INFO 529

HW1-Section 3

Authors: Chia-Hsuan Chou Venkata Prudhvi Raj Indana Jing Wang

January 29, 2016

Contents

I.	Goal	3
II.	Procedure	3
III.	Result	4
IV.	Appendix	7
V	Author Contributions	9

I. Goal

Our goal is to find potential E.coli genes from a given DNA sequence file based on E.coli codon-usage.

II. Procedure

First, we collected the entire *E.coli* genome from NCBI and selected 1000 gene sequences from the genome as our input (attach in the folder which called "gene_1000.fasta"). Next, we used these gene sequences to compute the codon usage table and build a probabilistic model (see Appendix A) in order to reference it later. Then, we input a given DNA sequence and the program predicted all possible open reading frames (ORFs) in that given DNA sequence. Moreover, it calculated the log-likelihood ratio of $P_i = p_i / p_0$ for every predicted ORF. Here notices that we assume every codon is independent to its adjacent codons so that we can multiply every probability of codon based on the codon-usage dictionary we made and the outcome will be p_i . Also, for the random model of coding DNA, we assume the probability of every codon is 1/64 = 0.0156, so the multiplication of every codon in every predicted ORF will be p_0 . Finally, the result will be the likelihood of every predicted ORF and the program gives us two FASTA files with predicted genes and the translated protein respectively (see Appendix B and Appendix C).

III. Result

```
Chouse

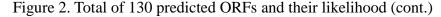
Please find two output files with predicted genes as genes.fasta and protein.fas
ta in the folder.

ORF 1: -87.8205889295
ORF 2: 3.70501803067
ORF 3: -83.352400405
ORF 4: -73.7828603597
ORF 5: -73.065695434
ORF 6: -66.618850132
ORF 7: -65.5764894475
ORF 8: 3.78084879014
ORF 9: -62.2147436207
ORF 10: -61.2682356888
ORF 11: -57.2452264888
ORF 12: -56.9963062655
ORF 13: -56.9963062655
ORF 13: -56.59639312278
ORF 15: -56.9963091278
ORF 15: -56.9963091278
ORF 16: -59.3076503758
ORF 17: -54.2596220311
ORF 16: -6.93.9076503758
ORF 17: -54.2596220312
ORF 20: -42.2610172431
ORF 21: -21.947974199
ORF 22: -24.9402284146
ORF 20: -42.2610172431
ORF 22: -24.9402284146
ORF 26: -17.8025912835
ORF 26: -17.8025912835
ORF 26: -17.8025912835
ORF 27: 3.348222412
ORF 28: -11.237014812
ORF 28: -11.237014812
ORF 28: -11.237014812
ORF 28: -11.237014813
ORF 31: -26.552082386
ORF 32: -22.5209018643
ORF 32: -22.5209018643
ORF 33: -26.5326737312
ORF 33: -26.5326737312
ORF 33: -26.5326737312
ORF 33: -26.552082386
ORF 33: -27.5503943641
ORF 39: -15.5533945771
ORF 39: -15.5533945771
ORF 39: -15.5533945771
ORF 39: 5.64330992198
```

Figure 1. Total of 130 predicted ORFs and their likelihood (cont.)

```
Chouse

Chouse
```



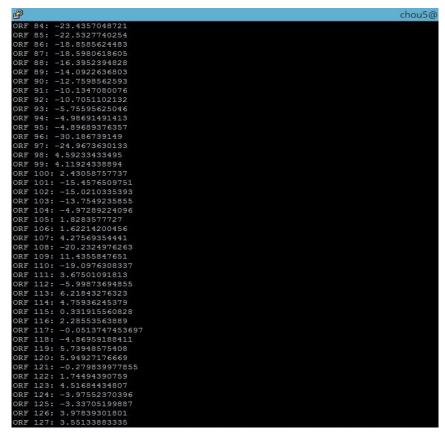


Figure 3. Total of 130 predicted ORFs and their likelihood (cont.)

```
ORF 128: 1.03604196086
ORF 129: 0.285977347663
ORF 130: 0.780194343133
```

Figure 4. Total of 130 predicted ORFs and their likelihood

For our program, there are 130 predicted ORFs (Figure 1, 2, 3, 4) and 46 predicted genes (see Appendix B) which belong to *E.coli*. The two output files contain with the predicted genes and proteins which belongs to *E.coli* (see Appendix B and Appendix C).

```
>gil545778205lgblU00096.3l Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide
          9306
                 15 + 1 12.15
orf00001
                 2799 +1 12.06
orf00002
           343
                 3733 +2 11.17
          2801
orf00003
orf00005
          3734
                  5020 + 2
                            14.13
          5288
                  5530 + 2
                            3.57
orf00006
orf00009
          6459
                  5683 -1 12.06
          7959
                            9.23
orf00011
                  6529 -1
orf00012
          8307
                 9191 +3 14.29
```

Figure 5. The 8 predicted genes from Glimmer

gil545778205 gblU00096.3 Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide 2799 -3294,262806 + 0 gene id 1	GeneMark.hmm	CDS	337
gil545778205 gblU00096.3 Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide 3733 -1244.454019 + 0 gene id 2	GeneMark.hmm	CDS	2801
gil545778205lgblU00096.3l Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide 5020 -1710.714178 + 0 gene id 3	GeneMark.hmm	CDS	3734
gil545778205 gblT00096.3 Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide 5534 -236.846201 + 0 gene id 4	GeneMark.hmm	CDS	5349
gil545778205 gb U00096.3 Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide	GeneMark.hmm	CDS	5683
gil545778205 gb U00096.3 Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide	GeneMark.hmm	CDS	6529
gil545778205 gb U00096.3 Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide	GeneMark.hmm	CDS	8238
9191 -1258.031760 + 0 gene_id 7 gil545778205lgblU00096.3l Escherichia coli str. K-12 substr. MG1655, first 9800 nucleotide 9800 -662.860988 + 0 gene_id 8	GeneMark.hmm	CDS	9306

Figure 6. The 8 predicted genes from GeneMark

Moreover, we use Glimmer and GeneMark to predict our given DNA and they both predict 8 genes which are belong to *E.coli* (Figure 5, 6). By comparing our results with the two commercial available software, we found that our program contains many false positive results. Further optimization of the program to add more threshold to the score and prediction are required.

IV. Appendix

A. Codon-Usage Table (frequency per 1000 bp)

```
'CTT': 10.30967530779820, 'ATG': 27.736331578153628, 'ACA': 6.282709073348543, 'ACG': 14.471408539735407,
'ATC': 25.06987027072123, 'AAC': 21.873967187444848, 'ATA': 3.6643788364474137, 'AGG': 1.026796171333776.
'CCT': 6.741558612413324, 'ACT': 8.4614421993974000, 'AGC': 16.621263023465502, 'AAG': 10.07222869317726,
'AGA': 1.761597181444510, 'CAT': 12.616757955263774, 'AAT': 16.220170769038244, 'ATT': 29.55568604423566,
'CTG': 55.20312916133214, 'CTA': 3.757432239474537, 'CTC': 10.887248154173445, 'CAC': 9.972757814079301,
'AAA': 32.73554543733495, 'CCG': 24.453792567920964, 'AGT': 8.028262564615963, 'CCA': 8.272126655307735,
'CAA': 15.16449595538571, 'CCC': 5.6088740859107520, 'TAT': 16.034063962984000, 'GGT': 24.51796732862932,
'TGT': 4.935039098472962, 'CGA': 3.3627574611181170, 'CAG': 30.139676366681748, 'TCT': 8.025053826580544,
'GAT': 31.61890460100947, 'CGG': 5.5992478718044980, 'TTT': 21.931724472082376, 'TGC': 6.484859569579879,
'GGG': 11.34930643127365, 'TAG': 0.2085679723021733, 'GGA': 7.556578073409509, 'TAA': 1.9380777733925025,
'GGC': 30.61777833395903, "TAC': 12.764359904893006, "TTC': 16.425530003305000, "TCG': 9.055058735949737,
"TTA": 12.908753116486817, "TTG": 13.239253134134877, "TCC": 8.660383957593318, 'ACC': 24.367156640964673,
'TCA': 6.436728499048609, 'GCA': 19.656729204970976, 'GTA': 10.707558824190034, 'GCC': 26.093457704019585,
'GTC': 15.254340620377413, 'GCG': 35.48543393368822, 'GTG': 27.46679758317851, 'GAG': 18.078030091545298,
'GTT': 17.57425821998466, 'GCT': 15.129199836996108, 'TGA': 0.9080728640233083, 'GAC': 19.84925348709606,
'CGT': 21.44078755266342, 'TGG': 15.745277539796373, 'GAA': 39.233239959056505, 'CGC': 22.63122936380351
```

B. Output FASTA file for the predicted genes (Real file: output_gene.fasta)

```
| Note 1: | A recommendation of the control of the
```

```
>ORF 62:
>ORF 63:
ATGCCAATCATCTGGACAATCCCTTGCGCTGCCGGATGCCGGAGGCCCAGGACGCCGCTGCCGCTGCCGCTTTGGCGTCGAACCCATTCCCGCCTCATTGGAAAACATACTGCGCTGA
>ORF 64:
ATGACGTGGG
            AGTTGCCCGATATTCATTACGCAAATTACCAGGCTGGTCAGTACCCAGATTATCGCCATCAACGGGACAAAGCCCTGCATGAGCCGGGCGACGCCATGA
>ORF 66:
ATGAGCAAAAAGACGCGAACAGAACGCCCATCCAGCGCATCCCAGCCCGGCGCGCCCATATACCATGCCGGTCCGCCACGAACTGCCCATTGACGTCACGTTCTTTATAA
ATGCCGATAA
>ORF 68:
ATGAATGC
>ORF 69:
ATGACTITO
>ORF 74:
>ORF 79:
ATGGTATATO
                      GGATGCGCTGGATGGGCGTTCTGTTCGCCGTCTTTTTGCTCATCGCCTATGGCATAATTTTCAGCGGAGTTCAAGCGAACGCCGTTGCCCGCGCCCTGAGTTTTTCTTTTGA
AGE 98:
AIGACTGGCAGCCAGATTICACGCCGGCGAATGCCCGCCAGGCGATTCGGCGTTTAAAGGTGATGTCTACACCGGCTTGCAGGCCGAAAACCTTCAGCGAAGACGATTTCGACACAACAGCATTTGCGAATGCTTTCGGCTTGATGGCGTATCCCCGCTTGAAGGTGATGCCTACACCGGCTTGAAGGTGATTCCGAGCCGAAACCTTCAG
>ORF 99:
ATGTCTACACC
>ORF 105:
ATGAAGATTO
>ORF 106:
ATGGCGAACTGGTGT
>ORF 107:
```

```
>ORF 111:
ATGCCAGCGTTGGCC
        CGTGGAACAATTCCAGACAACCGACATCGCTTTCAACATTGGCGACCGGAGCCGGGAAGGCAAACGCCGCGCGCCACGCGCTCTTCCAGGATTTCCTGTGGGATTTCATCACCAATAA
>ORF 113:
>ORF 114:
ATGCAATCC
>ORF 115:
ATGAATGCT
>ORF 119:
ATGGCGACAATGTTGATATTGGCGCGGGCCAGTGCGGCAAAGAATTTCGCCGAGATCCCACGCAAGGTGCGCATACCATCACCTACCACCAGATAATGGCCAGCCGTTCCGTCACTGCCAGCGGCTCCAGTAAGCCTTCTTTCAGTTCCAGGTAG
ATGTTGATATTGG
        GGCCAGTGCGGCAAAGAATTTCGCCGAGATCCCACGCAAGGTGCGCATACCATCACCTACCACGAGATAATGGCCAGCCGTTCCGTCACTGCCAGCGGCTCCAGTAAGCCTTCTTTCAGTTCCAGGTAG
>ORF 122:
ATGCCGACCAT
SORE 123.
>ORF 126:
>ORF 127:
ATGTACCGCC
>ORF 128:
        ATGGTGGTGGTAA
```

C. Output FASTA file for the predicted proteins (Real file: output_protein.fasta)

```
MQNVFCVLPIFWKAMPGRGRWPPSSLPPPKSPTTWWR_
>Protein 8:
MKKANWWCLDATVPTTLLRCWLPVYAPIVARFGRTLTGSIPATRVRCPMRGC_
>Protein 18:
MVCAPCVGSRRNSLPHWPAPISTLSPLLRDLLNAQSLSW
>Protein 27:
MLPLLWRICHNSTISLPRAWRRPVMKEKFCAMLAILMKMASAA
>Protein 29:
MAKTPWPSIATIISRCRWYCADMVRAMTLQLPVSLLICYVPSHGS
>Protein 39:
MKSCGCMRIRGLKSRRQKPGLFYRRSIAARIALRTGDIWQASFTPAIPVSLSLPRS_
>Protein 43:
MKSHRKSWKSACARRLPSRLRSPMLKAMSVVWNCSTGQRWHLKISAVALWHKC_
MACSSAPRIRRNLKRAWKRFSVKRWICQKSWQNVLIYPCFHIICPPILLRCVN_
>Protein 49:
MEKNDREKGEILNKCGNLEIRIAENNNRRSHRVISGYRPITGNDKRSNL
>Protein 50:
MIKGVICEKDAIYRTRTFPGSGRSHGSTGCGNYVSPVSKITDRRS_
>Protein 55:
MILNLPFFSSRNTGLIISAFNFFGFTDLKYSSDARLITTLSPCAASASFSLLVMISPON
MAVKLVIAMQAIMISASCHISGRLSKVPPMTVAKVQIRKHQIALGLFRRKKKRLFSA
MQAIMISASCHISGRLSKVPPMTVAKVQIRKHQIALGLFRRKKKRLFSA_
MISASCHISGRLSKVPPMTVAKVQIRKHQIALGLFRRKKKRLFSA
```

```
33 >Protein 62:
34 MTRVSIKMPIIWTIPCAAGCGGQDAAAAAAFGVEPIPASLENILR
35 >Protein 63:
36 MPIIWTIPCAAGCGGQDAAAAAAFGVEPIPASLENILR
    >Protein 64:
38 MTWGSCPIFITQITRLVSTQIIAINGTKPCMSRATP_
40 MSKKTANRTPIQRIPSPRAIYHAGPPRNCPLTSRSL
42 MPINAATHIQKTAPGPPAVIASATPARLPLPTRAARLVHND
44 MNAIFKTFAKLADVAKLHKSGAKSEPTTCAEEQVNHYRSPKDAVNEGEKIWHAYPSYCRSR
46 MTFLCKRYQLFKKYIRYEVLDNFCVAYMRFCIKMAGDINAVSEIRNSLAWR
    >Protein 74:
48 MSAILPAQGCLYPDKRRWMRDQRWKLSQYVKYSSRFFRLNSLVKQNDTLRRE_
    >Protein 79:
     MVYGARAGDALDGRSVRRLFAHRLWHNFQRSSSERRCPRPEFFF_
    >Protein 98:
    MTGSQISRRRMPARRFWRLKVMSTPACRPKPSAKTISILPNSICECFPACMAYSARSI_
    >Protein 99:
54 MPARRFWRLKVMSTPACRPKPSAKTISILPNSICECFPACMAYSARSI
   >Protein 100:
MSTPACRPKPSAKTISILPNSICECFPACMAYSARSI_
58 MKIPPAMANWCINATSSVKFAVNAGCGVNALSCLØKHADSIYCRINVGLIRRASVASGAECRIVTRHLSFSGDVCLDHDRRDDHEELSYGGAAVAAVRVGASDCLAIHNVVSTSRGRASDLHPSNSHDYHDRLSVILLTGLT
    MANWCLNATSSVKFAVNAGCGVNALSCLQKHADSIYCRINVGLIRRASVASGAECRIVTRHLSFSGDVCLDHDRRDDHEELSYGGAAVAAVRVGASDCLAIHNVVSTSRGRASDLHPSNSHDYHDRLSVILLIGLT
62 MGATRTRESASTIDCIFFTGYSFYHCPLWVDIRRLLDENGGCYSPQS
64 MVIGLSPAIWVSICAIKRPPKSLNASVGPWNNSRQPTSLSTLATGAGKANAARTRSSRISCGISSPINAERIFALRVIKSSFSISSISVRLNSGRSCGKNSP_
66 MPALARGTIPDNRHRFQHWRPEPGRQTPRARALPGFPVGFHHQ
    >Protein 113:
68 MKPARCRPCAMQSWRRYCAGKIALASAVETLIPGYASTHSHSSNPGTCWLMMSFSSIINCIPPRKHGATLS_
    >Protein 114:
    MOSWRRYCAGKIALASAVETLIPGYASTHSHSSNPGTCWLMMSFSSIINCIPPRKHGATLS
    >Protein 115:
MNAAGDTAFQLAHQSQQTSVIKRLAAVFIHRHQRRDHRTGTGA_
    MLEEIRPSSSPIKASKRVSLSGLPQCSFIAISAATTEQALEPKPEPIGIFFSRVIATGICLPSSWQKRSQH_
>Protein 119:
MAIMLILARASAAKNFAEIPRKVRIPSPITEIMASRSVTASGSSKPSFSSR_
     >Protein 120:
     MLILARASAAKNFAEIPRKVRIPSPTTEIMASRSVTASGSSKPSFSSR
    >Protein 122:
MPTIPFIPGPETLNIAMLFRLEMPLTGNSSSSRLAPMSVPGA
    >Protein 123:
MAIDIFSPRQISAALMLSGHCPNKLMPCRTCFIWANS_
    WYSIIATRUVILAGAERTVATCPCLALLSRISATRKKRSAFATDVPPNFNTRMVVTSLPLVEKKSPHCQVRAFFCVSCTRQPAPLPVVMVMVVVMVVLMRFMDVVYSVIFICLCAMPILVKVFSDLSQ_
>Protein 127:
MYRRISTLAWLLPRYLMSKKKARTVRCGLFSVFPVRVSPHRYLW_
     {\tt MVVTSLPLVEKKSPHCQVRAFFCVSCTRQPAPLPVVMVMVVVMVVLMRFMDVVYSVIFICLCAMPILVKVFSDLSQ\_}
90 MVMVVVMVVLMRFMDVVYSVIFICLCAMPILVKVFSDLSQ
    >Protein 130:
MVVVMVVLMRFMDVVYSVIFICLCAMPILVKVFSDLSQ
```

V. Author Contributions

- Chia-Hsuan Chou: did coding for finding ORF and calculating likelihood and wrote report.
- Venkata Prudhvi Raj Indana: did coding for building codon usage table and debugging.
- Jing Wang: did coding for combining the above two codes, implemented Glimmer and GeneMark and made power point.
- Our group members are contributed evenly to the project.