

'Central Dogma' is the process by which the instructions in DNA are converted into a functional product.

```
raw_data
```

[illegible]

Preprocessing steps:

1. Remove the extra header
2. Remove newlines and tabs("\n\t")

Logic : DNA always starts with the ATG codon. Keeping this in mind let's extract the DNA from the given data

```
#Required data for the central dogma process
#to convert DNA template strand to RNA
mappings = {'A':'T','T':'A','C':'G','G':'C'}
#to convert RNA to protein
gencode = {
    'ATA':'I', 'ATC':'I', 'ATT':'I', 'ATG':'M',
    'ACA':'T', 'ACC':'T', 'ACG':'T', 'ACT':'T',
    'AAC':'N', 'AAT':'N', 'AAA':'K', 'AAG':'K',
    'AGC':'S', 'AGT':'S', 'AGA':'R', 'AGG':'R',
    'CTA':'L', 'CTC':'L', 'CTG':'L', 'CTT':'L',
    'CCA':'P', 'CCC':'P', 'CCG':'P', 'CCT':'P',
    'CAC':'H', 'CAT':'H', 'CAA':'Q', 'CAG':'Q',
    'CGA':'R', 'CGC':'R', 'CGG':'R', 'CGT':'R',
    'GTA':'V', 'GTC':'V', 'GTG':'V', 'GTT':'V',
    'GCA':'A', 'GCC':'A', 'GCG':'A', 'GCT':'A',
    'GAC':'D', 'GAT':'D', 'GAA':'E', 'GAG':'E',
    'GGA':'G', 'GGC':'G', 'GGG':'G', 'GGT':'G',
    'TCA':'S', 'TCC':'S', 'TCG':'S', 'TCT':'S',
    'TTC':'F', 'TTT':'F', 'TTA':'L', 'TTG':'L',
    'TAC':'Y', 'TAT':'Y', 'TAA':'_', 'TAG':'_',
    'TGC':'C', 'TGT':'C', 'TGA':'_', 'TGG':'W'}
```

```
class centralDogma:
    def __init__(self,mappings,gencode):
        self.mappings = mappings
        self.gencode = gencode
    def preprocess(self,raw_data,start_point):
        data = ''.join(raw_data[start_point:].split())
        return data
    '''def transcription(self,data):
        rna = data.translate(str.maketrans(mappings))
        return rna'''
    def translation(self,rna):
        protein = [gencode[rna[i:i+3]] for i in range(0,len(rna),3)]
        return ''.join(protein)
```

```
#Search for the position of 'ATG'
start_point = re.search('ATG',raw_data).start()
cd = centralDogma(mappings,gencode)
data = cd.preprocess(raw_data,start_point)
print('\nPreprocessed data :\n',data)
protein = cd.translation(data)
print('\nProtein :\n',protein)
```

Preprocessed data :

```
ATGTTTGTTTTCTTGTATTATGCCACTAGTCTCTAGTCAGTGTGTTAATCTTACAACCAAGAACTCAATTACCCCTGCATACCTAATCTTTCACACGTGGTGTATTACCCTGACAAAGTT
TTCAGATCCTCAGTTTTACATTTCAACTCAGGACTTGTCTTACCTTTCTTCCAAATGTTACTTGGTTCATGCTATACATGCTCTGGGACCAATGGGTACTAAGAGGTTTGATAACCTGTGCTTAC
CATTTAATGATGGTGTATTATTTGCTTCCACTGAGAAAGTCTAACATAATAAGAGGCTGGATTTTGGTACTACTTTAGATTTCGAAGACCAGTCCCTACTTATTGTTAATAACGCTACTAATGTTGT
TATTAAGTCTGTGAATTTCAATTTTGTAAATGATCCATTTTGGGTGTTTATTACCACAAAACAACAAAGTTGGATGGAAAGTGAGTTCAGAGTTTATCTAGTGCGAATAATGCACTTTTGAA
TATGCTCTCAGCCTTTTCTTATGGACCTTGAAGGAAAAACAGGTAATTTCAAAAACTTAGGGAATTTGTTGTTAAGAAATATTGATGGTTATTTAAAAATATATCTTAAGCACACGCCCTATTAAAT
TAGTGCGTGATCTCCCTCAGGGTTTTTCGGCTTTAGAACCATTGGTAGATTTGCCAATAGGTATTAACATCACTAGGTTTCAAACCTTACTTGCCTTACATAGAAAGTTATTGACTCCTGGTGATTCT
TTCTTCAGGTTGGACAGCTGGTGTGTCAGCTTATTATGTGGGTATCTTCAACCTAGGACTTTTCTATTAAAAATAATGAAAAATGGAACCATACAGATGCTGTAGACTGTGCACCTTGACCCCTCTC
TCAGAAAAAAGTGTACGTTGAAATCTTCACTGTGAAAAAAGGAATCTATCAAACTTCTAACTTTAGAGTCCCAACCAACAGAATCTATTGTTAGATTTCCATAATTACAAACTTTGGCCCTTTTG
GTGAAGTTTTTAACGCCACCAGATTGTCATCTGTTTATGCTTGGAACAGGAAGAGAATCAGCAACTGTGTTGCTGATTATCTGCTCTATATAATCCGCATCATTTCCACTTTTAAAGTGTATTGG
AGTGCTCTCTACTAAATTAATGATCTCTGCTTTACTAATGTCTATGCAGATTCAATTTGTAATTAGAGGTGATGAAGTCAGACAAAATCGCTCCAGGGCAAACTGGAAGATTGCTGATTATAAATAT
AAATTACCAGATGATTTTACAGGCTGCGTTATAGCTTGGAAATCTAACAACTCTGATTCTAAGGTGGTGGTAATTATAATTACCTGTATAGATTGTTTAGGAAGTCTAATCTCAAACTTTTGAGA
GAGATATTTCAACTGAAATCTATCAGGCCGTAGCACACCTTGTAAATGGTGTGGAAGTTTTAATTTGTTACTTTCCCTTTACAATCATATGGTTTCCAAACCCACTAATGGTGTGGTTACCAACCATATA
CAGAGTAGTAGTACTTTCTTTTGAACCTTACATGCACACAGCAACTGTTTGTGGACCTAAAAAGTCTACTAATTTGGTTAAAAACAAATGTGCAATTTCAACTTCAATGGTTTAAACGGCACAGGT
GTTCTTACTGAGTCTAACAAAAAGTTTCTGCCCTTTTCAACAAATTTGGCAGAGACATTGTCTGACACTACTGATGCTGTCCGTGATCCACAGACACTTGAGATTCTTGACATTACACCATGTTCTTTTG
GTGGTGTGAGTGTATAACACAGGAACAAATACTTCTAACAGGTTGTGTTCTTTATCAGGATGTTAAGTGCACAGAAGTCCCTGTTGCTATTCTATGCAGATCAACTTACTCTACTTGGCGTGT
TTATCTACAGGTTCTAATGTTTTTCAAAACAGTGCAGGCTGTTTAAATAGGGGCTGAACATGTCAACAACTCATATGAGTGTGACATACCCATTGGTGCAGGTATATGCGCTAGTTATCAGACTCAG
ACTAATCTCCTCGCGGGGCAGTAGTGTAGCTAGTCAATCCATCATTTGCCCTACACTATGTGCACCTGGTGCAGAAAAATCAGTTGCTTACTCTAATAACTCTATTGCCATACCCACAAATTTTACTA
TTAGTGTACACAGAAATCTACAGTGTCTATGACCAAGACATCAGTAGATTGTACAATGTACATTTTGGTGTATTCACTGAATGCAGCAATCTTTGTTGCAATATGGCAGTTTGTGTACACA
ATTAACACCGTGTCTTAACTGGAAATAGCTGTTGAACAAGACAAAAACACCAAGAAAGTTTTGACACAAGTCAACAAATTTACAAAAACACCAATTAAGAGTTTTGGTGGTTTTAAATTTTTCACAA
ATATTACCAGATCCATCAAAACCAAGCAAGAGGTGATTTATTGAAGATCTACTTTTCAACAAAGTGACACTTGCAGATGCTGGCTTCATCAAAACAAATATGGTGTATGCTTGGTGATATTGCTGCTGA
GAGACCTCAATTTGTGCACAAAAAGTTTAAAGGCTTACTGTTTTGCCACCTTTGCTCACAGATGAAATGATTGCTCAATACACTTCTGCACCTGTTAGCGGGTACAATCACTTCTGGTTGGACCTTTGG
TGCAGGTGCTGCATTACAAATACCATTTGCTATGCAAAATGGCTTATAGGTTTAAATGGTATTGGAGTTACACAGAATGTTCTCTATGAGAACCAAAAAATGATTGGCAACCAATTTAATAGTGTCTATT
GGCAAAATTCAGACTCACTTTCTTCCACAGCAAGTGCACCTTGGAAAACTTCAAGATGTGGTCAACCAAAATGCACAAGCTTTAAACACGCTTGTGTTAAACAACTTAGCTCCAATTTTGGTGCATTTT
CAAGTGTTTTAAATGATATCCTTTTCAAGTCTTGCACAAAGTTGAGGCTGAAGTGCAAATTTGATAGGTTGATCACAGGCAGACTTCAAAAGTTTGCAGACATATGTGACTCAACAATTAATTAGAGCTGC
AGAAATCAGAGCTTCTGCTAATCTTGTCTGCTACTAAAATGTGAGAGTGTGACTTTGGACAATCAAAAGAGTTGATTTTGTGGAAAGGGCTATCATCTTATGCTCTTCCCTCAGTCAGCACCTCAT
GGTGTAGTCTTCTTGATGTGACTTATGTCCCTGCACAAGAAAAAGAACTTCAACAACGTCTCTGCCATTGTGATGATGGAAAAAGCACACTTCTCTGTGAAGGTGTCTTTGTTTCAAAATGGCACAC
ACTGGTTTGTAAACACAAAGGAATTTTTATGAACACAAATCATTACTACAGACAAACACATTTGTGCTGGTAACTGTGATGTTGTAATAGGAATTTGTCAACAACACAGTTTATGATCCTTTTGAACCC
TGAATTAGACTCAATTCAGGAGGAGTTAGATAAATATTTAAGAATCATACATCACCAGATGTTGATTTAGGTGACATCTCTGGCATTAATGCTTCAGTTGTAACATTTCAAAAAAGAAATTTGACCGC
CTCAATGAGGTGGCAAGAATTTAAATGAATCTCTCATCGATCTCCAAGAAGTGGAAAGTATGAGCAGTATATAAAATGGCCATGGTACATTGCGTAGGTTTTATAGCTGGCTGATTGGCCATAG
TAATGGTGACAAATTTATGCTTTGCTGTATGACCAAGTGTGCTAGTTGTGCTCAAGGGCTGTGTTGTTCTGTGGATCCTGCTGCAAAATTTGATGAAGACGACTCTGAGCGAGTGTCTCAAGGAGTCAAAAT
ACATTACACATAA
```

Protein :

```
MFVFLVLLPLVSSQCVNLTRTQLPPAYTNSFTRGVYYPDKVFRSSVLHSTQDLFLPFFSNVTWFHAIHVSQGTNGTKRFDNPVLPFNDGVVFASTEKSNIIRGWIFGTTLDSKTQSLIVNNATNV
VIKVFCEQFCNDPFLGVYHKNKNSWMESEFRVYSSANNCTFEYVSQPLFMDLEGKQGNFKNLREFVFKNIDGYFKIYKSHTPINLVRDLPGQFSALEPLVDLPIGINITRFQTLALHRSYLPDGD
SSSGWTAGAAAAYYGYLQPRFTLLKYNENGTITDAVDCALDPLSETKCTLSFTVEKGIYQTSNFRVQPTESIVRFPNIITNLCPFGEVFNATRFASVYAWNRRKISNCVADYSVLVNSASFSTFKCY
GVSPKTLNLDLCTNMYADSFVIRGDEVQRQIAPGQTKGIADYNYKLPPDFTGCVIAWNSNNLDSKVGNNYNYLYRLFRKSNLKPFERDISTEYQAGSTPCNGVEGFNCFYPLQSYGFGPTNGVGYQP
YRVVVLSEFLHAPATVCGPKKSTNLVKNKCVNFFNGLTGTGVLTESNKKFLPFQGFGRDIADTTDAVRDPQGLEILDITPCSFGGVSVITPGTNTSNQVAVLYQDVNCTEVPVAIHADQLTPTWR
VYSTGSMVFQTRAGCLIGAEHVNNSEYCDIPIGAGICASYQTQNSPRRARSVASQSI IAYTMSLGAENSVAYSNNNSIAIPTNFTISVTEILPVSMTKTSVDCTMYICGDSSTECNSLLQYGSFCT
QLNRLTGIAVEQDKNTQEVFAQVKQIYKTPPIKDFGGNFSSQILPDPSPKSPKRSFIEDLFFNKVTLADAGFIKQYGDCLGDIAARDLICAQKFNGLTVLPLLTDEMAIQYTSALLAGTITSGWTF
GAGAALQIPFAMQMAYRFNGIGVTQNVLYENQKLIANQFNSAIGKIQDSLSSASALGKLQDVVNQNAQALNTLVKQLSSNFGAISSVLNDILSRLDKVEAEVQIDRLITGRQLSLQTYVTQQLIRA
AEIRASANLAATKMSCEVLGQSKRVDFCGKGYHLMSPQSAHPGVVFLHVTYVPAQEKNFTTAPAICHDKAHFPREGVFSVNGTHMFTVQRNFYEPQIITTDNTFVSGNCDVIGIVNNNTVYDPLQ
PELDSFKEELDKEYKNHTSPDVLGDISGINASVWNIQKEIDRLNEVAKNLNESLIDLQELGKYEQYIKWPYIWLGFIAGLIAIMVMTIMLCMTSCCCLGCCSCGSCCKFDEDDSEPVLGKVK
LHYT_
```

```
print('Length of the DNA : {}'.format(len(data))\n      'Length of the protein : {}'.format(len(protein)))
```

Length of the DNA : 3822

Length of the protein : 1274