

Lecture 10 - String processing.

Reminder - no class on Thursday, 17th - There will be labs.

Let's start with some data

[illegible]

This is from the file:

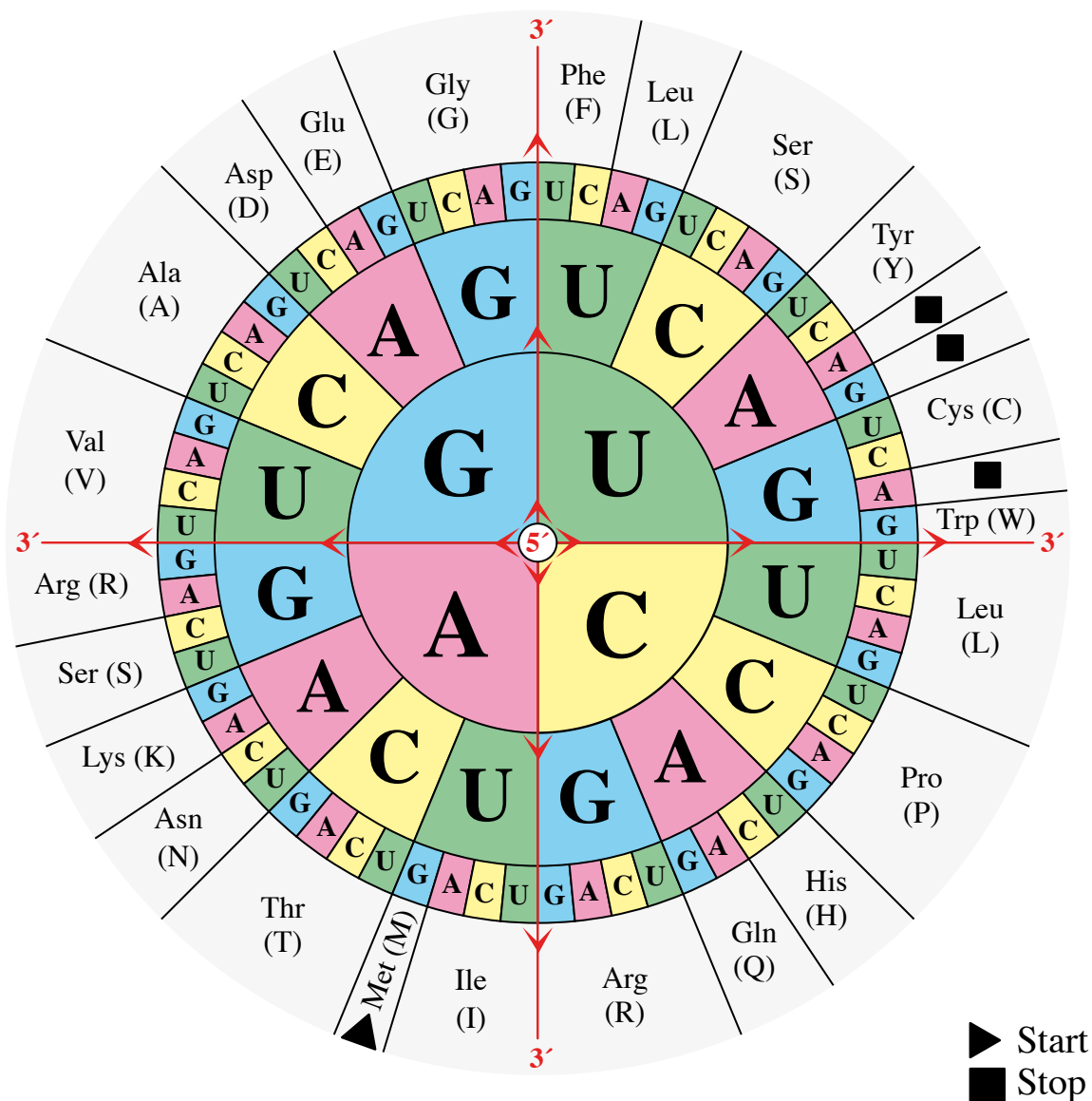
GCF_000001405.39_GRCh38.p41214_genomic.s44a-21.fna

What we want to do is convert this to Amino Acids.

1st thing to note is that the data is upper/lower case. Our life will be simpler if we make it all lower case.

2nd thing is that this is DNA data. It has "T" in it. Let's make our conversion from messenger RNA that uses U instead of T. So we need a function to start cleaning up the data.

Background on generation of Amino Acids from RNA



Amino Acid

Chart

Or if you prefer a table: https://en.wikipedia.org/wiki/DNA_and_RNA_codon_tables

Let's do a Google search and you might find:

https://github.com/T101J/Translating_RNA_to_Protein

this is cool. A click on the code, and ... It kind of is what is in the homework ... that has promise. What is the license? OOps! no LICENSE file. Look at the code - they could have put the license in the file. Noope! That means that this is proprietary code - "All Rights Reserved".

Nice chart - but not in the class notes! This is *Proprietary* code and the nice blue chart is also!

The code. Useful to read - but it is not exactly what we want for class. Remember that we want to scan across to find a start-codon, then take 3 at a time and look them up. This assumes that you have already found the start - and just takes 3 at a time. It has the processing for the stop codon. That is good. So you can't copy this code but you can read it and learn from it.

There are about 5 other places where I found the exact same code on the web. Nice articles like:

<https://towardsdatascience.com/starting-off-in-bioinformatics-rna-transcription-and-translation-aaa7a91db031> It has a good explanation of how to do this - but not exactly our requirements. (It also has the same chart - from what I can tell the chart is actually copyright a biomedical company.)

Note: going from DAN to RNA change 'T' to 'U'

First convert to one string

We have a list of strings. We need just one. "How do you convert a list of strings to a single string?"

Join them together. There is a string operation to do this called join. Join gives us the ability to join with a separator between each string.

```
a = [ "aa", "bb", "cc" ]

print ( ", ".join(a) )
```

We can also join with a 0 length separator (this is what we will want in our code).

```
a = [ "aa", "bb", "cc" ]

print ( "".join(a) )
```

So we start out with the data in the program:

```
dna_input = [
"taaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
"accctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
"cctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
"taaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
"caaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
"taaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
"AGCCGGCCGCCGCCGGGCTGACCTGAGGAGAACTGTGCTCCGCCTTCAGAGTACCACCGAAATCTGTGCAGAGGAc",
"aacgcagctccgcctcgcggtGCTctccgggtctgtgctgaggagaaacgCAACTCCGCCGTTGCAAAGGCGcgccgcgc",
"cggcgcaggcgagagaggcgccgcgcggcgaggcgagagaggcgUGAcgcgcggcgaggcgagagaggcg",
"gccgcgcggcgaggcgagagaggcgccgcgcggcgaggcgagagaggcgccgcgcggcgaggcgagag",
]

print ( dna_input )
```

And then join it to make 1 string.

How to get to a section of a string

Let's say we have a string

```
bb = "abc123def"
```

and we want to get to the 123 part of the string. We can pull out just those three characters with:

```
x3 = bb[3:6]
```

So a little program that will do this is:

```
bb = "abc123def"  
x3 = bb[3:6]  
print ( "x3 = {}".format(x3) )
```

Replace - modifying a string

Python has a builtin string function to replace strings.

```
1: s = "aBBcBBd"  
2: t = s.replace("B","-")  
3: print ( t )
```

Use replace to remove a string.

```
1: # Remove Bs  
2:  
3: s = "aBBcBBd"  
4: t = s.replace("B","")  
5: print ( t )
```

replace with a longer string

```
1: # Replace Bs with a longer string.  
2:  
3: s = "aBBcBBd"  
4: t = s.replace("B","--hi--")  
5: print ( t )
```

Replace just the 1st occurrence

```
1: # Replace just the 1st B with a longer string.  
2:  
3: s = "aBBcBBd"  
4: t = s.replace("B","--hi--",1)  
5: print ( t )
```

Convert to upper/lower case

To convert to lower case

```
1: a = "This is A Mixed Case"
2: b = a.upper()
3: print ( b )
```

To convert to upper case

Functions for Lab 05

First is a function to take our input, convert from upper/lower case to lower, then replace the t s with u s.

Second is to use a dictionary lookup.

To create a dictionary in Python

```
d = { "name": 1, "name2": 2 }
```

or to lookup other strings

```
d = { "in1": "out1", "in2": "out2" }
```

So we can lookup our 3 letter codon with a dictionary.

```
# RNA codon table
rna_codon = {
    "aaa" : "K", "aac" : "N", "aag" : "K", "aau" : "N", "aca" : "T", "acc" : "T", "acg" : "T",
    "acu" : "T", "aga" : "R", "agc" : "S", "agg" : "R", "agu" : "S", "aua" : "I", "auc" : "I",
    "aug" : "M", "auu" : "I", "caa" : "Q", "cac" : "H", "cag" : "Q", "cau" : "H", "cca" : "P",
    "ccc" : "P", "ccg" : "P", "ccu" : "P", "cga" : "R", "cgc" : "R", "cgg" : "R", "cgu" : "R",
    "cua" : "L", "cuc" : "L", "cug" : "L", "cuu" : "L", "gaa" : "E", "gac" : "D", "gag" : "E",
    "gau" : "D", "gca" : "A", "gcc" : "A", "gcg" : "A", "gcu" : "A", "gga" : "G", "ggc" : "G",
    "ggg" : "G", "ggu" : "G", "gua" : "V", "guc" : "V", "gug" : "V", "guu" : "V", "uaa" : ".",
    "uac" : "Y", "uag" : ".", "uau" : "Y", "uca" : "S", "ucc" : "S", "ucg" : "S", "ucu" : "S",
    "uga" : ".", "ugc" : "C", "ugg" : "W", "ugu" : "C", "uua" : "L", "uuc" : "F", "uug" : "L",
    "uuu" : "F"
}
```

We can see if a value is in the lookup table with an "if"

```
from rna_lookup import rna_codon

if "xyz" in rna_codon:
    print ( "found" )
else:
    print ( "not found" )
```

```
if "gac" in rna_codon:  
    print ( "found value = {}".format(rna_codon["gac"] ) )  
else:  
    print ( "not found" )
```

A Program to go from DNA to Amino Acids.

Program State

```

1: import conv_t_to_u
2: import rna_lookup
3:
4: dna_input = [
5: "taaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
6: "accctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacc",
7: "cctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacc",
8: "taaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
9: "caaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
10: "taaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta",
11: "AGCCGGCCCGCCCGGGTCTGACCTGAGGAGAACTGTGCTCCGCCTTCAGAGTACCACCGAAATCTGTGCAGAGGAC",
12: "aacgcagctccgccctcgcggtGCTCtccgggtctgtgctgaggagaacgCAACTCCGCGTTGCAAAGGCGcgccgcgc",
13: "cggcgcgaggcgagagaggcgccgcgcccggcgaggcgagagaggcgUGAcgcgccggcgaggcgagagaggcg",
14: "gccgcgcccggcgaggcgagagaggcgccgcgcccggcgaggcgagagaggcgcgccggcgaggcgagag",
15: ]
16:
17: amino_string = ""
18:
19: rna = conv_t_to_u.conv_t_to_u ( "".join(dna_input) )
20:
21: st = "before"
22: i = 0
23: while i < len(rna)-2:
24:     three = rna[i:i+3]
25:     # print ( "st {} three >{}< at i={}".format(st, three,i))
26:     amino = rna_lookup.rna_to_amino_acid(three)
27:     if amino == '!' :
28:         print ( "Invalid 3 letter sequence {} at {}".format(three, i ) )
29:         break
30:     if st == "before" and amino == "M" :
31:         st = "encode"
32:         i = i + 3
33:     elif st == "encode" and amino == ".":
34:         st = "before"
35:         print ("Protein : ", amino_string)
36:         i = i + 3
37:     elif st == "encode" :
38:         amino_string = amino_string + amino
39:         i = i + 3
40:     else:
41:         i = i + 1

```

Copyright

Copyright © University of Wyoming, 2021.