

Project BIO - Genetic Code in R

Carlos Adir **Ely Murussi Leite**

June 18, 2020

1 Introduction

This project has the objective of using the concepts in R learned during the classes to apply in biology: Make the processus of transcription, translation and alignment.

In the first moment there are some exercises to reminder of the language in R, the use of a GUI(Graphical User Interface) known as Shiny to help in data analysis.

The 2 last exercises(7.5 and 7.6, bonus) were not made.

2 Remider on R

Exercise 2.1

- a) Write a function **compte** which takes two arguments: a sequence of characters and a letter, and which returns the number of occurrences of the letter in the sequence.

We code the folowed function

```
1 compte <- function(sequence, letter){
2     return(sum(sequence == letter))
3 }
```

- b) Apply your function to find the number of occurrences of the letter “a” in the biological sequence

```
sequence<-c("a","a","t","g","a","g","c","t","a","g","c","t","g")
```

Using the code below, we get

```
1 sequence <- c("a","a","t","g","a","g","c","t","a","g","c","t","g")
2 qtt_a <- compte(sequence, "a")
3 print(qtt_a)
```

OUTPUT:

```
[1] 4
```

- c) Use a for loop to find the numbers of occurences of each letter in the sequence.

```
1 bases = c("a","c","g","t")
2 composition = rep(0,4)
3 for(i in 1:4){
4     composition[i] = compte(sequence, bases[i])
5 }
6 print(composition)
```

OUTPUT:

```
[1] 4 2 4 3
```

3 Dealing with sequences

Exercise 3.1

- a) What are the **seqtype** and **forceDNAtolower** options for?

If we read the documentation, we get:

seqtype: the nature of the sequence: DNA or AA, defaulting to DNA

forceDNAtolower: whether sequences with `seqtype == "DNA"` should be returned as lower case letters

- b) Print the first 50 nucleotides of your sequence:

For do that, we use the code:

```
1 library("seqinr")
2 filename <- "~/ECN/Biologie/13-TP-ProjR/Students/Mysterious_seq.txt"
3 sequences <- read.fasta(file = filename, seqtype = "DNA", forceDNAtolower = T)
4 mysequence <- sequences$seq0
5 print(mysequence[1:50])
```

OUTPUT:

```
[1] "g" "c" "g" "c" "c" "g" "g" "g" "a" "a" "g" "a" "a" "a" "g" "g" "a" "a" "c"
[20] "a" "t" "g" "g" "c" "t" "c" "c" "t" "g" "a" "g" "g" "c" "g" "c" "a" "c" "a"
[39] "g" "c" "g" "c" "c" "g" "a" "g" "c" "g" "c" "g"
```

- c) Print the length of the sequence.

To do so, we just need

```
1 print(length(mysequence))
```

OUTPUT

```
[1] 5080
```

- d) Print the number of occurrences of nucleotides (a, t, c, g).

To do so, we just need

```
1 bases <- c("a", "c", "g", "t")
2 composition <- matrix(integer(4), ncol=4)
3 for(i in 1:4){
4   composition[i] <- sum(mysequence == bases[i])
5 }
6 colnames(composition) <- bases
7 print(as.table(composition))
```

OUTPUT:

	a	c	g	t
A	1391	1079	1202	1408

Exercise 3.2

- a) Why is the computation of GC-content relevant in molecular biology?

It give us the stability of the chain, because as the link $G \equiv C$ is stronger than the link $A = T$. We call that *thermostability*. Its entirely correlated with the quantity of hidrogen bonds.

Other reason in use CG content instead of AT , is because we can use the same concept in RNA. That means, we just need to say CG content instead of AT/AU content.

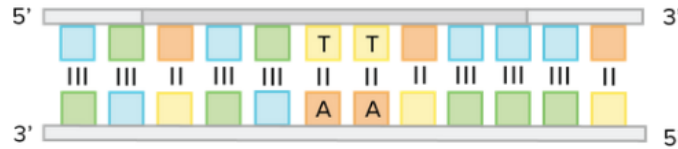


Figure 1: Links between the nucleotides
Source: Khan Academy

b) Compute the GC-content of your sequence.

```
1 GCcount = composition[2] + composition[3]
2 GCcontent = GCcount/length(mysequence)
3 print(paste("Value GC-content:", GCcontent))
```

OUTPUT:

```
[1] "Value GC-content: 0.449015748031496"
```

Exercise 3.3

a) By convention, in which direction ($5' \rightarrow 3'$ or $3' \rightarrow 5'$) do we write DNA sequences? What do $5'$ and $3'$ represent?

The convention is using the direction $5' \rightarrow 3'$.

As we know, the molecules are in 3D and it's hard to say what is below, what is above, and so on. For that, the convention is numerate using the number of the carbon that the structure is linked, as seen in the Figure (2).

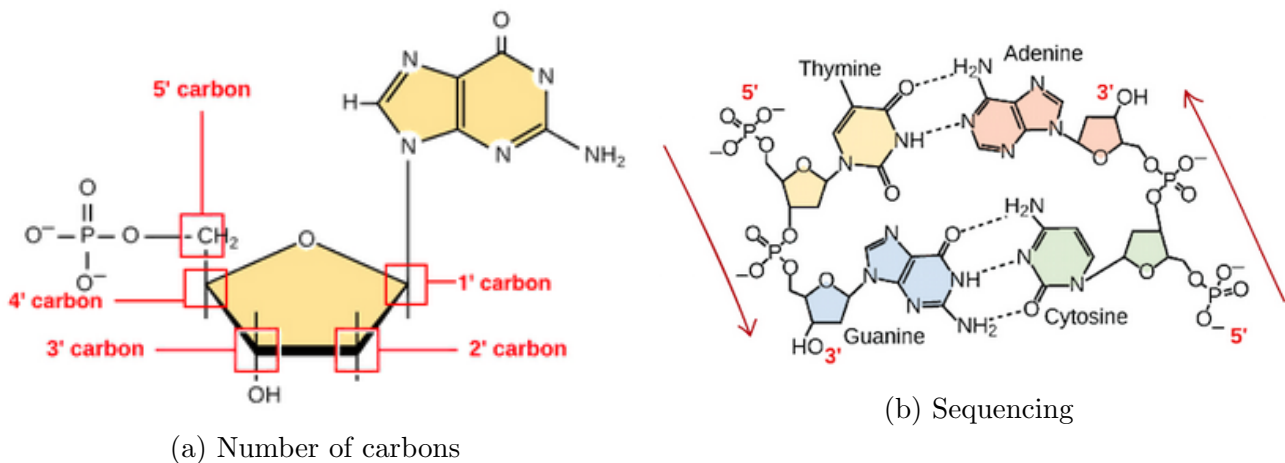


Figure 2: Structure
Source: Khan Academy

b) Write a function complementary that takes a sequence as input and returns its complementary sequence (i.e. a becomes t; t, a; c, g; and g, c).

```
1 complementary <- function(seq){
2   bases <- "acgt"
3   bases_rev <- "tgca"
4   complement_seq <- chartr(bases, bases_rev, seq)
5   return(rev(complement_seq)) # To adjust from 5->3 again
6 }
```

c) Compute and print the reverse complementary of the reference sequence.

```
1 mysequence_complement <- complementary(mysequence)
2 print("Complementary sequence:")
3 print(mysequence_complement[1:25])
```

OUTPUT:

```
[1] "Complementary sequence:"  
[1] "t" "t" "c" "t" "g" "c" "t" "a" "c" "t" "a" "t" "a" "t" "c" "a" "a" "g" "t"  
[20] "t" "t" "t" "t" "a" "t"
```

4 Introduction to Shiny

Exercise 4.1

- a) Add the function *compte* from *exercise 2.1* in the Shiny application (*util.R* file)

As we have already done the code in 2.1, it's only copy and paste.

- b) Run the application and test with *Mysterious_seq.txt*. Alternatively, you can test it with a dataset of 31 PKD2 mRNA sequences (*PKD2_mRNA_sequences.txt*).

When we run, we can select a file. If we select the dataset of 31 sequences given, we obtain the Figure (3)

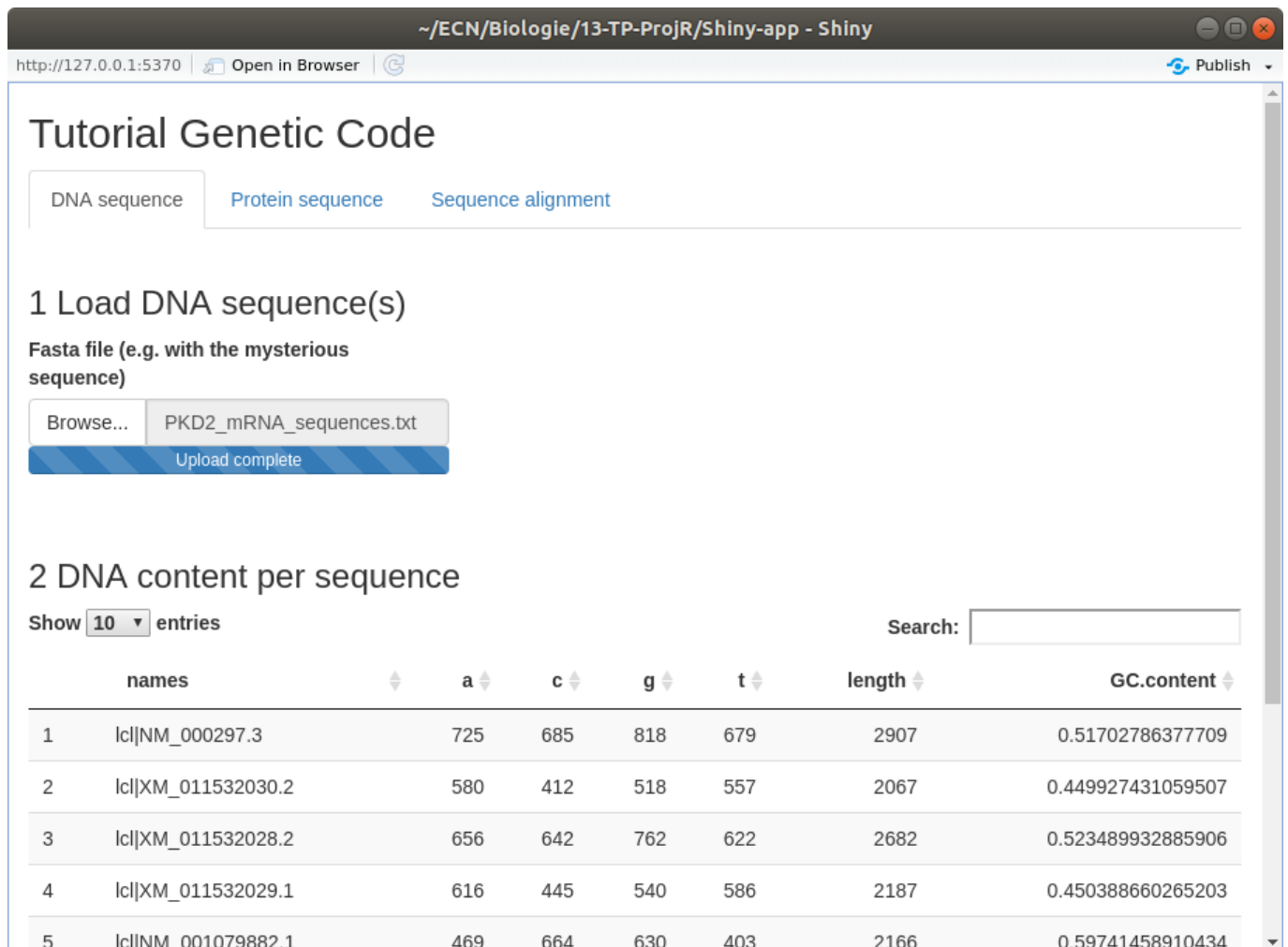


Figure 3: First screen

5 Transcription

Exercise 5.1

- a) What is the transcription mechanism for? Describe briefly the main steps of transcription.

The transcription is the processing of getting the DNA's sequence, copying it and making a RNA molecule.

As we can see in the Figure (4), it's the processus os creating the RNA.

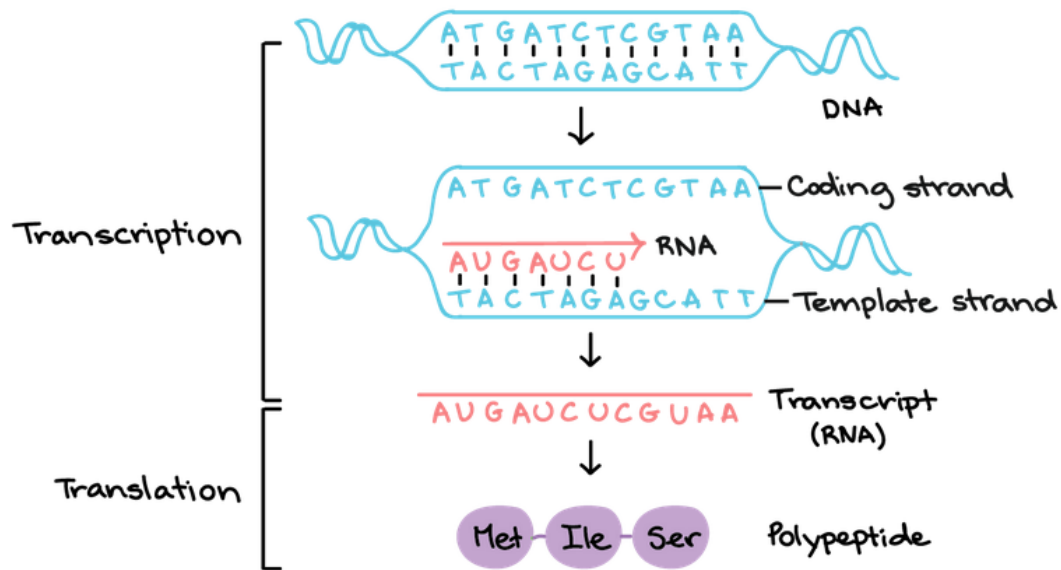


Figure 4: Transcription Mechanism
Source: Khan Academy

The processus follow the order:

Begin of the transcription The enzyme polymerase goes to the DNA in a region that we call **promoter**.

Elongation The enzyme begins the elongation: The RNA grows because the enzyme puts more and more nucleotides while it walks on the DNA, as we can see in the Figure (5).

End of transcription It keeps going until it reaches the signal to **stop**. When it happens, the polymerase writes the sequence known as terminator

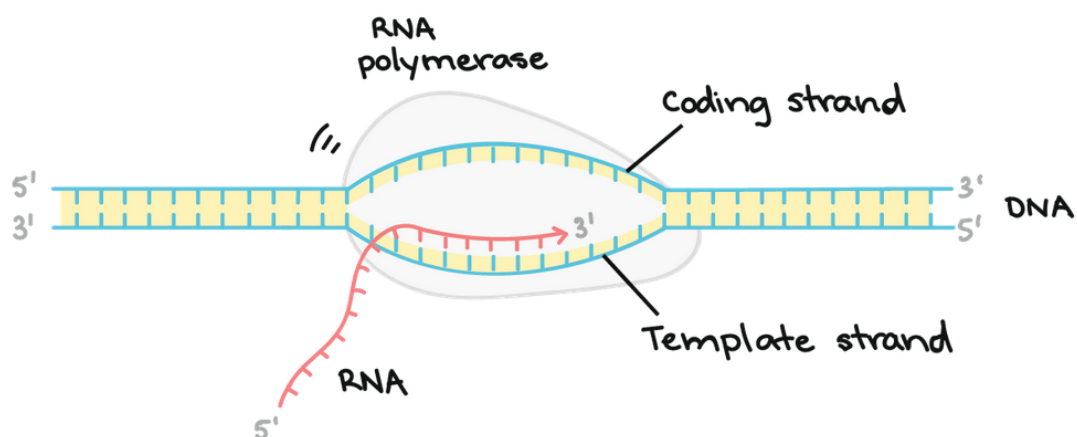


Figure 5: Polimerase produces pre-mRNA from DNA
Source: Khan Academy

b) What is(are) the difference(s) between **mRNA** and **cDNA**?

The first difference that we have, is that *RNA* has the base *U* instead of *T* that we find in the *DNA*. The second difference is that it's a single-stranded similar to the *DNA* sequence

We have that *mRNA* is the *RNA* produced by transcription, and *cDNA* is the *DNA* synthesized using the *mRNA* sequence.

c) What enzyme catalyzes the **mRNA** to **cDNA** reaction?

While we have the polymerase that produces *mRNA* from the *DNA*(Figure 5), there is the **enzyme reverse transcriptase** that gets the *mRNA* and transforms into *cDNA*(complementary *DNA*).

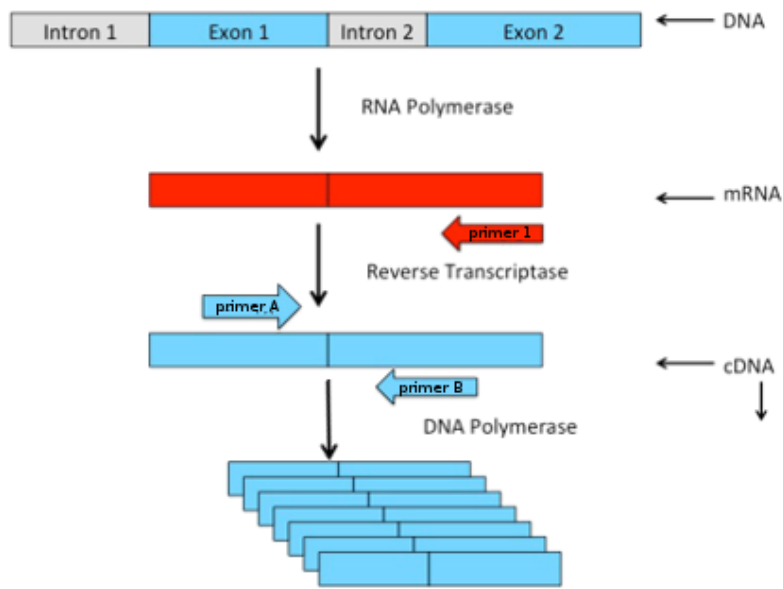


Figure 6: Reverse trans produces cDNA from mRNA

Source: Khan Academy

d) What is the difference between cDNA and the original DNA sequence (if any)?

During the process of transcription of the *DNA*, we get the **pre-mRNA**. During its maturation, introns are thrown away, remaining only the exons(as we can see in the Figure 7), resulting in the **mRNA** that is converted into **cDNA**

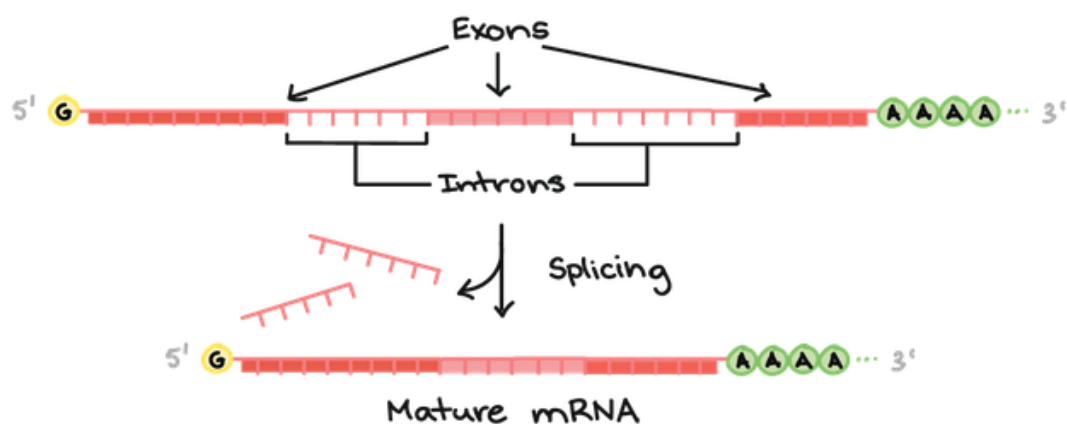


Figure 7: Slicing

Source: Khan Academy

So, we could say that cDNA is a mix of slices(in order) of the original *DNA*.

But we have to say that cDNA and DNA are different only in eukaryote cells, because the prokaryotes are not able to remove the introns.

- e) Write a function **transcribe** that converts the mysterious cDNA sequence into RNA.

```
1 transcribe <- function(cDNA){  
2   RNA <- chartr("t", "u", cDNA)  
3   return(RNA)  
4 }
```

6 Translation

Exercise 6.1

- a) What is the translation mechanism for? Describe briefly the main steps of translation.

In short, that's showed in the Figure (8)

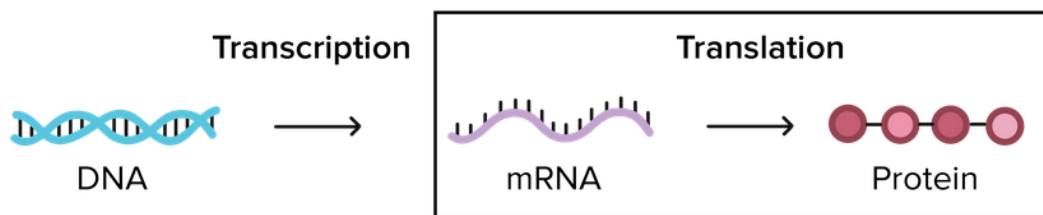


Figure 8: Translation Mechanism

Source: Khan Academy

We can say that the **translation mechanism** is the processus of getting the **mRNA** and transforming it into a chain of amino-acids that we call as a **protein**.

The processus is

Begin The ribosome gets the **mRNA** and the first tRNA(Transfer RNA).

Middle Ribosome gets amino-acids from tRNAs and link them to make a chain

End The sequence of amino-acids is released to make its job as protein.

- b) What is the genetic code? What are its main characteristics (number of nucleotides used for decoding, sliding properties of the ribosome, etc.)?

The genetic code are the rules to transform the genetic material(DNA or RNA) into proteins (sequence of peptides/amino-acid).

As there is 4 nucleotides(A, C, G, T/U), it's possible to create 64 combinations(called codons), as we can see in the Figure (9).

As 3 of them is a code to **Stop**, the 61 remaining are translated into 20 amino-acids(including the **Met**).

During the translation, the ribosomes can slide over the non coding sequence(but keeps in the **ORF** region) to synthesize a single protein. We call that process as translation bypassing.

- c) How many amino acids are there in the standard genetic code? Why is the genetic code defined as “redundant”?

As seen in the Figure (9), the same amino-acid can be given from different codons.

		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
	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 9: 64 codons, 20 amino-acids
Source: Khan Academy

That say us that some information is irrelevant. For example, if we have the sequence CU*, doesn't matter the value of *, because we already know that it's going to be translated into **Leu**. We can put them in a sequence and get:

- CU* → Leu
- CC* → Pro
- CG* → Arg
- UC* → Ser
- GU* → Val
- GC* → Ala
- GG* → Gly
- AC* → Thr

d) What is an “open reading frame” (ORF)? How many ORF is there for one sequence?

The **ORF** is a part of the **DNA** that when translated into amino-acids, doesn't have any **X**(stop codon).

As the amino-acid is built with 3 nucleotides sequence, it's possible to:

Read in the original sequence order and

- Not slide
- Slide 1 nucleotide
- Slide 2 nucleotides

Read in the reverse sequence order and

- Not slide
- Slide 1 nucleotide
- Slide 2 nucleotides

So, there is 6 possibilities to read and so get 6 different ORF for the same sequence of DNA.

e) What is the amino acid acting as the Start codon? Which codon correspond to this amino acid?

The **Start** codon is named by **Met**, or just **M**

It's Constitution, as we can see in the Figure (9) is **AUG**.

f) Which are the Stop codons? Which one-letter code is used here to describe Stop codons?

The **End** codon is named by **Stop**, or just **X**

It's constitution, as we can see in the Figure (9) can be one of three: **UAA**, **UAG** or **UGA**.

Exercise 6.2

In your working R script, Write a function *dna2peptide* that translates a DNA sequence into an amino acid chain.

```
1 codon_filename <- "~/ECN/Biologie/13-TP-ProjR/Code_shortcut/Genetic_code.txt"
2 codon.table <- read.table(file = codon_filename, col.names = c("codon", "aa", "letter"),
3                           stringsAsFactors = F)
4 codon.table$codon <- tolower(codon.table$codon) # use only lower case characters
5 rownames(codon.table) <- codon.table$codon # use codons as table keys
6 head(codon.table)
7
8 dna2peptide <- function(sequence, frame=0, sens='F'){
9   if(sens == 'R'){
10     sequence <- complementary(sequence)
11   }
12   step = 3
13   len <- length(sequence)
14   l <- step*((len - frame)%/%step)
15   init = frame + 1
16   end = l + frame
17   seq_coupe <- split(sequence[init:end], ceiling(seq(l)/step))
18   for( i in 1:length(seq_coupe) ){
19     seq_coupe[i] <- paste(unlist(seq_coupe[i]), collapse='')
20   }
21   peptides <- codon.table[unlist(seq_coupe), "letter"]
22   return( paste(peptides, collapse = "") )
23 }
24 sequence <- c('a', 't', 'g', 't', 't', 'c', 't', 't', 't', 'a')
25 dna2peptide(sequence, 0, 'F')
```

OUTPUT:

```
[1] "MFF"
```

Exercise 6.3

Define a function *find.mysterious.proteins* function that takes a cDNA sequence as input and returns a list of proteins as defined above. Apply this function to the mysterious sequence, save the resulting sequences in a text file and join it to your report.

```
1 find.mysterious.proteins <- function(cDNA){
2   output_file <- "~/ECN/Biologie/13-TP-ProjR/Students/proteins.txt"
3   pattern = 'M[ACDEFGHIKLMNPQRSTVWY]{80,}(X|$)'
4   proteins <- c()
5   step <- 3
6   for( frame in 0:(step-1) ){
7     for( sens in c("L", "R") ){
8       supposed_peps <- paste(dna2peptide(cDNA, frame, sens), collapse = "")
9       reg <- gregexpr(pattern, supposed_peps)
10      match <- regmatches(supposed_peps, reg)[[1]]
11      proteins <- c(proteins, match)
12    }
13  }
14  write(proteins, file = output_file)
15  return(proteins)
16 }
17
18 library("seqinr")
19 filename <- "~/ECN/Biologie/13-TP-ProjR/Students/Mysterious_seq.txt"
20 sequences <- read.fasta(file = filename, seqtype = "DNA", forcedNAtolower = T)
21 mysequence <- sequences$seq0
22 proteins <- find.mysterious.proteins(mysequence)
23 print(proteins)
```

OUTPUT:

```
[1] MVNSSRVQPQPGDAKRPPAPRAPDPGRLMAGCAA VGASLAAPGGLCEQRGLEIEMQRIRQAAARDPPAGAAASP
SPPLSSCSRQAWSRDNPGFEAE EEEVEEGEGMVVEMDVEWRPGSRRSAASSAVSSVGARSRLGGYHGAGHP
SGRRRRREDQGPPCPSVGGG DPLHRHLPLEGQPPRVAWAERLVRGLRGLWGTRLMEESSTNREKYLKSVLREL
V TYLLFLIVLCILTYGMSSN VYYYTRMMSQLFLDTPVSKTEKTNFKTLSSMEDFWKFTEGSLLDGLYWKMQPSNQ
TEADNRSFIFYENLLLGVPRI RQLRVNRGSCSIPQDLRDEIKECYDVYSVSSEDRAPFGPRNGTAWIYTSEKDLN
GSSHWGIATYSGAGYYLDLS RTREETA AQVASLKKNVWLDRGTRATFIDFSVYNANINLFCVVRLLEVPATGG
VIPSWQFQPLKLIRYVTTFD FLAACEIIFCFIFIVVVEEILEIRIHKLHYFRSFWNCLDVVIVVLSVVAIGIN
IYRTSNVEVLLQFLEDQNT FPNFEHLAYWQIQFN NIAAVTVFFVWIKLFKFINFNRMTMSQLSTTMSRCAKDLFGF
AIMFFIIFLAYAQLAYLVFG TQVDDFSTFQECIFTQFRIILGDINF AEIEEANRVLGP IYFTTFVFFMFILLNM
```

```

FLAIINDTYSEVKSDLAQKKAEMELSDLIRKGYHKALVKLKLKNTVDDISESLRQGGGKLNFDLRQDLKGKGH
TDAEIEAIFTKYDQDGDQELTEHEHQMRDDLEKEREDLDLDHSSLPRPMSSRSFPRSLDDSEDDDEDSGHSSR
RRGSISSGVSYEEFQVLVRRVDRMEHSIGSIVSKIDAVIVKLEIMERAKLRREVLGRLLDGAEDERLGRDSEI
HREQMERLVREELERWESDDAASQISHGLGTPVGLNGQPRPRSSRPSSSQSTEGMEGAGGNGSSNVHVX
[2] MWYPCVSTMCHISIMVTRKSLYDTFYCTVFHFTCKILQNSSFPLINYIYFSSLMENSPPHLFPIIFPCVLVLLK
RHYIRKSFSTIKNVINVX
[3] MQTLVCSNKELCVAYSRESKDRDEAMRKNLQKPSMFHFQPGKVYLLMFSLEIGLTFYQLNEQLVKLKYSIKRQDS
LGSPSEGYSTALHALLLISAQNVTLNLSVKEDFFFX
[4] MQRVPAADWAGARRALVLSAPPPARVARAVVAPKPPAPRAHGAHGRGGGRPPAARAPLYVHLHHSAFFPFHLL
LLLGLLEAGVIAAPRLPGARRERRRRRGGSGRGVPRGRPLDALHLDLQAPLLAEAARGGEAGAHGRAASHQPAR
VRRAGRGRPLGVPGLLRHATGVHHRGHRPARGCAAPRSALCASGAMFLSSRR
[5] MHKTMRRRYVTSSRKTLRLYFSRLVLLSSMSLVPQRPRSPRTSLSAQATRGGCPSRGRWRCSGSPPTGLGHGG
PWSSRRRLPLGWPAWXX

```

Exercise 6.4

Fill the declaration of `output$a.a.sequences` in `server.R` so that if you load a sequence file via shiny, the protein sequences appear in the main panel.

```

1 output$a.a.sequences <- renderUI({
2   inFile.2 <- input$dna.file.2
3   if (is.null(inFile.2))
4     return(NULL)
5   dna.sequences <- read.fasta(file = inFile.2$datapath,
6                               seqtype = "DNA", forceDNAtolower = T)
7   list.of.proteins <- find.mysterious.proteins(dna.sequences$seq0)
8   HTML(proteins.to.html(list.of.proteins))
9 })

```

The screenshot shows a Shiny web application titled "Tutorial Genetic Code". The interface has three tabs: "DNA sequence", "Protein sequence", and "Sequence alignment". The "DNA sequence" tab is active. It contains a "Fasta file (e.g. with the mysterious sequence)" section with a "Browse..." button and a text input "Mysterious_seq.txt", and an "Upload complete" button. Below this is a "Minimum length for a protein" slider ranging from 1 to 200, with a value of 80 selected. The "Protein sequence" tab is also visible, showing a "Compute proteins" section with a long list of protein sequences. The browser address bar shows "http://127.0.0.1:7386".

Figure 10: Proteins

7 Alignment

Exercise 7.1

a) What are the pros and cons of aligning amino acids instead of nucleotides?

The pros are

- It's easier to treat, because if we have a sequence of $3n$ nucleotides, when we translate we get a sequence of n amino-acids
- We cut off the redundancies. As we know, from the 64 codons, we have only 21 amino-acids.

And the cons are

- There is not only one sequence, but 6 sequences of amino-acids that we have to pay attention, because it's possible to slide over the sequence.
- The alphabet of amino-acids is 21, while the alphabet of nucleotides is only 4.

b) Back to your working R script (not in Shiny), align the longest mysterious translated protein from the previous part against the two reference proteins - you generated the PKD1 and PKD2 sequences in the fasta format during the Unix session (exercise 2). What was the mysterious sequence ?

```
1 requireNamespace("Biostrings", quietly = TRUE)
2 library("seqinr")
3
4 TPfolder <- "~/ECN/Biologie/13-TP-ProjR/"
5 mysterious_filename <- paste(TPfolder, "Students/Mysterious_seq.txt", sep = "")
6 codon_filename <- paste(TPfolder, "Code_shortcut/Genetic_code.txt", sep = "")
7 reference_filename <- paste(TPfolder, 'Students/PKD1-PKD2_seq.txt', sep = "")
8 proteins_filename <- paste(TPfolder, 'Students/Protein_sequences.txt', sep = "")
9 alignment_filename <- paste(TPfolder, 'Students/alignment.txt', sep = "")
10
11 ref.prot.sequences <- Biostrings::readAAStringSet(reference_filename)
12 target.prot.sequences <- Biostrings::readAAStringSet(proteins_filename)
13 myst_sequences <- read.fasta(file = mysterious_filename, seqtype = "DNA", forceDNAtolower = T)
14
15 myst_sequence <- myst_sequences$seq0
16 myst_proteins <- find.mysterious.proteins(myst_sequence)
17 biggest_myst_protein <- myst_proteins[which.max(nchar(myst_proteins))]
18 pair.alignment <- Biostrings::pairwiseAlignment(pattern = ref.prot.sequences,
19                                                  subject = biggest_myst_protein,
20                                                  substitutionMatrix = "BLOSUM62",
21                                                  type = "global")
22 if( which.max(pair.alignment@score) == 1 ){
23   print("[ PKD1 - 4303aa ] is the sequence")
24 }else{
25   print("[ PKD2 - 968aa ] is the sequence")
26 }
```

OUTPUT:

```
[1] "[ PKD2 - 968aa ] is the sequence"
```

c) Align your pair's sequence against the two reference proteins and identify which of the two reference proteins was assigned to you (use pairwiseAlignment). Write the result in a file (writePairwiseAlignments).

My sequence is the sequence 12. So, I compared them.

```
1 mysequence <- target.prot.sequences$seq12
2 pair.alignment <- Biostrings::pairwiseAlignment(pattern = ref.prot.sequences,
3                                                  subject = mysequence,
4                                                  substitutionMatrix = "BLOSUM62",
5                                                  type = "local")
6 Biostrings::writePairwiseAlignments(pair.alignment, alignment_filename)
7 if( which.max(pair.alignment@score) == 1 ){
8   print("[ PKD1 - 4303aa ] is the sequence")
9 }
```

```

9  }else{
10     print("[ PKD2 - 968aa ] is the sequence")
11  }

```

OUTPUT:

```
[1] "[ PKD1 - 4303aa ] is the sequence"
```

- d) By looking at the alignment, identify by eye the mutation/deletion/insertion between your sequence and the reference sequence for your protein by looking at the alignment result. Mutations are annotated in this way : G7602T means that a Glycine is mutated into a Threonine at the 7602th reference position.

When written in the file, we could see well the comparison of the proteins. The sequence matches with the proteins PKD1 almost entirely between the values 3641 and 4196, except the cases:

- At position **3778**, **seq12** has **R**(or Arg = CG*/AGA/AGG) while **PKD1** has **S**(or Ser = TC*/AGT/AGC)
- Between the positions **4152** and **4155**, **seq12** has **PQFR** while **PKD1** has **---**

So, there is a mutation **S3778R**.

Exercise 7.2

- a) Compute the pairwise alignment for each of the 32 AA sequences against all the others in *Protein_sequences.txt*. Build a 32 by 32 scores matrix. Use heatmap to plot the clustered scores matrix and comment the results.

```

1  len = length(target.prot.sequences)
2  matrix.scores = matrix(0, len, len)
3  for(i in 1:len){
4      for(j in 1:len){
5          pair.alignment <- Biostrings::pairwiseAlignment(pattern = target.prot.sequences[i],
6                                                         subject = target.prot.sequences[j],
7                                                         substitutionMatrix = "BLOSUM62",
8                                                         type = "global")
9          matrix.scores[i, j] = pair.alignment@score
10     }
11 }
12 heatmap(matrix.scores, symm = T, main = "Alignment of the scores of 32 proteins")

```

OUTPUT:

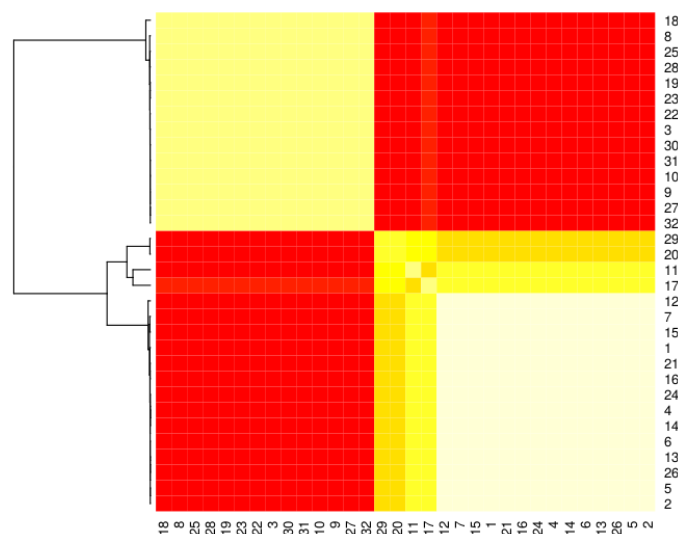


Figure 11: Heatmap

b) Do you identify clusters? How can you retrieve the sequences? What is your interpretation?

Exercise 7.3

a) Why do we use substitution matrices to deal with amino acids changes instead of constant penalties ?

Exercise 7.4

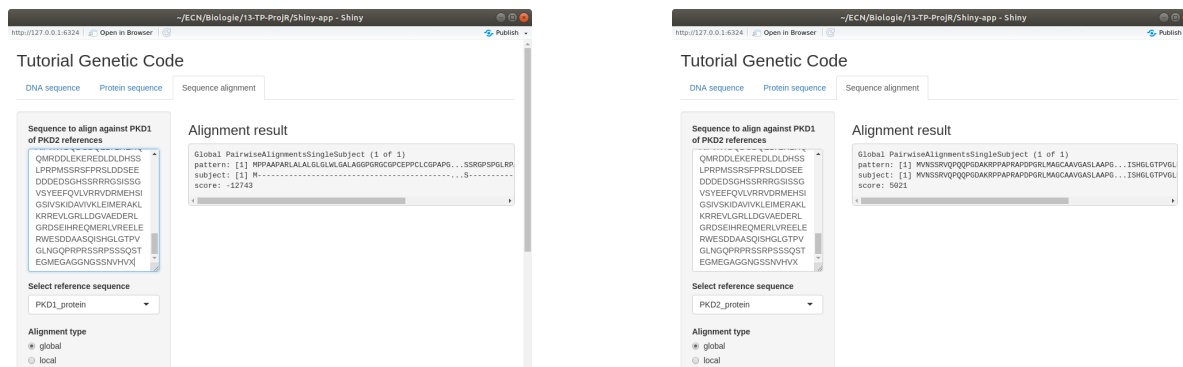
a) Fill the `output$align.out` part of the `server.R` code.

```
1 output$align.out <- renderPrint({
2   ref.sequence <- clean.sequence(ref.prot.sequences[input$select.ref.seq])
3   in.seq <- clean.sequence(input$subject.seq)
4   if(is.null(in.seq))
5     return(NULL)
6   pair.alignment <- pairwiseAlignment(attn = ref.sequence,
7                                     subject = in.seq,
8                                     substitutionMatrix = input$align.matrix,
9                                     type = input$align.type)
10  return(pair.alignment)
11 })
```

b) Add two `sliderInput` objects in the UI, one to pilot the gap opening cost (default=10) from `pairwiseAlignment`, the other one to pilot the gap extension cost (default=4). Their min and max will be 0 and 100.

```
1 tabPanel("Sequence alignment",
2   br(),
3   sidebarLayout(sidebarPanel(
4     textAreaInput(inputId = "subject.seq", value = "
5       MARVPRVRPPHGFALFLAKEEARKVKRLHGMRLSLLVYMLFLLVTLTLLASYGDASCHGHAYRLQSAIKQELH",
6       label = "Sequence to align against PKD1 of PKD2 references",
7       resize = 'both',
8       cols = 80,
9       rows = 10),
10    selectInput("select.ref.seq", "Select reference sequence", c("PKD1_protein", "
11      PKD2_protein"), "PKD1_protein", multiple = F),
12    radioButtons("align.type", "Alignment type", c("global", "local"), "global"),
13    selectInput("align.matrix", "Substitution matrix", c("BLOSUM45", "BLOSUM50", "
14      BLOSUM62", "BLOSUM80", "BLOSUM100", "PAM30", "PAM40", "PAM70", "PAM120", "PAM250")
15      , "BLOSUM62"),
16    sliderInput("slider1", "LabelSlider1", 0, 100, 10),
17    sliderInput("slider2", "labelSlider2", 0, 100, 4)
18  ),
19  mainPanel( h3("Alignment result"), verbatimTextOutput("align.out") )
20 )
```

c) Copy and paste the mysterious sequence from panel 2 and align it to the two reference sequences.



(a) Comparison with PKD1

(b) Comparison with PKD2

Figure 12: Comparison of Mysterious Sequence

d) Comment the effects of every parameter.

Exercise 7.5 (bonus)

Combine the **sapply** lines from the second chunk of code **util.R** into a single call.

Exercise 7.6 (bonus)

In the graphical rendering of the shiny application, color each aminoacid letter according to its physico-chemical properties.