, samples must be collected directly from the rumen and then assembled using metagenomic techniques like High-C binning. If these samples are able to be easily and accurately metagenomically assembled then it would improve our ability to sequence, analyse, and classify more genomes and employ the knowledge gained from those processes in developing novel medical techniques.

It is hard to assess the quality of metagenomic assembly because the generated Metagenome-33 Assembled Genomes (MAGs) cannot be easily aligned against a reference genome using a BLAST+ 34 search algorithm (Camacho et al., 2009) as the MAG may contain the genome for a prokaryote, the host 35 cattle, or a combination of the two. A solution would be to identify Open Reading Frames (ORFs) and 36 compare those to ORFs that are known to produce proteins. These verified ORFs can then be used to 37 approximate a gene density for the genome. This density can be used to evaluate the MAG. This paper lays 38 out a pipeline to solve this problem. The pipeline works to: identify possible ORFs; align the ORFs using 39 BLAST+ against a protein database (The UniProt Consortium, 2020); filter the results of the alignment 40 to exclude any alignments with an e-value beneath a given threshold; calculate an approximation for 41 the number of genes in, and therefore the gene density of, the genome. This pipeline provides a novel 42 technique that, with further development, could be used on a wider scale in crowdsourced citizen science 43 projects as it is capable of running smoothly on consumer grade hardware running Unix style operating 44 systems. 45

⁴⁶ This paper will use the gene density and related metrics calculated from four of the genomes (RUG005,

47 RUG154, RUG431, RUG466) from the Stewart et al. (2018) to analyse the accuracy of the High-C binning

⁴⁸ method used to produce 913 Metagenome-Assembled Genomes. I will compare the four MAGs to the

49 S. Cerevisiae genome released by Engel et al. (2022), which will be processed using the same pipeline

⁵⁰ described in Section . It will then be possible to compare the result of this pipeline with known metrics of

other prokaryotes and validate the ability of both the proposed pipeline and the High-C binning technique

⁵² used by Stewart et al. (2018).

53 MATERIALS AND METHODS

To calculate the gene density of a FastA file with nucleotide sequences from an unknown genome, the FastA must first be processed to identify potential ORFs. These ORFs are then aligned with a BLAST+ search against the Uni-Prot Swiss-Prot protein database (The UniProt Consortium, 2020). The Pairwise format output from the BLAST+ search is processed with several proprietary Python scripts that refine and filter the results. A Python file is then used to calculate the approximate number of genes for each given genome, allowing for the calculation of gene density.

The Prodigal (Hyatt et al., 2010) v2.6.3 Command Line program was used to identify unique ORFs from each genome, creating a corresponding .gff3 file. Used in combination with Bedtools (Quinlan and Hall, 2010) v2.31.0 these ORFs are extracted from the original genome FastA file into a new FastA file containing only the identified ORFs.

A BLAST+ search was performed on the UniProt Swiss-Prot Database (The UniProt Consortium, 64 2020) to compare these identified ORFs to known protein Coding Sequences (CDSs). The program found 65 in Appendix C, Listing 6, is used to extract just the query names and the lists of significant alignments 66 and write them to a new file. The code found in Appendix B filters out any matches that have e-values 67 greater than a given value (1e-5). The filtering of e-values is performed after the BLAST+ search rather 68 than using a threshold during the BLAST+ search. This is because during a BLAST+ search filtering does 69 not occur during the final stage of the search but rather during the earlier scanning stage of the search 70 algorithm (Camacho et al., 2009). 71 72 A method described by Santos-Magalhaes and de Oliveira (2015) provides a way to approximate

the number of genes in a given genome based on the average protein length expected for a genome. Equation 2 represents a modification to an equation proposed by Santos-Magalhaes and de Oliveira (2015), substituting the average length of amino acid residues in bacteria, \bar{L} . Santos-Magalhaes and de Oliveira (2015) proposed 300 as an average value for bacterial genomes. This pipeline instead

racalculates an average length of amino acid CDS directly from the genome. This provides a more accurate

⁷⁸ approximation for the number of genes in a genome.

$$R = \frac{\text{Length of Identified ORFs (BP)}}{\text{Length of Genome (BP)}} \quad (1) \qquad g = \frac{\text{Length of Genome (BP)}}{\left(\frac{\text{Average Length Of Protein CDS}}{R}\right)} \quad (2)$$

Equation 1 describes the ratio between the number of Base Pairs (BPs) in an identified ORF and the number of Base Pairs in the genome. This ratio is used in Equation 2 to calculate an approximate number of genes (g) in the genome, given an average length of protein CDS measured in BP.

The code found in D describes a program which takes as input the genome FastA file, the FastA file of ORFs generated by Prodigal and Bedtools, and the Blast+ output processed by the programs in Appendix B and Appendix C. It uses these inputs to identify ORFs which have been significantly aligned to at least one protein CDS. The program then calculates an average length for all identified protein CDS and applies this information to Equations 1 and 2 to calculate an approximate number of genes for the genome FastA file.

The approximate number of genes can then be divided by the total length of all sequences in the genome FastA file measured in Mega Base Pairs to calculate an approximate gene density, δ (Equation 3).

$$\delta = \frac{g}{\text{Length of genome (MBP)}} \tag{3}$$

Appendix E describes examples of commands needed in order to run the pipeline. These commands are presented in the appropriate order.

94 RESULTS

⁹⁵ Table 1 shows the total number of significant alignments that obtained an *e*-value lower than a threshold

value of 1e-5 and the number of nonsignificant alignments which scored higher than than the threshold

value. The genomes with shorter overall lengths of ORFs (RUG005 and RUG466) see a higher percent of

significant alignments compared to the others, although the there is no significant difference between the

genomes. As shown in Table 2, RUG005 and RUG466 also had shorter ORFs compared to RUG154 and RUG431, with the ORFs identified in RUG005 constituting the smallest percent of the genome (80.75%).

Genome	RUG005	RUG154	RUG431	RUG466
Significant Alignments	251,004	309,417	230,145	229,716
Nonsignificant Alignments	85,190	116,488	97,314	84,401
Percent Significant / %	74.66	72.65	70.28	73.13

Table 1. The number and percentage of alignments with an e-value of <1e-5.

100

103

Genome	RUG005	RUG154	RUG431	RUG466
Total Genome Length / BP	2,345,410	2,774,557	2,770,634	1,925,880
Total ORF Length / BP	1,893,858	2,504,499	2,562,720	1,711,224
Percent ORFs / %	80.75	90.27	92.50	88.85

Table 2. The length of ORFs and wha	percentage of the original	genome file they make up.
-------------------------------------	----------------------------	---------------------------

Table 3 lists the average protein CDS length and the approximate gene count for each genome, calculated by the program found in Appendix D. Combining this information with the Total Genome Length in MBP (Table 2), an approximate gene density is calculated for each genome.

Genome	RUG005	RUG154	RUG431	RUG466
Average Protein CDS Length	1045.54	1060.75	1261.08	1029.03
Approximate Gene Count	1811.36	2361.06	2032.16	1662.95
Approximate Gene density / gMBP ⁻¹	503.69	850.97	733.46	863.48

Table 3. Approximate Gene Count and Gene Density (in Genes per Mega BPs) given the Average length of Protein CDS for each genome.

This pipeline can be verified by processing reference genomes with known metrics. The genome of S. Cerevisiae (Engel et al., 2022) as well as a curated file of Coding Sequences are available through the National Center for Biotechnology Information (NCBI) database (Sayers et al., 2022). The results of processing the full genome with Prodigal and Bedtools to extract ORFs can be compared against the FastA file containing CDS provided by NCBI. Both the ORF and CDS files can be aligned using a BLAST+ search against the Swiss-Prot database. The output of that alignment search can be processed through the Python pipeline described above to gather an approximate gene count and gene density (Table 4).

Due to a limitation of available hardware, the S. Cerevisiae genome FastA file was separated into five sequentially numbered files that were then processed through the gene_calculator.py script in Appendix D. The results of these five operations were then combined and are shown in Table 4. The same

¹¹⁴ process was applied to the S. Cerevisiae CDS FastA file.

115 DISCUSSION

¹¹⁶ Using metrics calculated from the four MAGs and comparing those to metrics of a known genome (S.

117 Cerevisiae), the pipeline calculated the gene density to within $\approx 8\%$ of the expected value. These metrics

allow for poorly assembled MAGs to be removed from classification training data, improving the ability

¹¹⁹ of classification. This in turn makes the process of sampling, assembling, and identifying the MAGs ¹²⁰ easier and more accurate.

One major limitation of this paper is the imprecision of the method used to calculate the number of genes for a genome. The equations put forth by Santos-Magalhaes and de Oliveira (2015) originally use

Genome	S. Cerevisiae	S. Cerevisiae CDS
Significant Alignments	531,043	535,281
Nonsignificant Alignments	275,416	260,454
Percent Significant / %	65.85	67.27
Total Genome Length / BP	12,157,105	8,826,477
Total ORF Length / BP	9,154,275	8,770,452
Percent ORFs / %	75.30	99.37
Average Protein CDS Length	1397.78	1459.62
Approximate Gene Count	6549.14	6008.73
Approximate Gene Density / $gMBP^{-1}$	538.71	494.26 †

Table 4. Metrics produced by the developed pipeline for the S. Cerevisiae genome.

 † Value calculated using the Total Length of the S. Cerevisiae genome, not the Total Length of the S. Cerevisiae CDS genome, which is already a subset of the full S. Cerevisiae genome.

average values for the length of protein CDS of a taxonomic domain. This paper alters these equations to instead utilise the average length of protein CDS of a genome. While this alteration does allow for a more appropriate value it is still hampered by using the average length of protein CDS across the genome. By using the average lengths it is likely that the true number of genes in a genome is not calculated, which makes it harder to accurately calculate the gene density of a genome and harder again to accurately evaluate the High-C binning used by Stewart et al. (2018) to create the MAGs in the first place. This can be seen with the results of processing the S. Cerevisiae genome as seen in Table 4

Table 4 shows the results of processing the S. Cerevisiae genome in the same manner as the genomes from Stewart et al. (2018). The calculated approximate gene count of 6,549.14 is greater than the actual protein producing gene count of 6,014 recorded by Engel et al. (2022). This difference is likely due to the threshold value set during the filtering stage of the process. Decreasing the threshold from an *e*-value of 1²⁴ 1*e*-5 to a value of 1*e*-10 may reduce the number of ORFs that found significant alignments during the BLAST+ search, thus decreasing the number of identified protein producing Coding Sequences. This in turn would reduce the number of approximate genes in a genome.

This larger number of genes will also create a larger than expected gene density. Using the gene count of 6,014 combined with total genome length of 12,071,326 BP we can calculate a more accurate gene density of 498.21 $gMBP^{-1}$, a difference of 40.5 $gMBP^{-1}$ from the approximation produced by the pipeline. This difference of $\approx 8\%$ is large, although it does not invalidate the results generated for the four RUG genomes as it would not significantly alter the genomic density for the RUG genomes.

The NCBI RefSeq database also provides the ability to download a curated FastA file that contains only protein CDS. Table 4 shows the metrics produced by the pipeline when processing this CDS FastA file. The pipeline is able to align 99.37% of the curated protein CDS to the Swiss-Prot protein database (The UniProt Consortium, 2020). The approximate gene density calculated for the CDS FastA of 494.26 falls within the expected margin of error ($\approx 8\%$) for the pipeline with these parameters.

Given an error margin of $40.5 \ gMBP^{-1}$ for the gene densities of the RUG genomes, the gene density of all four RUG genomes falls within the expected range of 500-1,000 $gMBP^{-1}$ for prokaryotic genomes (O'Leary et al., 2016). These gene densities show that the High-C binning method employed by Stewart et al. (2018) has been able to produce genomes similar in gene composition to validated prokaryotic genomes in the RefSeq database O'Leary et al. (2016).

152 CONCLUSION

The method proposed in this paper is capable of calculating the gene density of an unknown genome given the alignment of the genome to a protein database. This allows for the analysis of Metagenome-Assembled Genomes by calculating the gene density of reference genomes, either manually using statistics calculated by organisations like NCBI or by processing the reference genomes through the same pipeline as the

157 unknown genome.

The four MAGs produced by Stewart et al. (2018) were found to be within the expected range of gene density for prokaryotic genomes. This similarity shows the High-C binning technique used to assemble the four MAGs is able to accurately assemble reads from the rumen of cattle into distinct genomes.

ADDITIONAL GENE INVESTIGATIONS

162 Code Overview

The code consists of two Python scripts, part_one.py and similarity.py. They have been written using Python 3.11 and are designed to be run from the command line. Both scripts should be located in the same directory, FastA files may be located elsewhere, as long as the file path passed to the script is correct. Output files may be sent directly to subdirectories of the outputs directory provided the subdirectory already exists.

168 part_one.py

part_one.py should then be run with a single FastA input file, an output filename, and a K value for *K*-Mer composition operations. An example can be found in Appendix A, Listing 1.

This will run the part_one.py script on the genome_1.fa file, with the statistics generated by the script written to ./outputs/genome_1_output.txt. The outputs directory will be created in the directory the command was executed from. Running the example command (Appendix A, Listing 1) will also create a file called genome_1_dict.csv which is a Comma Separated Values (CSV) file containing all unique *K*-Mers found in the corresponding FastA file. Each line will begin with the *K*-Mer

in question followed by the number of occurrences of that *K*-Mer in the provided FastA file.

part_one.py generates the GC-Content Mean, GC-Content Standard Deviation (SD), N50, N90,
and L50 statistics. These are written to the given output file, with the name of the provided FastA file on
line 1, then the GC-Content Mean, GC-Content SD, N50, N90, and L50 statistics, each on a new line.
The file ends on a blank line. An example of the output file can be seen in Appendix A, Listing 2.

181 similarity.py

To compare the Manhattan Distance between any number of genomes, similarity.py will be required. The similarity.py script takes the _dict outputs of part_one.py and uses them to calculate the distance between genomes. The similarity.py script can be passed a filename for the results to be written too, and two or more _dict files and will return a list of each comparison, and the distance between the compared genomes. An example can be found in Appendix A, Listing 3.

The example from Appendix A, Listing 3 uses similarity.py to calculate the distance between: gl_dict and g2_dict; gl_dict and g3_dict; g2_dict and g3_dict. These will be returned to the user as a CSV file in the outputs directory with the name which was given as the first argument in the example command (distances.csv). An example section of this file can be found in Appendix A Listing 4

A, Listing 4.

192 Third Party Libraries

The part_one.py script makes use of Pysam (John Marshall, 2023) to read in the sequence names and sequence strings from a given FastA file. Both scripts also utilise NumPy (Harris et al., 2020) for some operations.

196 Statistics

197 GC-Content

As shown in Table 5, the mean GC-Contents for each genome vary by 9.014%. Comparing the two GC-Content SDs that are most different - RUG005 (2.4090) and RUG154 (1.1743) - with the Hartley's F-Max test gives an F_{Max} value of 4.10. With over one hundred degrees of freedom, the critical value for this comparison would be 1.00. This means that the GC-Content SD between the two most divergent are

statistically heterogenous. This also applies to the RUG466 genome, which is statistically heterogenous

 $_{\rm 203}$ $\,$ from the RUG154 genome but not the RUG005 or RUG431 genomes.

Genome	RUG005	RUG154	RUG431	RUG466
GC-Content Mean / %	48.935	44.288	51.058	57.949
GC-Content SD	2.4090	1.1743	1.9024	2.3720

Table 5. Table showing the GC-Content means and Standard Deviations per genome.

204 Metrics

²⁰⁵ While it is not possible to compare the N50 and L90 metrics found in Table 6, because of the differing

numbers of contigs in each FastA file, it is possible to compare the L50 metrics. The RUG466 genome

has the highest L50 value, as well as the lowest N50 and N90 values. This implies that the RUG466

²⁰⁸ genome has been assembled from shorter and therefore less accurate reads compared to the others. This is ²⁰⁹ consistent with genome having the highest mean average GC-Content (Table 5), which might be expected

consistent with genome having the highest mean average GC-Content (Table 5), which
 given the Hi-C-based proximity-guided assembly employed by Stewart et al. (2018).

Genome	RUG005	RUG154	RUG431	RUG466
N50	46061	43578	52738	16904
N90	12229	15257	22961	5244
L50	12	24	15	31

Table 6. Table showing the N50, N90, and L50 metrics for each genome.

Genome Similarity

Table 7 shows the Manhattan Distance between each of the four genomes. The least similar genomes

are the RUG466 and RUG154 genomes, with a distance of 1,607,019. This mirrors the findings of

- other metrics which, using GC-Content, showed the RUG466 and RUG154 genomes to be statistically
- ²¹⁵ heterogenous.

Genome	RUG005	RUG154	RUG431	RUG466
RUG005	0	541788	547265	628785
RUG154	541788	0	456163	1067019
RUG431	547265	456163	0	908798
RUG466	628785	1067019	908798	0

Table 7. The Manhattan distance between each of the four genomes with regards to 2-Mer frequencies.

216 A ADDITIONAL GENETIC INVESTIGATIONS EXAMPLES

217 \$ python3 part_one.py data/gene_1.fa gene_1_output.txt 2

Listing 1. Example command line input to run part_one.py



225 \$ python3 similarity.py distances.csv gene_1_dict.csv gene_2_dict.csv gene_3_dict.csv
Listing 3. Example command line input to run similarity.py

226	алалала, 494
227	ААААААС, 329
228	AAAAAG,791
229	алалаат, 777
230	алаласа, 347
231	AAAAACC,275
232	AAAAACG,248
233	AAAAACT,434
234	AAAAAGA, 852
235	AAAAAGC, 654
236	AAAAAGG,670
237	AAAAAGT, 424
238	алалата, 625
239	AAAAATC,663
240	AAAAATG,544
241	AAAAATT, 540
242	ААААСАА, 211
243	AAAACAC,94
244	



245 **B FILTERER.PY CODE**

```
246
    import sys
    import re
247
248
249
    def filter_bad_matches(filename):
250
251
        """Removes any alignments from a blastx alignment that are above a given threshold
252
        writes a file with the name of the input file with the '_filtered.txt' suffix
253
254
255
        :param filename: filename of blastx alignments to be filtered by e-value
        :return: good_matches: number of alignments below the threshold value
256
        :return: bad_matches: number of alignments above the threshold"""
257
258
        # alter this float to desired value
259
        threshold = float(1e-5)
260
        # pattern matches the pairwise alignments produced from a blastx search
261
        pattern = re.compile(r"^([A-Z0-9]+) +([^\n]+) +(\d+\.?\d+) +((^\n]+) $", re.
262
263
        MULTILINE)
        good_matches = 0
264
        bad_matches = 0
265
        with open(filename, 'r') as infile:
266
            outfile_name = filename.split('.')[0] + '_filtered.txt'
267
268
            with open(outfile_name, 'w') as outfile:
                 # could be converted to tqdm (or other) progress bar
269
270
                 print("Filtering file ...")
                 first_line = True
271
272
                 for i, line in enumerate(infile):
273
                     if line.startswith('Query='):
                         # we don't want to start the file with a blank line
274
275
                          if first_line:
                             outfile.write(line)
276
                              first_line = False
277
278
                         else:
                              # but we do want subsequent queries to be separated by a blank
279
280
                             outfile.write("\n" + line)
281
282
                     for match in re.finditer(pattern, line):
                          \ensuremath{\texttt{\#}} once we have matched a query line, then write all its good
283
284
285
                          if float(match[4]) < threshold:</pre>
                              outfile.write(match.group() + "\n")
286
287
                              good_matches += 1
                         else:
288
289
                              bad_matches += 1
290
291
        return good_matches, bad_matches
292
293
    if ___name___ == '__
                      _main__':
294
        good, bad = filter_bad_matches(sys.argv[1])
295
        print("good: "+str(good) +
296
                 "\nbad: "+str(bad))
297
        print("Percent Good = "+str('{0:.2f}'.format(good/(good+bad)*100)))
298
299
```

Listing 5. filterer.py code

300 C TRIMMER.PY CODE

```
import sys
301
    import re
302
303
304
    def trim_file(filename):
305
306
        """Trims the Pairwise output file of a blastx search and extracts only the query
307
308
        and the list of significant alignments and writes them to a new file
        :param filename: filename of the alignment file to be trimmed"""
309
310
        pattern = re.compile(r"^([A-Z0-9]{6,})\s+([^{n}]+)\s+([d+\.?\d+)\s+([^{/}])([^{n}]+)$"
311
        , re.MULTILINE)
312
        with open(filename, 'r') as infile:
313
            outfile_name = filename.split('.')[0] + '_trimmed.txt'
314
            with open (outfile_name, 'w') as outfile:
315
                print("Trimming file ...")
316
317
                 first_line = True
                 for i, line in enumerate(infile):
318
                     if line.startswith('Query='):
319
320
                          if first_line:
                             outfile.write(line.strip() + "\n")
321
                              first_line = False
322
323
                         else:
                             outfile.write("\n" + line)
324
325
                     for match in re.finditer(pattern, line):
                         outfile.write(match.group().strip() + "\n")
326
327
        print("File trimmed!")
328
329
330
    if ___name___ == '___main__':
331
332
        trim_file(sys.argv[1])
333
```

Listing 6. trimmer.py code

D GENOMIC_CALCULATOR.PY CODE

```
import sys
335
336
    import pysam
337
    import part_one as po
338
    import re
339
340
    def process_gene_orf_inputs(genome_path, orf_path):
341
         """Calls part_one.py to process the FastA file and the corresponding FastA file of
342
343
344
        :param genome_path: path to the FastA genome file, passed as the first cli
345
        :param orf_path: path to the FastA genome ORF file, passed as the second cli
346
347
348
        :return: genome_sequence_names: list of sequence names from the genome FastA file
        :return: orf_sequence_names: list of ORF sequence names from the ORF FastA file
349
350
351
        :return: orf_file: Pysam object representing the ORF FastA file"""
352
        genome_sequence_names = po.get_sequence_names(genome_path)
353
354
        orf_sequence_names = po.get_sequence_names(orf_path)
        genome_file = pysam.FastaFile(genome_path)
355
356
        orf_file = pysam.FastaFile(orf_path)
357
358
        return genome_sequence_names, orf_sequence_names, genome_file, orf_file
359
360
361
    def process_filtered_orf_inputs(filtered_orf_path):
        """Processes the filtered file of alignments and returns a list of ORFs that had
362
363
364
365
        :param filtered_orf_path: path to the filtered file of alignments produced by
366
367
        :return: protein_producers: list of identified ORFs that had at least one filtered
368
        significant alignment
369
370
        alignments"""
371
372
        pattern = re.compile(r"^([A-Z0-9]+)\s+([^\n]+)\s+(\d+\.?\d+)\s+([^\n]+)$", re.
373
        MULTILINE)
374
        protein_producers = []
375
376
        chance orfs = []
        with open(filtered_orf_path, 'r') as infile:
377
            lines = infile.readlines()
378
            for index, line in enumerate(lines):
379
                 if line.startswith('Query=') and lines[index] != lines[len(lines)-1]:
380
                     next_line = lines[index + 1]
381
382
                     if re.match(pattern, str(next_line)):
                         protein_producers.append(line[7:].strip())
383
384
                     else:
385
                         chance_orfs.append(line[7:].strip())
386
        return protein_producers, chance_orfs
387
388
389
    def get_average_length_of_protein_cds(protein_producers, orf_file, orf_sequence_names)
390
391
        ""Calculates the average length of protein producing ORFs across a genome
392
        :param protein_producers: list of protein producing ORFs, from
393
394
395
396
397
        process_filtered_orf_inputs()
        :return: float value for the average length of protein coding sequences"""
398
399
        length_of_protein_cds = 0
400
401
        print("orf sequences
                                            = " + str(len(orf_sequence_names)))
```

```
print("protein sequences = " + str(len(protein_producers)))
402
403
        for sequence in protein_producers:
            for orf in orf_sequence_names:
404
405
                 if orf.find(sequence) != -1:
                     read = orf_file.fetch(orf)
406
407
                     length_of_protein_cds += len(read)
408
        return length_of_protein_cds / len(protein_producers)
409
410
411
412
    def calc_approx_gene_count(genome_length, average_cds_length, orf_length):
413
        """Calculates an approximate gene count given the number of protein producing
414
        and the total length of all identified ORFs and the total length of the genome in
415
416
        :param genome_length: total length of the genome in base pairs, from
417
418
        get_num_useful_sequences()
        :param average_cds_length: average length of protein CDS, from
419
420
        get_average_length_of_protein_cds()
421
422
        :return g: float value of the approximate gene count"""
423
424
425
        g = (genome_length / (average_cds_length / (orf_length / genome_length)))
426
        return g
427
428
429
    def get_sequence_lengths(genome_sequence_names, orf_sequence_names, genome_file,
430
        orf file):
        """Calculates the total length of ORFs and the total length of the whole genome,
431
432
        both in base pairs
        :param genome_sequence_names: list of sequence names from the genome FastA, from
433
434
        :param orf_sequence_names: list of sequence names from the ORF FastA, from
435
436
437
        :param genome_file: Pysam object representing the genome FastA file, from
        process_gene_orf_inputs()
438
439
        :param orf_file: Pysam object representing the ORF FastA, from
440
        :return genome_length: total length of the genome in base pairs
441
        :return orf_length: total length of all identified ORFs in base pairs"""
442
443
444
        genome_length = 0
        orf_length = 0
445
446
        for sequence_name in genome_sequence_names:
447
448
            read = genome_file.fetch(sequence_name)
449
            genome_length += len(read)
450
        for sequence_name in orf_sequence_names:
451
            read = orf_file.fetch(sequence_name)
452
453
            orf_length += len(read)
454
        return genome_length, orf_length
455
456
457
        _name__ == '__main__
458
                             1 :
        (genome_sequence_names, orf_sequence_names,
459
            genome_file, orf_file) = process_gene_orf_inputs(sys.argv[1], sys.argv[2])
460
461
        genome length, orf length = get sequence lengths (genome sequence names,
462
        orf_sequence_names,
463
                                                               genome_file, orf_file)
464
        \# prints the different metrics, aligned to each other by the '=' sign for
465
466
        # could be modified to output these in a csv file instead (like part_one.
467
468
        write_statistics() does)
        print("total genome length
                                            = " + str(genome_length))
469
470
        print("total orf length
                                            = " + str(orf_length))
471
```

```
filtered_orf_path = sys.argv[3]
472
        protein_producers, chance_orfs = process_filtered_orf_inputs(filtered_orf_path)
473
474
        average_cds_length = get_average_length_of_protein_cds(protein_producers, orf_file
475
        , orf_sequence_names)
476
477
        print("Average protein CDS length = "+str(average_cds_length))
478
        gene_count = calc_approx_gene_count(genome_length, average_cds_length, orf_length)
479
       print("Gene count
                                           = " + str(gene_count))
480
```

Listing 7. genomic_calculator.py code

481 E EXAMPLE PIPELINE COMMANDS

This supplemental file shows an example workflow to follow when using the pipeline. These commands are applicable to the MacOS zsh 2.14 (452) Terminal, but are translatable to other Unix derived operating systems.

485 Third Party Command Line Programs

486 prodigal -i genome_1.fa -o genome_1.gff3 -f gff

Listing 8. Example Prodigal Command

487 bedtools getfasta -fi genome_1.fa -bed genome_1.gff3 -s -fo genome_1.orfs.fa

Listing 9. Example Bedtools Command

488 time blastx -query genome_1.fa -db uniprot_sprot -out genome_1_matches.txt

Listing 10. Example Blast+ Command

- $_{\tt 489}$ $\,$ Inclusion of the time command is optional, but handy when running larger BLAST+ alignments on less
- ⁴⁹⁰ powerful hardware. the uniprot_sprot database is the Swiss-Prot protein database.

491 The Pipeline

- ⁴⁹² The proprietary Python scripts developed for the pipeline must be located in the same level directory, in
- this example they are all located in the / [my_pipeline]/code/ directory. They are all run from the
- 494 top level directory (/[my_pipeline]/).
- 495 python3 code/trimmer.py outputs/genome/genome_1_matches.txt

Listing 11. Example trimmer.py Command

496 python3 code/filterer.py outputs/genome/genome_1_matches_trimmed.txt

Listing 12. Example filterer.py Command

497 python3 code/genomic_calculator.py data/genome/genome_1.fa data/genome/genome_1.orfs. 498 fa outputs/genome/genome_1_matches_trimmed_filtered.txt

Listing 13. Example genomic_calculator.py Command

- ⁴⁹⁹ The genomic_calculator.py script takes three command line arguments. The first is the path from
- the user's present working directory to the FastA file containing the genome. The second is the path to the
- ⁵⁰¹ FastA file of Open Reading Frames generated by Prodigal and Bedtools. The final argument is the path to
- the trimmed and filtered file of alignments generated by the BLAST+ alignment search. The

503 Optional Pipeline Files

504 python3 code/splitter.py data/genome/genome_1.fa

Listing 14. Example splitter.py Command

splitter.py will take a FastA file and divide it into a given number of sequentially numbered FastA 505 files - i.e. genome_1.fa becomes: genome_1_1.fa, genome_1_2.fa. This may be needed if the original 506 FastA genome file is too large to be processed on the user's hardware. Through testing a FastA filesize of 507 < 5MB was desirable. Above this threshold the Pysam third party library may run out of memory and 508 report an unexpected end of file error. If the original FastA file must be processed with the splitter then 509 each sequential FastA file must be run through the genomic_calculator.py individually, although 510 the second and third arguments need not change. 511 python3 code/clean_query.py outputs/genome/genome_1_matches_trimmed_filtered.txt 512

Listing 15. Example clean_query.py Command

 $_{\tt 513}$ The <code>clean_query.py</code> and script should be run if there is a mismatch in the sequence names between

the three different files that are input to the genomic_calculator.py script. clean_query.py

 $_{\tt 515}$ $\,$ takes the processed <code>matches_trimmed_filtered.txt</code> file and renames the sequences to remove

⁵¹⁶ improper characters, such as any "[]" characters and any characters in between these. This is needed as

 $_{\tt 517}$ the <code>genomic_calculator.py</code> must compare the sequence names for some operations. To change

 $_{\tt 518}$ $\,$ what the offending characters are, open <code>clean_query.py</code> with a preferred editor and update the RegEx

519 pattern variable on line 5.

REFERENCES 520

- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T. L. (2009). 521 Blast+: architecture and applications. *BMC Bioinformatics*, 10(1):421. 522
- Engel, S. R., Wong, E. D., Nash, R. S., Aleksander, S., Alexander, M., Douglass, E., Karra, K., Miyasato, 523
- S. R., Simison, M., Skrzypek, M. S., Weng, S., and Cherry, J. M. (2022). New data and collaborations 524
- at the saccharomyces genome database: updated reference genome, alleles, and the alliance of genome 525

- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, 527
- E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M., 528
- Haldane, A., del Río, J. F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., 529
- Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. (2020). Array programming with NumPy. 530 Nature, 585(7825):357-362. 531
- Hyatt, D., Chen, G.-L., LoCascio, P. F., Land, M. L., Larimer, F. W., and Hauser, L. J. (2010). Prodigal: 532 prokaryotic gene recognition and translation initiation site identification. BMC bioinformatics, 11:1-11. 533
- John Marshall, A. H. (2023). Pysam. https://github.com/pysam-developers/pysam. 534
- O'Leary, N. A., Wright, M. W., Brister, J. R., Ciufo, S., Haddad, D., McVeigh, R., Rajput, B., Robbertse, 535
- B., Smith-White, B., Ako-Adjei, D., Astashyn, A., Badretdin, A., Bao, Y., Blinkova, O., Brover, V., 536
- Chetvernin, V., Choi, J., Cox, E., Ermolaeva, O., Farrell, C. M., Goldfarb, T., Gupta, T., Haft, D., 537
- Hatcher, E., Hlavina, W., Joardar, V. S., Kodali, V. K., Li, W., Maglott, D., Masterson, P., McGarvey, 538
- K. M., Murphy, M. R., O'Neill, K., Pujar, S., Rangwala, S. H., Rausch, D., Riddick, L. D., Schoch, C., 539
- Shkeda, A., Storz, S. S., Sun, H., Thibaud-Nissen, F., Tolstoy, I., Tully, R. E., Vatsan, A. R., Wallin, 540
- C., Webb, D., Wu, W., Landrum, M. J., Kimchi, A., Tatusova, T., DiCuccio, M., Kitts, P., Murphy, 541
- T. D., and Pruitt, K. D. (2016). Reference sequence (refseq) database at ncbi: current status, taxonomic 542
- expansion, and functional annotation. Nucleic Acids Res., 44(D1):D733-45. 543
- Quinlan, A. R. and Hall, I. M. (2010). Bedtools: a flexible suite of utilities for comparing genomic 544 features. Bioinformatics, 26(6):841-842. 545
- Santos-Magalhaes, N. S. and de Oliveira, H. M. (2015). Of protein size and genomes. arXiv preprint 546 arXiv:1502.03732. 547
- Sayers, E. W., Bolton, E. E., Brister, J. R., Canese, K., Chan, J., Comeau, D. C., Connor, R., Funk, K., 548 Kelly, C., Kim, S., Madej, T., Marchler-Bauer, A., Lanczycki, C., Lathrop, S., Lu, Z., Thibaud-Nissen, 549
- F., Murphy, T., Phan, L., Skripchenko, Y., Tse, T., Wang, J., Williams, R., Trawick, B. W., Pruitt, K. D., 550
- and Sherry, S. T. (2022). Database resources of the national center for biotechnology information. 551 Nucleic Acids Res., 50(D1):D20–D26.
- Stewart, R. D., Auffret, M. D., Warr, A., Wiser, A. H., Press, M. O., Langford, K. W., Liachko, I., Snelling, 553 T. J., Dewhurst, R. J., Walker, A. W., Roehe, R., and Watson, M. (2018). Assembly of 913 microbial 554
- genomes from metagenomic sequencing of the cow rumen. Nature Communications, 9(1):870. 555
- The UniProt Consortium (2020). UniProt: the universal protein knowledgebase in 2021. Nucleic Acids 556
- Research, 49(D1):D480-D489. 557

resources. Genetics, 220(4). 526