

COM.e PROJECT: Investigating DNA in Python

Utsav Choudhury L24i

TARGET ASSESSMENT LEVEL: 3

1 Preface

[Jump ahead to the specification if this is boring.]

In Marjo Petäjäaho's BIO4 course, we are studying DNA and protein synthesis. Interestingly, DNA is the body's way of conveying information in code. After realizing that this project's purpose was text processing, I decided that processing theoretical DNA strands might be interesting. Of course, the project is highly theoretical, not taking mutations and specific complex processes into account; hence, it cannot be used in real-world applications. Nevertheless, it offers quite an interesting insight into how the body might theoretically process DNA.

2 Background Theory

DNA is a very complex molecule; however, the way it conveys information is quite understandable. DNA strands contain the organic bases adenine, cytosine, guanine, and thymine, commonly abbreviated as "A," "C," "G," and "T." These organic bases exhibit a characteristic: complementary pairing. A always bonds with T, and G always bonds with C.

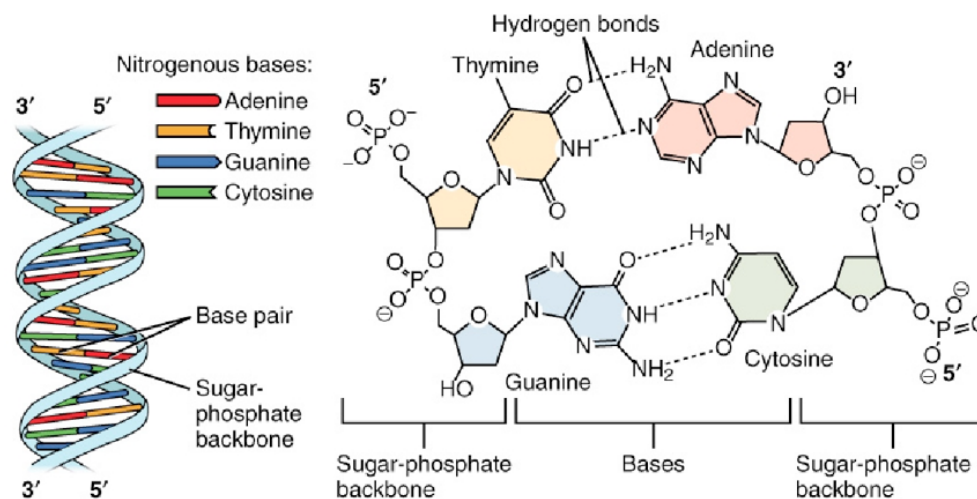


Figure 1: DNA Molecule [LibreTexts]

Complementary pairing can be used to our advantage. When DNA replicates, its two strands unwind. When one strand becomes split, transcription can take place. The singular strand pairs with its complementary counterpart. However, thymine is replaced by uracil (T becomes U). The newly formed strand is called mRNA, which is composed of many different triplets called codons. These codons are either code for a unique amino acid or punctuation. "codon tables" organize this information (Figure 3).

Punctuation means that codons can indicate "START" or "STOP," signaling when to start or stop the amino acid chain. Translation is the process where amino acids connect to their codon counterparts and form amino acid chains. When these amino acids link up, they form proteins, which are fundamental to all bodily processes. This whole process is simply an algorithm where characters and strings are processed, making it relevant for a CS project [Britannica].

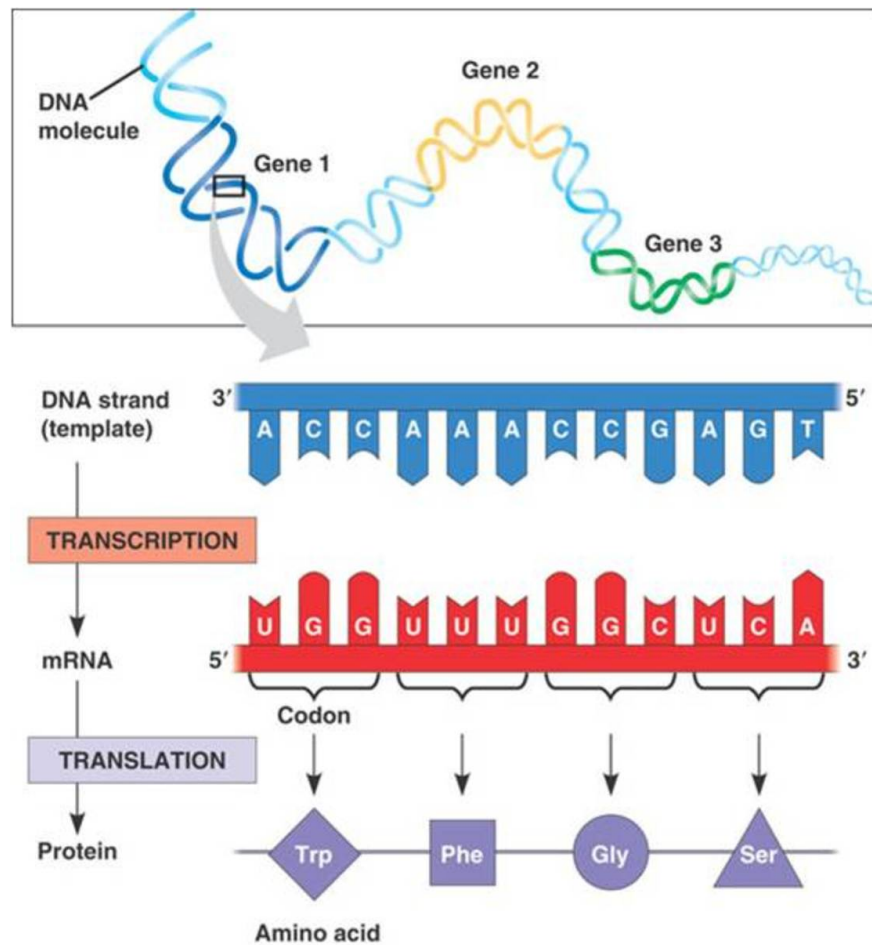


Figure 2: Transcription and Translation [Owlcation]

		second letter					
		U	C	A	G		
first letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA stop UAG stop	UGU } Cys UGC } UGA stop UGG Trp	U C A G	third letter
	C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }	U C A G	
	A	AUU } AUC } Ile AUA } AUG Met	ACU } ACC } Thr ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G	
	G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } Gly GGA } GGG }	U C A G	

Figure 3: Codon Table [Vedantu]

3 Specification

3.1 What does the program do?

The program takes in different pieces of data. First, it takes in a filename where output will be stored (one for text in a .txt file). Then, the program asks the user whether they want a random DNA sequence or a custom one. If the user wants a random sequence, they will be asked for a length. A random sequence will then be generated. If the user chooses custom, they will be prompted to enter a custom sequence. In both cases, the user must enter their input in a specific format, following certain criteria. The random or custom sequence will then be transcribed (see Figure 2). During transcription, start codons are added, and stop codons are deleted. After transcription, the mRNA will be translated to its amino acid counterparts (see Figure 3). It does this with a dictionary in the file `Codon_Table.py`. Finally, the amino acids will be strung together. All of these conversions are documented and written into a .txt file with the name that the user inputted.

3.2 Data Format

The user can select any file name for their report (except "" or the absence of value). File extensions do not matter.

The user must enter "custom" or "random", where capitalization does not matter.

If "random" is chosen, a length l_1 must be selected such that:

$$9 \leq l_1 \leq 30 \quad \text{and} \quad l_1 \equiv 0 \pmod{3}.$$

If "custom" is chosen, the string must be composed only of the letters "A," "C," "G," and "T," and a length l_2 must be selected such that:

$$9 \leq l_2 \leq 30 \quad \text{and} \quad l_2 \equiv 0 \pmod{3}.$$

Capitalization does not matter in any case.

4 Correctness

4.1 Ideal Test Case

In this project's GitHub repository, I have linked some files documenting input and output cases. However, it must be noted that input occurs through the Python console, not through file uploads. Obviously, the date will vary, as it changes.

4.1.1 Example Console Input

```
1 C:\Users\utsav.choudhury\PycharmProjects\CS_PROJECT\.venv\Scripts\python.exe
   C:\Users\utsav.choudhury\PycharmProjects\CS_PROJECT\CS_PROJECT_COPY.py
2 Enter a filename for your report
3 No file extension needed. However, please note that if a file with this name
   already exists, it will be overwritten: test
4 Enter 'CUSTOM' for a custom sequence or 'RANDOM' for a random sequence: custom
5 Enter your custom DNA sequence (A, C, G, T). Length must be 9-30 and divisible by
   3: AGTCCCTCT
6 Process completed. Check test.txt for information.
7
8 Process finished with exit code 0
```

4.1.2 Example Text File Output

```
1 REPORT
2
3 DATE: 2025-01-31
4
5 RANDOM BASE ARRANGEMENT: AGTCCCTCT
6
7 RNA sequence: UCAGGGAGA
8
9 Codons: ['AUG', 'UCA', 'GGG', 'AGA']
10
11 Amino Acids: ['Met', 'Ser', 'Gly', 'Arg']
12
13 FINAL CHAIN: Met-Ser-Gly-Arg
14
15 NOTE: NOT A FULLY ACCURATE REPRESENTATION. NOT APPLICABLE TO REAL LIFE CONTEXTS.
16
17 UTSAV CHOUDHURY 2025
```

4.2 Resource Management

The input file is opened using a `with-statement`, and will therefore be closed automatically.

References

- [1] Britannica. *DNA: Discovery, Function, Facts, and Structure*. 2018.
<https://www.britannica.com/science/DNA>
- [2] Britannica. *Metabolism: The Synthesis of Macromolecules*. <https://www.britannica.com/science/metabolism/The-synthesis-of-macromolecules>
- [3] LibreTexts. *9.1: The Structure of DNA*. 2017. https://bio.libretexts.org/Bookshelves/Introductory_and_General_Biology/Concepts_of_Molecular_Biology/9.01:_The_Structure_of_DNA
- [4] Rhys Baker. *Protein Production: A Simple Summary of Transcription and Translation*. 2012.
<https://owlcation.com/stem/protein-production-a-step-by-step-illustrated-guide>
- [5] Vedantu. *How Do You Read a Codon Table?* <https://www.vedantu.com/question-answer/how-do-you-read-a-codon-table-class-11-biology-cbse-60ed1177a183f842ef1654c2>