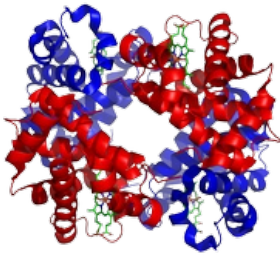


Infer mRNA from protein

Just as nucleic acids are polymers of nucleotides, **proteins** are chains of smaller molecules called **amino acids**. 20 amino acids commonly appear in every species. Just as the primary structure of a nucleic acid is given by the order of its nucleotides, the primary structure of a protein is the order of its amino acids. Some proteins are composed of several subchains called polypeptides, while others are formed of a single polypeptide. Proteins power every practical function carried out by the cell, and so presumably, the key to understanding life lies in interpreting the relationship between a chain of amino acids and the function of the protein that this chain of amino acids eventually constructs. The field devoted to the study of proteins is called proteomics.



The human hemoglobin molecule consists of 4 polypeptide chains; α subunits are shown in red and β subunits are shown in blue.

How are proteins created? The **genetic code** — discovered throughout the course of a number of ingenious experiments in the late 1950s — details the translation of an RNA molecule called **messenger RNA** (mRNA) into amino acids for protein synthesis. The apparent difficulty in translation is that somehow 4 RNA bases must be translated into a language of 20 amino acids. In order for every possible amino acid to be created, we must translate 3-nucleobase strings (called **codons**) into amino acids. Note that there are $4^3=64$ possible codons, so that multiple codons may encode the same amino acid. Two special types of codons are the **start codons** (AUG in the standard genetic code), which code for the amino acid methionine and always indicate the start of translation, and the stop codons (UAA, UAG, UGA in the standard genetic code), which do not code for an amino acid and cause translation to end.

The notion that protein is always created from RNA, which in turn is always created from DNA, forms the central dogma of molecular biology. Like all dogmas, it does not always hold. However, it offers an excellent approximation of the truth. A eukaryotic organelle called a ribosome creates peptides by using a helper molecule called **transfer RNA** (tRNA). A single tRNA molecule possesses a string of three RNA nucleotides on one end (called an anticodon) and an amino acid at the other end. The ribosome takes an mRNA molecule transcribed from DNA and examines it one codon at a time. At each step, the tRNA possessing the complementary anticodon bonds to the mRNA at this location, and the amino acid found on the opposite end of the tRNA is added to the growing peptide chain before the remaining part of the tRNA is ejected into the cell, and the ribosome looks for the next tRNA molecule.

Not every RNA base eventually becomes translated into a protein, and so an interval of RNA (or an interval of DNA translated into RNA) that does code for a protein is of great biological interest. Such an interval of DNA or RNA is called a gene. Because protein creation drives cellular processes, genes differentiate organisms and serve as a basis for heredity, or the process by which traits are inherited.

Assignment

The 20 commonly occurring amino acids are abbreviated by using 20 letters from the English alphabet (all letters except for B, J, O, U, X, and Z). Protein strings are constructed from these 20 symbols. An **RNA codon table** dictates the details regarding the encoding of specific codons into the amino acid alphabet. The codon table shown below gives the mapping used by the standard genetic code. However, there are alternative genetic codes that show slight variations on the mapping scheme. Stop codons do not code for an amino acid, and are indicated by the word `Stop` instead of an amino acid symbol.

UUU F	CUU L	AUU I	GUU V
UUC F	CUC L	AUC I	GUC V
UUA L	CUA L	AUA I	GUA V
UUG L	CUG L	AUG M	GUG V
UCU S	CCU P	ACU T	GCU A
UCC S	CCC P	ACC T	GCC A
UCA S	CCA P	ACA T	GCA A
UCG S	CCG P	ACG T	GCG A
UAU Y	CAU H	AAU N	GAU D
UAC Y	CAC H	AAC N	GAC D
UAA Stop	CAA Q	AAA K	GAA E
UAG Stop	CAG Q	AAG K	GAG E
UGU C	CGU R	AGU S	GGU G
UGC C	CGC R	AGC S	GGC G
UGA Stop	CGA R	AGA R	GGA G
UGG W	CGG R	AGG R	GGG G

When researchers discover a new protein, they would like to infer the strand of mRNA from which this protein could have been translated, thus allowing them to locate genes associated with this protein on the genome. Unfortunately, although any RNA string can be translated into a unique protein string, reversing the process yields a huge number of possible RNA strings from a single protein string because most amino acids correspond to multiple RNA codons. For a given protein string, determine the total number of different mRNA strings from which the protein could have been translated. In order to do this, you proceed as follows:

- Write a function `codontable` that takes the location of a text file as an argument. This text file must contain the mapping from codons to amino acids as used by a genetic code, in the format as shown in the table above. The function must read the mapping from the text file, and return it as a dictionary (that maps each of the 64 possible codons into their corresponding amino acid). Stop codons should be mapped onto an asterisk (*).
- Write a function `reverse_codontable` that takes a dictionary having the form of the dictionaries returned by the function `codontable`. The function must return a new dictionary that maps each of the 20 amino acids and the stop codon (represented by an asterisk: *) onto the set of codons that translate to this amino acid/stopcodon.
- Write a function `mRNA` that takes two arguments: a protein string and a codon table. The codon table must be passed as a dictionary having the form of the dictionaries returned by the function `codontable`. The function must return the total number of different mRNA strings from which the protein could have been translated. This number can be written as the product of the number of codons that gets mapped onto each amino acid in the protein string. Take into account that protein synthesis ends at a stop codon (which is not explicitly indicated at the end of the protein string that is passed to the function). This means you also have to multiply by the number of stop codons. An example: according to the codon table of the standard genetic code there 12 mRNA strings that translate to the protein MA, because there is a single codon translating to M, four to A and there are three stop codons: $1 \times 4 \times 3 = 12$

