

# 1 Some Molecular Biology

## 1.1 Basics

- Four bases: adenine (A), cytosine (C), guanine (G), and thymine (T)
- Two complementary base pairs: A — T and G — C
- DNA is double stranded, RNA is single stranded
  - Written from the 5' to 3' direction, e.g. 5' ACCTGAC 3'

## 1.2 Central Dogma

- Central dogma of information flow, from Crick (1958)

The central dogma states that once 'information' has passed into protein it cannot get out again. The transfer of information from nucleic acid to nucleic acid, or from nucleic acid to protein, may be possible, but transfer from protein to protein, or from protein to nucleic acid, is impossible. Information means here the precise determination of sequence, either of bases in the nucleic acid or of amino acid residues in the protein.

- DNA → RNA → Protein
  - DNA → DNA = replication
  - DNA → RNA = transcription
  - RNA → Protein = translation

## 1.3 The Genetic Code

- 20 amino acids are used to create protein sequences, as determined by Crick.
  - Assume sequences are read as blocks of nucleotides, and cannot be less than 3 letters long (since  $4$  and  $4^2 < 20$  and  $4^3 > 20$ ).
  - Therefore, if amino acids are encoded by triplets of nucleotides (called *codons*), and the code is comma-free, each reading frame is:

$$\underbrace{x_1 x_2 x_3}_{R_1} \quad \underbrace{x_4 x_5 x_6}_{R_2} \quad \underbrace{x_7 x_8 x_9}_{R_3}$$

and not

$$\underbrace{x_1 x_2 x_3 x_4}_{R_1} \quad \underbrace{x_6 x_7 x_8}_{R_2} \quad \underbrace{x_8 x_9 x_{10}}_{R_3}$$

etc...

- Each nucleotide needs a reading frame.
- All possible amino acid sequences are possible. However, AAA, TTT, GGG, and CCC are not (since in AAAAAA, for example, there is no obvious reading frame, and four different places to begin reading AAA). So we're left with  $4^3 - 4 = 60$  combinations.
- As to the others, let XYZ be a codon. To have a comma-free code, XYZXYZ must be read unambiguously. So if XYZ is a codon, YZX and ZXY cannot be. So we're left with  $60 \times \frac{1}{3} = 20$ .
- There are three stop codons, UAA, UAG, and UGA.
- There are many pairs of codons which code for the same amino acid that differ only in the third position, but relatively few which code for the same amino acid while differing in the first or second.
- Assume a sequence UUUUACUGCGGC...
  - There are three reading frames in the 5' to 3' direction and three in the opposite direction from the complementary DNA strand. So a possible of six reading frames for double-stranded DNA.
- Let  $\mathbf{N} = \{A, C, G, U\}$  be the set of nucleic acids,  $\mathbf{C} = \{(x_1 x_2 x_3) : x_i \in \mathbf{N}\}$ , and  $\mathbf{A}$  be the set of amino acids and termination codon. The genetic code is then just the map  $g : \mathbf{C} \rightarrow \mathbf{A}$  (see Table ??).

		2 <sup>nd</sup>			
		U	C	A	G
1 <sup>st</sup>					
					3 <sup>rd</sup>
U	Phe	Ser	Tyr	Cos	U
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	TC	TC	A
	Leu	Ser	TC	Trp	G
C	Leu	Pro	His	Arg	U
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	U
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	U
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

Table 1: Genetic code. Termination codons indicated by TC.

#### 1.4 Transfer RNA and Protein Sequences

- mRNA (messenger RNA) is read to make proteins.
- Amino acids are linked to tRNA (transfer RNA), which then brings them to the mRNA by interacting with codons on the mRNA.
- Since RNA is single stranded, it tends to fold back on itself to form helical regions (Figure ??).

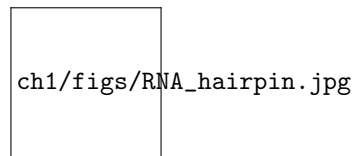


Figure 1: Example of RNA folding on itself to create a hairpin.