# Jacek Radajewski
# Student number: 43612772
# SCIE2100 Prac 7

## Question 1

### *Results*

| Genome | Mean | STD DEV | Upper Bound | Lower Bound |
|--------|------|---------|-------------|-------------|
| genome1 | 0.38 | 0.09 | 0.56 | 0.2 |
| genome2 | 0.64 | 0.02 | 0.68 | 0.6 |
| genome3 | 0.39 | 0.04 | 0.47 | 0.31 |

### *Code*

```python
'''
Created on 13/05/2014

@author: jacekrad
'''
from sequence import *

fasta_files = ["Genome1.fasta", "Genome2.fasta", "Genome3.fasta"]

def get_gc_count(sequence):
    return sequence.count('G') + sequence.count('C')

def get_gc_fraction(sequence):
    return round(float(get_gc_count(sequence)) / len(sequence), 2)

def get_mean(values):
    return round(float(sum(values)) / len(values), 2)

def get_standard_deviation(values):
    vals = []
    mean = get_mean(values)
    for i in range(len(values)):
        vals.append((values[i] - mean) ** 2)
    standard_deviation = math.sqrt((1 / float(len(values))) * sum(vals))
    return round(standard_deviation, 2)

#dictionary in which we'll save the contigs
contigs = {}

for fasta_file in fasta_files:
    contig_list = []
    contigs.update({fasta_file:contig_list})
    sequences = readFastaFile(fasta_file, DNA_Alphabet)
    for sequence in sequences:
        contig_list.append(get_gc_fraction(sequence.sequence))
    mean = get_mean(contig_list)
    standard_deviation = get_standard_deviation(contig_list)
    upper_bound = mean + (2 * standard_deviation)
```

```
    lower_bound = mean - (2 * standard_deviation)
    print fasta_file, ": ", mean, standard_deviation, upper_bound,
lower_bound
```

## Unit Test

```python
import unittest
import question1 as q1

class Q1Test(unittest.TestCase):

    data1 = [1, 2, 3, 4, 5]
    data2 = [1, 2, 3, 4, 5, 6]

    def setUp(self):
        pass

    def tearDown(self):
        pass

    def test__mean__1(self):
        self.assertEquals(3, q1.get_mean(self.data1), "")

    def test__mean__2(self):
        self.assertEquals(3.5, q1.get_mean(self.data2), "")

    def test__standard_deviation__1(self):
        self.assertEquals(1.41, q1.get_standard_deviation(self.data1), "")

    def test__standard_deviation__2(self):
        self.assertEquals(1.71, q1.get_standard_deviation(self.data2), "")

if __name__ == "__main__":
    # import sys;sys.argv = ['', 'Test.testName']
    unittest.main()
```

## Output

```
Genome1.fasta :  0.38 0.09 0.56 0.2
Genome2.fasta :  0.64 0.02 0.68 0.6
Genome3.fasta :  0.39 0.04 0.47 0.31
```
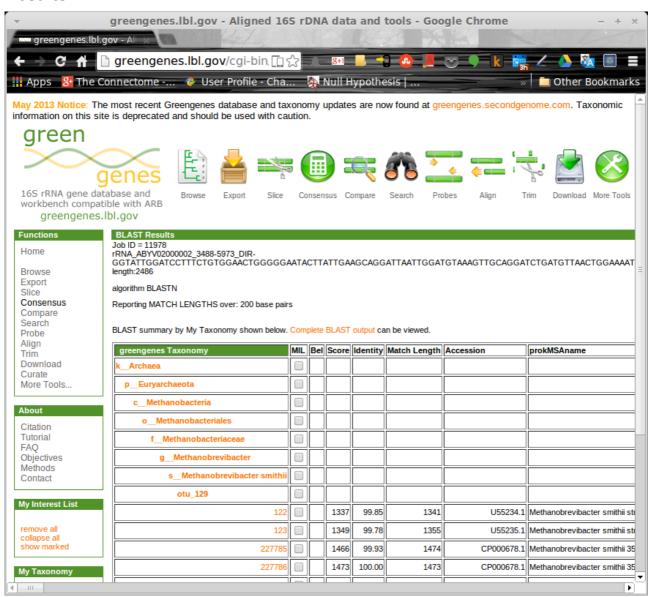
# Question 2

| | Genome1 | Genome2 | Genome3 |
|---|---|---|---|
| Kingdom | Archaea | Bacteria | Bacteria |
| Phylum | Euryarchaeota | Actinobacteria | Firmicutes |
| Class | Methanobacteria | Actinobacteria (class) | Clostridia |
| Order | Methanobacteriales | Actinomycetales | Clostridiales |
| Family | Methanobacteriaceae | Mycobacteriaceae | Veillonellaceae |
| Genus | Methanobrevibacter | Mycobacterium | Veillonella |
| Species | s__Methanobrevibacter smithii | Unclassified | Veillonella dispar |

## *Genome1*

## 16S rRNA

```
>rRNA_ABYV02000002_3488-5973_DIR- /molecule=16s_rRNA /score=612.8
GGTATTGGATCCTTTCTGTGGAACTGGGGGAATACTTATTGAAGCAGGATTAATTGGATG
TAAAGTTGCAGGATCTGATGTTAACTGGAAAATGAAAAATGGAAGTGCAATCAATTTGGA
TTACTGTGGAATAACTGATTATAGAACATTCAATGTTGATGTACGTGAACTTAAAATGTA
TGAAAAGGTAGATAGTGTAGTTACAGACCCTCCTTATGGAATATCTACTTCAACTGGGGA
TATTGAGGGTGATGAGATTTTCAATGAGTTTTTCCATTCAATTTATGATAATATGAAAGA
TGATGCCTACTTGTGTATGGCTAGTCCTCATTATGTTGATTTAAATCCTATGATTAAGGA
AGTTGGATTTGAATTAGTTGAACAATATGGAATCAAAATGCATAGAAGTTTAACAAGGAT
AATTTCAGTTATTCGTAAGAAAAATGTTTAATTTTTTTTATTTATTTATAAATAAATAGG
TAAGTTTATATATTACTTAGTAATCGATTTTACTTATTTATATATTATATTTGCTTTTTT
TCCAATTTATTTATTTTAAAAGAGATTTTTTACTTTTGCTTTTCATTTTGACTTTTAAAA
TTTTCATTAGGTCTATTTGATTAAATTTTTCATTTTTATGAATCATTAGTTTAACTATTG
TTTTTTTATAAATTAAAGATGTTTCATTAATATTTTTTTATTTGATTTATTTTTGCAATA
ATTTTTAATTGAATTTGTTGATATAATTTCCTTGTCTATTAGATGTGCTATGTGTTATAT
GGCGTTGGTCTAAGTTACATTGTATTGACAATTATAACTATGATGACTGCTTTTCCACAT
GTAGCTAAACATATTTAATGAGCACTTGAGTTTTTTGAGTGTGATGTTGGTTTTGTAGAT
GTGGTGAATTTGATTACATTATTGTTATCAAATCAGTGTATTTACTCGTCTATATTTTTT
TAGCGTACTTCATAAATTTTAGCTTTTTTTATGATTATCTGCTTTTTTCATTCAATTCTGT
TTGATCCTGGCAGATGCTACTGCTATTGGGATTCGATTAAGCCATGCAAGTCGAACGAGT
TTAGGCTCGTGGCGTACGGCTCAGTAACACGTGGATAACCTACCCTTAGGACTGGGATAA
CCCTGGGAAACTGGGGATAATACTGGATAGGCAATTATTCCTGTAATGGTTTTTTGTTTA
AATGTTTTTTCGCCTAAGGATGGGTCTGCGGCCGATTAGGTAGTTGGTTAGGTAATGGCT
TACCAAGCCTTTGATCGGTACGGGTTGTGAGAGCAAGAGCCCGGAGATGGAACCTGAGAC
AAGGTTCCAGGCCCTACGGGGTGCAGCAGGCGCGAAACCTCCGCAATGTGAGAAATCGCG
ACGGGGGGATCCCAAGTGCCATTCTTAACGGGATGGCTTTTCATTAGTGTAAAGAGCTTT
TGGAATAAGAGCTGGGCAAGACCGGTGCCAGCCGCCGCGGTAACACCGGCAGCTCTAGTG
GTAGCAGTTTTTATTGGGCCTAAAGCGTCCGTAGCCGGTTTAATAAGTCTCTGGTGAAAT
CCTGCAGCTTAACTGTGGGAATTGCTGGAGATACTATTAGACTTGAGATCGGGAGAGGTT
AGAGGTACTCCCAGGGTAGAGGTGAAATTCTGTAATCCTGGGAGGACCGCCTGTTGCGAA
GGCGTCTGACTGGAACGATTCTGACGGTGAGGGACGAAAGCTAGGGGCGCGAACCGGATT
AGATACCCGGGTAGTCCTAGCTGTAAACGATGCGGACTTGGTGTTGGGGTGGCTTTGAGC
TGTCCCAGTGCCGAAGGGAAGCTGTTAAGTCCGCCGCCTGGGAAGTACGGTCGCAAGACT
GAAACTTAAAGGAATTGGCGGGGGAGCACCACAACGCGTGGAGCCTGCGGTTTAATTGGA
TTCAACGCCGGACATCTCACCAGAGGCGACAGCTGTATGATAGCCAGGTTGATGACTTTG
CTTGACTAGCTGAGAGGAGGTGCATGGCCGCCGTCAGCTCGTACCGTGAGGCGTCCTGTT
AAGTCAGGCAACGAGCGAGACCCACGCTCTTAGTTACCAGCGGATCCTTTTTTGGATGCC
GGGCACACTAAGGGGACCGCCTATGATAAATAGGAGGAAGGAGTGGACGACGGTAGGTCC
GTATGCCCCGAATCCTCTGGGCAACACGCGGGCTACAATGGCTGAGACAATGGGTTCCGA
CGCCGAAAGGCGGAGGTAATCCTCTAAACTTAGTCGTAGTTCGGATTGAGGACTGTAACT
CGTTCTCATGAAGCTGGAATGCGTAGTAATCGCGTGTCACAATCGCGCGGTGAATACGTC
CCTGCTCCTTGCACACACCGCCCGTCACGCCACCCAAAAAGGGATTGGATGAGGATGTAA
TGTTTTGTTATATTCGAATCTAGTTTTTTTAAGGAGGGCGAAGTCGTAACAAGGTAGCCG
TAGGGGAACCTGCGGCTGGATCACCT
```

## Results



## Genome2

## 16S rRNA

```
>rRNA_Genome2_Contig1_1475623-1477142_DIR+ /molecule=16s_rRNA /score=1904.5
AGAGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCTTAACACATGCAAGTCGAAC
GGAAAGGTCTCTTCGGAGATACTCGAGTGGCGAACGGGTGAGTAACACGTGGGTGATCTG
CCCTGCACTTCGGGATAAGCCTGGGAAACTGGGTCTAATACCGGATAGGACCACGGGATG
CATGTCTTGTGGTGGAAAGCGCTTTAGCGGTGTGGGATGAGCCCGCGGCCTATCAGCTTG
TTGGTGGGGTGACGGCCTACCAAGGCGACGACGGGTAGCCGGCCTGAGAGGGTGTCCGGC
CACACTGGGACTGAGATACGGCCCAGACTCCTACGGGAGGCAGCAGTGGGGAATATTGCA
CAATGGGCGCAAGCCTGATGCAGCGACGCCGCGTGGGGGATGACGGCCTTCGGGTTGTAA
ACCTCTTTCACCATCGACGAAGGTCCGGGTTCTCTCGGATTGACGGTAGGTGGAGAAGAA
GTACCGGCCAACTACGTGCCAGCAGCCGCGGTAATACGTAGGGTGCGAGCGTTGTCCGGA
ATTACTGGGCGTAAAGAGCTCGTAGGTGGTTTGTCGCGTTGTTCGTGAAATCTCACGGCT
TAACTGTGAGCGTGCGGGCGATACGGGCAGACTAGAGTACTGCAGGGGAGACTGGAATTC
CTGGTGTAGCGGTGGAATGCGCAGATATCAGGAGGAACACCGGTGGCGAAGGCGGGTCTC
TGGGCAGTAACTGACGCTGAGGAGCGAAAGCGTGGGGAGCGAACAGGATTAGATACCCTG
GTAGTCCACGCCGTAAACGGTGGGTACTAGGTGTGGGTTTCCTTCCTTGGGATCCGTGCC
GTAGCTAACGCATTAAGTACCCCGCCTGGGGAGTACGGCCGCAAGGCTAAAACTCAAAGG
AATTGACGGGGGCCCGCACAAGCGGCGGAGCATGTGGATTAATTCGATGCAACGCGAAGA
ACCTTACCTGGGTTTGACATGCACAGGACGCGTCTAGAGATAGGCGTTCCCTTGTGGCCT
GTGTGCAGGTGGTGCATGGCTGTCGTCAGCTCGTGTCGTGAGATGTTGGGTTAAGTCCCG
CAACGAGCGCAACCCTTGTCTCATGTTGCCAGCACGTAATGGTGGGGACTCGTGAGAGAC
TGCCGGGGTCAACTCGGAGGAAGGTGGGGATGACGTCAAGTCATCATGCCCCTTATGTCC
AGGGCTTCACACATGCTACAATGGCCGGTACAAAGGGCTGCGATGCCGCGAGGTTAAGCG
AATCCTTAAAAGCCGGTCTCAGTTCGGATCGGGGTCTGCAACTCGACCCCGTGAAGTCGG
AGTCGCTAGTAATCGCAGATCAGCAACGCTGCGGTGAATACGTTCCCGGGCCTTGTACAC
```

```
ACCGCCCGTCACGTCATGAAAGTCGGTAACACCCGAAGCCAGTGGCCTAACCCTCGGGAG
GGAGCTGTCGAAGGTGGGATCGGCGATTGGGACGAAGTCGTAACAAGGTAGCCGTACCGG
AAGGTGCGGCTGGATCACCT
```

# Results



# *Genome3*

## 16S rRNA

```
>rRNA_AEDS01000059_36-1586_DIR+ /molecule=16s_rRNA /score=1871.2
AGAGTTTGATCCTGGCTCAGGACGAACGCTGGCGGCGTGCTTAACACATGCAAGTCGAAC
GAAGAGCGATGGAAGCTTGCTTCTATCAATCTTAGTGGCGAACGGGTGAGTAACGCGTAA
TCAACCTGCCCTTCAGAGGGGGACAACAGTTGGAAACGACTGCTAATACCGCATACGATC
CAATCTCGGCATCGAGACTGGATGAAAGGTGGCCTCTATTTATAAGCTATCACTGAAGGA
GGGGATTGCGTCTGATTAGCTAGTTGGAGGGGTAACGGCCCACCAAGGCGATGATCAGTA
GCCGGTCTGAGAGGATGAACGGCCACATTGGGACTGAGACACGGCCCAGACTCCTACGGG
AGGCAGCAGTGGGGAATCTTCCGCAATGGACGAAAGTCTGACGGAGCAACGCCGCGTGAG
TGATGACGGCCTTCGGGTTGTAAAGCTCTGTTAATCGGGACGAATGGTTCTTGTGCGAAT
AGTGCGAGGATTTGACGGTACCGGAATAGAAAGCCACGGCTAACTACGTGCCAGCAGCCG
CGGTAATACGTAGGTGGCAAGCGTTGTCCGGAATTATTGGGCGTAAAGCGCGCGCAGGCG
GATTAGTTAGTCTGTCTTAAAAGTTCGGGGCTTAACCCCGTGATGGGATGGAAACTGCTG
ATCTAGAGTATCGGAGAGGAAAGTGGAATTCCTAGTGTAGCGGTGAAATGCGTAGATATT
AGGAAGAACACCAGTGGCGAAGGCGACTTTCTGGACGAAAACTGACGCTGAGGCGCGAAA
```

```
GCCAGGGGAGCGAACGGGATTAGATACCCCGGTAGTCCTGGCCGTAAACGATGGGTACTA
GGTGTAGGAGGTATCGACCCCTTCTGTGCCGGAGTTAACGCAATAAGTACCCCGCCTGGG
GAGTACGACCGCAAGGTTGAAACTCAAAGGAATTGACGGGGGCCCGCACAAGCGGTGGAG
TATGTGGTTTAATTCGACGCAACGCGAAGAACCTTACCAGGTCTTGACATTGATGGACAG
AACTAGAGATAGTTCCTCTTCTTCGGAAGCCAGAAAACAGGTGGTGCACGGTTGTCGTCA
GCTCGTGTCGTGAGATGTTGGGTTAAGTCCCGCAACGAGCGCAACCCCTATCTTATGTTG
CCAGCACTTCGGGTGGGAACTCATGAGAGACTGCCGCAGACAATGCGGAGGAAGGCGGGG
ATGACGTCAAATCATCATGCCCCTTATGACCTGGGCTACACACGTACTACAATGGGAGTT
AATAGACGGAAGCGAAACCGCGAGGTGGAGCAAACCCGAGAAACACTCTCTCAGTTCGGA
TCGTAGGCTGCAACTCGCCTACGTGAAGTCGGAATCGCTAGTAATCGCAGGTCAGCATAC
TGCGGTGAATACGTTCCCGGGCCTTGTACACACCGCCCGTCACACCACGAAAGTCGGAAG
TGCCCAAAGCCGGTGGGGTAACCTTCGGGAGCCAGCCGTCTAAGGTAAAGTCGATGATTG
GGGTGAAGTCGTAACAAGGTAGCCGTATCGGAAGGTGCGGCTGGATCACCT
```
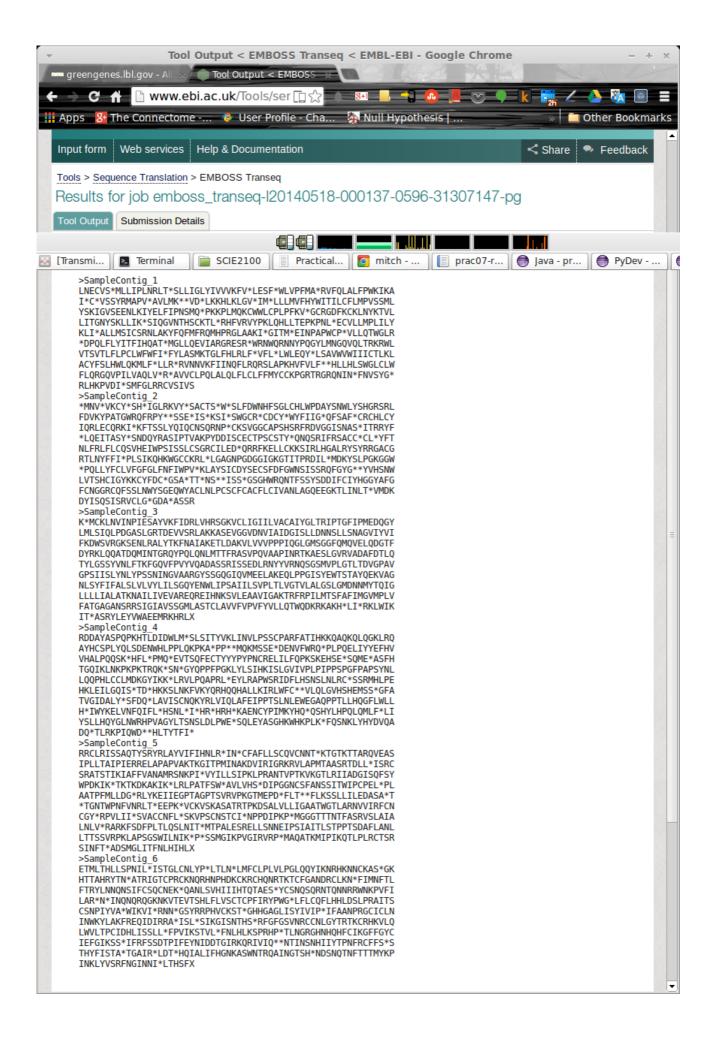
## Results

# Question 3

Out of the six reading frames (shown in screen grab) the following one is the correct translation because it starts with the correct START codon (M), ends with a STOP codon (*), and has a sufficient length to be a protein (over a hundred of residues).

>SampleContig_3
K***MCKLNVINPIESAYVKFIDRLVHRSGKVCLIGIILVACAIYGLTRIPTGFIPMEDQGY LMLSIQLPDGASLGRTDEVVSRLAKKASEVGGVDNVIAIDGISLLDNNSLLSNAGVIYVI FKDWSVRGKSENLRALYTKFNAIAKETLDAKVLVVVPPPIQGLGMSGGFQMQVELQDGTF DYRKLQQATDQMINTGRQYPQLQNLMTTFRASVPQVAAPINRTKAESLGVRVADAFDTLQ TYLGSSYVNLFTKFGQVFPVYVQADASSRISSEDLRNYYVRNQSGSMVPLGTLTDVGPAV GPSIISLYNLYPSSNINGVAARGYSSGQGIQVMEELAKEQLPPGISYEWTSTAYQEKVAG NLSYFIFALSLVLVYLILSGQYENWLIPSAIILSVPLTLVGTVLALGSLGMDNNMYTQIG LLLLIALATKNAILIVEVAREQREIHNKSVLEAAVIGAKTRFRPILMTSFAFIMGVMPLV FATGAGANSRRSIGIAVSSGMLASTCLAVVFVPVFYVLLQTWQDKRKAKH**\*LI*RKLWIK IT*ASRYLEYVWAEEMRKHRLX

greengenes.lbl.gov - Al    Tool Output < EMBOSS
www.ebi.ac.uk/Tools/ser
Apps    The Connectome -...    User Profile - Cha...    Null Hypothesis | ...    Other Bookmarks

Input form   Web services   Help & Documentation                    Share    Feedback

Tools > Sequence Translation > EMBOSS Transeq

## Results for job emboss_transeq-I20140518-000137-0596-31307147-pg

Tool Output   Submission Details

[Transmi...]   Terminal   SCIE2100   Practical...   mitch - ...   prac07-r...   Java - pr...   PyDev - ...

```
>SampleContig_1
LNECVS*MLLIPLNRLT*SLLIGLYIVVVKFV*LESF*WLVPFMA*RVFQLALFPWKIKA
I*C*VSSYRMAPV*AVLMK**VD*LKKHLKLGV*IM*LLLMVFHYWITILCFLMPVSSML
YSKIGVSEENLKIYELFIPNSMQ*PKKPLMQKCWWLCPLPFKV*GCRGDFKCKLNYKTVL
LITGNYSKLLIK*SIQGVNTHSCKTL*RHFVRVYPKLQHLLTEPKPNL*ECVLLMPLILY
KLI*ALLMSICSRNLAKYFQFMFRQMHPRGLAAKI*GITM*EINPAPWCP*VLLQTWGLR
*DPQLFLYITFIHQAT*MGLLQEVIARGRESR*WRNWQRNNYPQGYLMNGQVQLTRKRWL
VTSVTLFLPCLWFWFI*FYLASMKTGLFHLRLF*VFL*LWLEQY*LSAVWVWIIICTLKL
ACYFSLHWLQKMLF*LLR*RVNNVKFIINQFLRQRSLAPKHVFVLF**HLLHLSWGLCLW
FLQRGQVPILVAQLV*R*AVVCLPQLALQLFLCLFFMYCCKPGRTRGRQNIN*FNVSYG*
RLHKPVDI*SMFGLRRCVSIVS
>SampleContig_2
*MNV*VKCY*SH*IGLRKVY*SACTS*W*SLFDWNHFSGLCHLWPDAYSNWLYSHGRSRL
FDVKYPATGWRQFRPY**SSE*IS*KSI*SWGCR*CDCY*WYFIIG*QFSAF*CRCHLCY
IQRLECQRKI*KFTSSLYQIQCNSQRNP*CKSVGGCAPSHSRFRDVGGISNAS*ITRRYF
*LQEITASY*SNDQYRASIPTVAKPYDDISCECTPSCSTY*QNQSRIFRSACC*CL*YFT
NLFRLFLCQSVHEIWPSISSLCSGRCILED*QRRFKELLCKKSIRLHGALRYSYRRGACG
RTLNYFFI*PLSIKQHKWGCCKRL*LGAGNPGDGGIGKGTITPRDIL*MDKYSLPGKGGW
*PQLLYFCLVFGFGLFNFIWPV*KLAYSICDYSECSFDFGWNSISSRQFGYG**YVHSNW
LVTSHCIGYKKCYFDC*GSA*TT*NS**ISS*GSGHWRQNTFSSYSDDIFCIYHGGYAFG
FCNGGRCQFSSLNWYSGEQWYACLNLPCSCFCACFLCIVANLAGQEEGKTLINLT*VMDK
DYISQSISRVCLG*GDA*ASSR
>SampleContig_3
K*MCKLNVINPIESAYVKFIDRLVHRSGKVCLIGIILVACAIYGLTRIPTGFIPMEDQGY
LMLSIQLPDGASLGRTDEVVSRLAKKASEVGGVDNVIAIDGISLLDNNSLLSNAGVIYVI
FKDWSVRGKSENLRALYTKFNAIAKETLDAKVLVVVPPPIQGLGMSGGFQMQVELQDGTF
DYRKLQQATDQMINTGRQYPQLQNLMTTFRASVPQVAAPINRTKAESLGVRVADAFDTLQ
TYLGSSYVNLFTKFGQVFPVYVQADASSRISSEDLRNYYVRNQSGSMVPLGTLTDVGPAV
GPSIISLYNLYPSSNINGVAARGYSSGQGIQVMEELAKEQLPPGISYEWTSTAYQEKVAG
NLSYFIFALSLVLVYLILSGQYENWLIPSAIILSVPLTLVGTVLALGSLGMDNNMYTQIG
LLLLIALATKNAILIVEVAREQREIHNKSVLEAAVIGAKTRFRPILMTSFAFIMGVMPLV
FATGAGANSRRSIGIAVSSGMLASTCLAVVFVPVFYVLLQTWQDKRKAKH*LI*RKLWIK
IT*ASRYLEYVWAEEMRKHRLX
>SampleContig_4
RDDAYASPQPKHTLDIDWLM*SLSITYVKLINVLPSSCPARFATIHKKQAQKQLQGKLRQ
AYHCSPLYQLSDENWHLPPLQKPKA*PP**MQKMSSE*DENVFWRQ*PLPQELIYYEFHV
VHALPQQSK*HFL*PMQ*EVTSQFECTYYYPYPNCRELILFQPKSKEHSE*SQME*ASFH
TGQIKLNKPKPKTRQK*SN*GYQPPFPGKLYLSIHKISLGVIVPLPIPPSPGFPAPSYNL
LQQPHLCCLMDKGYIKK*LRVLPQAPRL*EYLRAPWSRIDFLHSNSLNLRC*SSRMHLPE
HKLEILGQIS*TD*HKKSLNKFVKYQRHQQHALLKIRLWFC**VLQLGVHSHEMSS*GFA
TVGIDALY*SFDQ*LAVISCNQKYRLVIQLAFEIPPTSLNLEWEGAQPPTLLHQGFLWLL
H*IWYKELVNFQIFL*HSNL*I*HR*HRH*KAENCYPIMKYHQ*QSHYLHPQLQMLF*LI
YSLLHQYGLNWRHPVAGYLTSNSLDLPWE*SQLEYASGHKWHKPLK*FQSNKLYHYDVQA
DQ*TLRKPIQWD**HLTYTFI*
>SampleContig_5
RRCLRISSAQTYSRYRLAYVIFIHNLR*IN*CFAFLLSCQVCNNT*KTGTKTTARQVEAS
IPLLTAIPIERRELAPAPVAKTKGITPMINAKDVIRIGRKRVLAPMTAASRTDLL*ISRC
SRATSTIKIAFFVANAMRSNKPI*VYILLSIPKLPRANTVPTKVKGTLRIIADGISQFSY
WPDKIK*TKTKDKAKIK*LRLPATFSW*AVLVHS*DIPGGNCSFANSSITWIPCPEL*PL
AATPFMLLDG*RLYKEIIEGPTAGPTSVRVPKGTMEPD*FLT**FLKSSLLILEDASA*T
*TGNTWPNFVNRLT*EEPK*VCKVSKASATRTPKDSALVLLIGAATWGTLARNVVIRFCN
CGY*RPVLII*SVACCNFL*SKVPSCNSTCI*NPPDIPKP*MGGGTTTNTFASRVSLAIA
LNLV*RARKFSDFPLTLQSLNIT*MTPALESRELLSNNEIPSIAITLSTPPTSDAFLANL
LTTSSVRPKLAPSGSWILNIK*P*SSMGIKPVGIRVRP*MAQATKMIPIKQTLPLRCTSR
SINFT*ADSMGLITFNLHIHLX
>SampleContig_6
ETMLTHLLSPNIL*ISTGLCNLYP*LTLN*LMFCLPLVLPGLQQYIKNRHKNNCKAS*GK
HTTAHRYTN*ATRIGTCPRCKNQRHNPHDKCKRCHQNRTKTCFGANDRCLKN*FIMNFTL
FTRYLNNQNSIFCSQCNEK*QANLSVHIIIHTQTAES*YCSNQSQRNTQNNRRWNKPVFI
LAR*N*INQNQRQGKNKVTEVTSHLFLVSCTCPFIRYPWG*LFLCQFLHHLDSLPRAITS
CSNPIYVA*WIKVI*RNN*GSYRRPHVCKST*GHHGAGLISYIVIP*IFAANPRGCICLN
INWKYLAKFREQIDIRRA*ISL*SIKGISNTHS*RFGFGSVNRCCNLGYTRTKCRHKVLQ
LWVLTPCIDHLISSLL*FPVIKSTVL*FNLHLKSPRHP*TLNGRGHNHQHFCIKGFFGYC
IEFGIKSS*IFRFSSDTPIFEYNIDDTGIRKQRIVIQ**NTINSNHIIYTPNFRCFFS*S
THYFISTA*TGAIR*LDT*HQIALIFHGNKASWNTRQAINGTSH*NDSNQTNFTTTMYKP
INKLYVSRFNGINNI*LTHSFX
```

# Question 4

The first result in the list was a hypothetical protein lpp2580 (shown in screen grab below). Uniprot search for this protein tells us that it is from **Legionella pneumophila (strain Paris)** species and it's taxonomic lineage is: Bacteria > Proteobacteria > Gammaproteobacteria >

Legionellales › Legionellaceae › Legionella

# Question 5

Top 10 Pathways were calculated based on the counts (scores) provided. For each pathway the scores from each genome were added to give the final score (see code and screenshot below). Pathways with top 10 scores/counts, from highest to lowest were:

| Code | Score | Name |
|---|---|---|
| 03010 | 161 | Ribosome |
| 00230 | 134 | Purine metabolism |
| 00240 | 110 | Pyrimidine metabolism |
| 02010 | 91 | ABC Transporters |
| 00860 | 86 | Porphyrin and chlorophyll metabolism |
| 00680 | 86 | Methane metabolism |
| 00720 | 67 | Carbon fixation pathways in prokaryotes |
| 00910 | 55 | Nitrogen metabolism |
| 00970 | 49 | Aminoacyl-tRNA biosynthesis |
| 00190 | 43 | Oxidative phosphorylation |

Note that the above list includes pathways with highest counts and some of these do not appear in all 3 genomes. Those pathways have been highlighted in cyan.

Above (non-highlighted) pathways are conserved because they are essential to basic function of the cell. For example:

- Ribosome: mRNA translation without which proteins could not be produced

- Purine and Pyrimidine metabolism is what creates nucleotides without which the cell would not be able to function

```
jacekrad@z400 ~/var/github/prac07/prac_7 $ python question5.py | sort -nr | head -10
161 03010
134 00230
110 00240
91 02010
86 00860
86 00680
67 00720
55 00910
49 00970
43 00190
jacekrad@z400 ~/var/github/prac07/prac_7 $
```

```python
'''
Created on 19/05/2014

@author: jacekrad
'''


genome1 = {"00680":86, "03010":59, "00230":36, "00240":35, "00860":24, \
           "00970":24, "00720":19, "00400":18, "02010":18, "00250":16}
genome2 = {"03010":53, "00230":51, "00190":43, "02020":40, "00240":37, \
           "00910":34, "02010":33, "00860":30, "00720":27, "00330":27}
genome3 = {"03010":49, "00230":47, "02010":40, "00240":38, "00860":32, \
           "00970":25, "00270":24, "00720":21, "00910":21, "00400":19}

genomes = [genome1, genome2, genome3]

scores = {}

for genome in genomes:
    for pathway in genome:
        score = genome.get(pathway)
        existing = scores.get(pathway)
        if existing == None:
            scores.update({pathway:score})
        else:
            scores.update({pathway:(existing + score)})

for key in scores:
    print scores.get(key), key
```

# Question 6

## *Methanotrophic archaeon*

Genome 1 is the most likely candidate for this organism as it is the genome that shows methane metabolism.

Feng-Ping Wang, Yu Zhang, Ying Chen, Ying He, Ji Qi, Kai-Uwe Hinrichs, Xin-Xu Zhang, Xiang Xiao and Nico Boon, *Methanotrophic archaea possessing diverging methane-oxidizing and electron-transporting pathways*, The ISME Journal (2014) 8, 1069-1078

## *Veilonella*

Genome3 is the most likeley match for this organism based  on Phenylalanine, tyrosine and tryptophan biosynthesis.

http://patricbrc.org/portal/portal/patric/CompPathwayMap?
cType=genome&cId=168093&dm=feature&feature_info_id=41230233&map=00970&algorithm=PATRIC&ec_number=

## *Mycobacterium*

Genome2 is the most likely candidate for this organism based on the Arginine and proline metabolism.

Anjali Seth and Nancy D. Connell, *Amino Acid Transport and Metabolism in Mycobacteria: Cloning, Interruption, and Characterization of anL-Arginine/γ-Aminobutyric Acid Permease inMycobacterium bovis BCG*, February 2000, Journal of Bacteriology