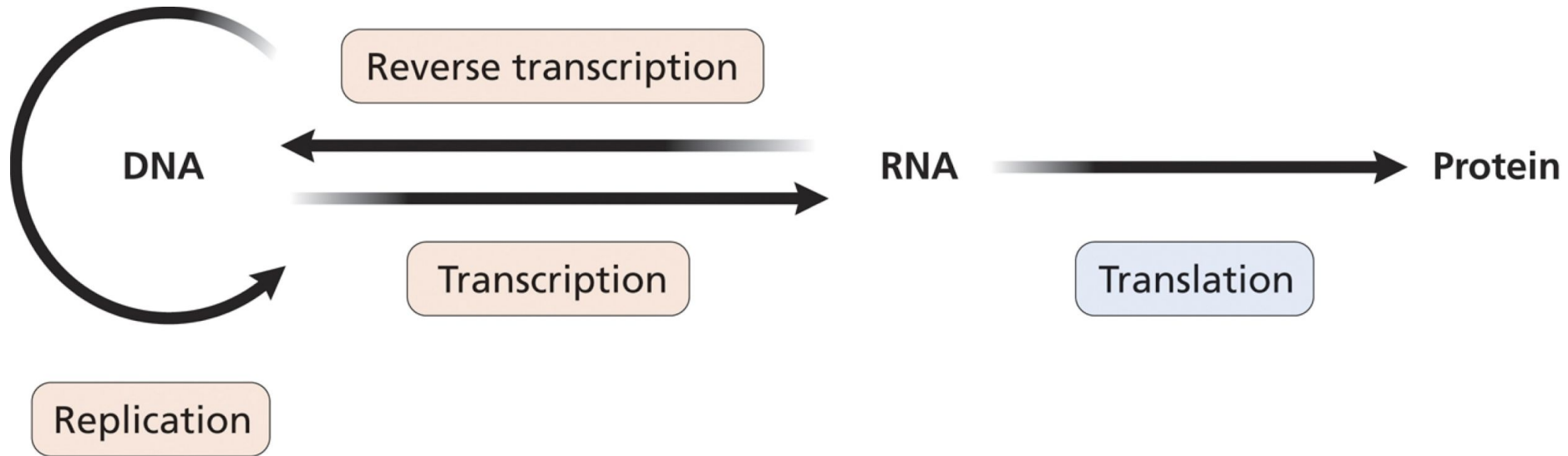


Bio 208FS: Representing the Central Dogma of Molecular Biology in Python

Paul Magwene

Central Dogma of Molecular Biology



DNA, RNA, and Proteins are biopolymers

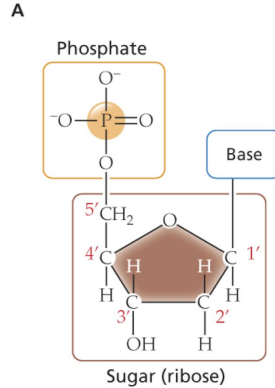
- Polymer -- a material made of many repeating subunits (monomers)
- DNA and RNA (polynucleotides)
 - monomeric units are nucleotides
- Proteins (polypeptides)
 - monomeric units are amino acids
- While the molecules adopt 3D shapes that are important for their function, their basic structure can be represented as linear sequences of their respective monomers

Using Python strings to represent biopolymers

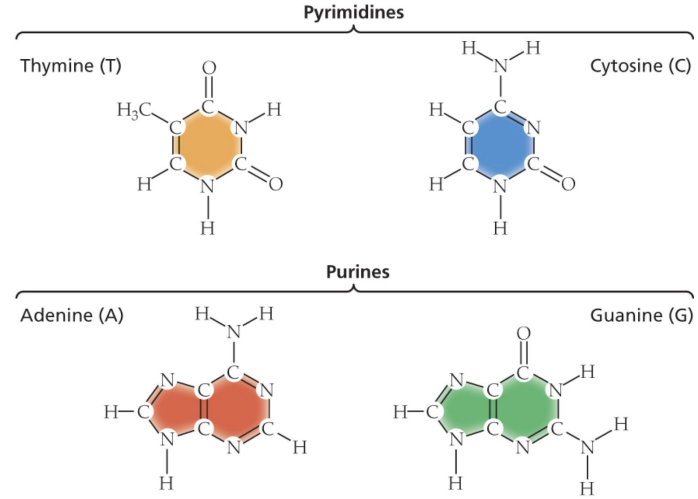
- Python strings are ordered sequences of characters (unicode in Python 3)
- Strings can be indexed and sliced like lists
- Unlike lists, strings are immutable
 - Efficiency reasons
 - Changing a string implies you're creating a new string
- Strings can be concatenated
- A variety of string methods (functions "attached to" an object) for manipulating strings are part of the Python standard implementation
 - see the [Python Documentation on Strings](#)

DNA monomers

- The size of the "alphabet" of DNA monomers (nucleotides) is 4
 - Adenine (A)
 - Guanine (G)
 - Thymine (T)
 - Cytosine (C)

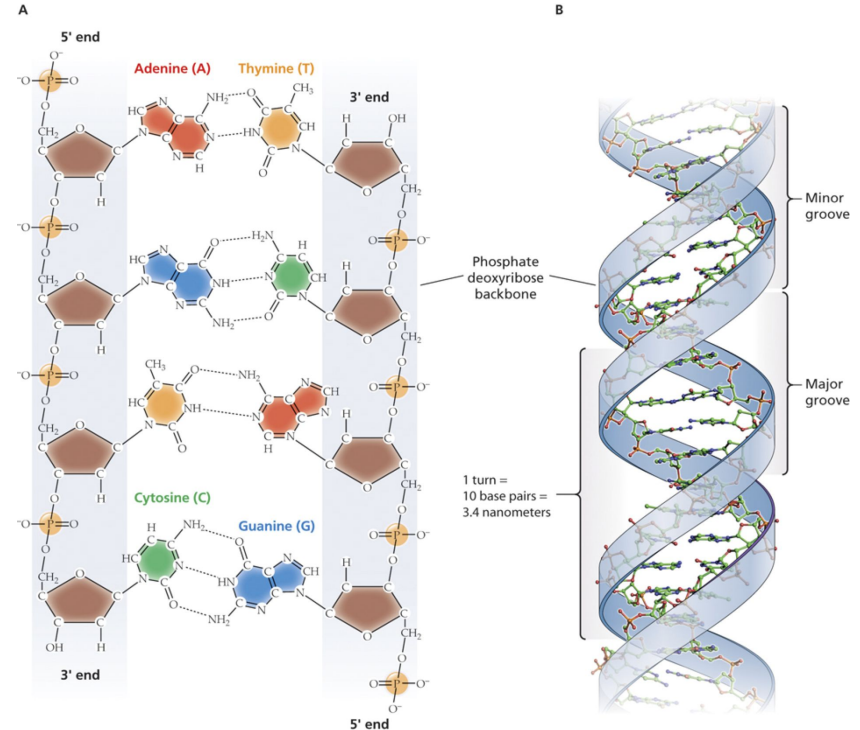


B Nitrogenous bases of DNA



DNA polymeric structure

- The backbone of the polymeric DNA structure is formed by chains of nucleotides joined by covalent bonds
 - Strings of characters from the DNA alphabet
- DNA molecules have a polarity: 5' ("5 prime") to 3' ("3 prime")
 - By convention we'll represent left (5') to right (3')
- DNA in the genome is usually in a "double stranded" state
 - The two strands are "antiparallel" -- complementary nucleotides from each strand for hydrogen bonds with each other (base pairing) -- A/T and G/C



c 5' AGTC
3' TCAG

Image from Momand and McCurdy, 2016

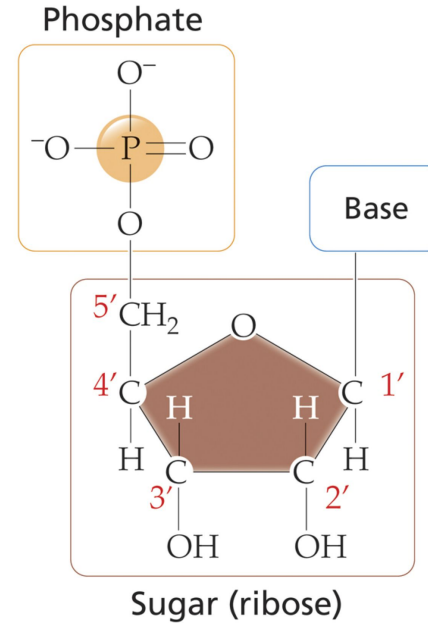
Python string representations of DNA and functions on DNA

- See hands on
- Implement functions for:
 - Reverse
 - Complement
 - Reverse complement
- Programming concepts
 - control flow statements
 - if, if-else
 - for loops
 - while
 - list comprehensions
 - Python dictionaries

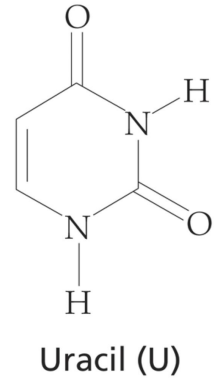
RNA

- Differences from DNA
 - Sugar is ribose instead of deoxyribose
 - Uracil is the nitrogenous base instead of thymine
- Often single stranded
 - But can form complementary interactions with DNA (or other RNAs)

A



B

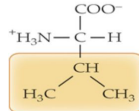


Python string representations of RNA

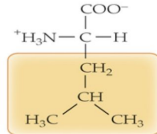
- See hands on
 - RNA as strings
 - Implement function for simulating transcription

Protein monomeric units are amino acids

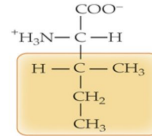
Amino acids with hydrophobic side groups



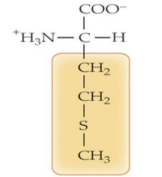
Valine
(Val) **V**



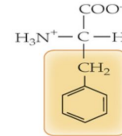
Leucine
(Leu) **L**



Isoleucine
(Ile) **I**

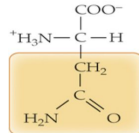


Methionine
(Met) **M**

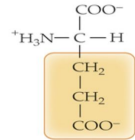


Phenylalanine
(Phe) **F**

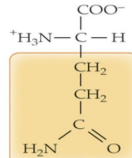
Amino acids with hydrophilic side groups



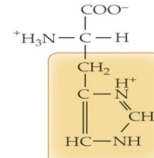
Asparagine
(Asn) **N**



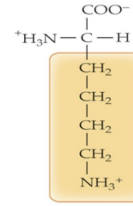
Glutamic acid
(Glu) **E**



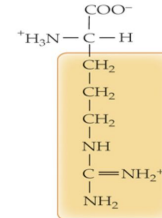
Glutamine
(Gln) **Q**



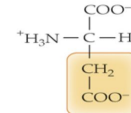
Histidine
(His) **H**



Lysine
(Lys) **K**

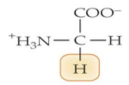


Arginine
(Arg) **R**

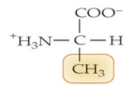


Aspartic acid
(Asp) **D**

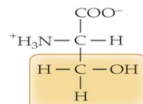
Amino acids with side groups that are in between



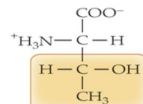
Glycine
(Gly) **G**



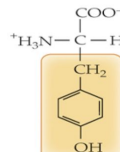
Alanine
(Ala) **A**



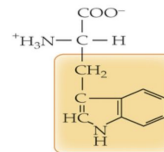
Serine
(Ser) **S**



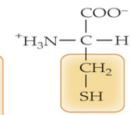
Threonine
(Thr) **T**



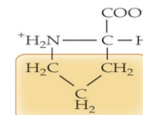
Tyrosine
(Tyr) **Y**



Tryptophan
(Trp) **W**

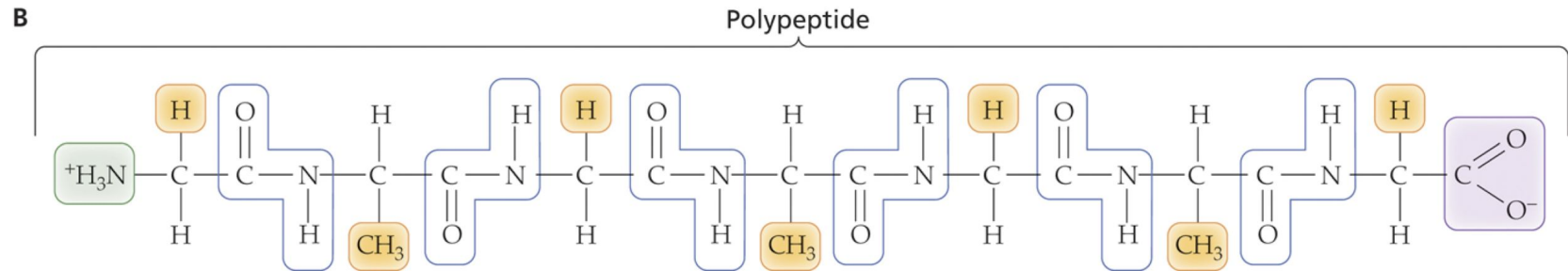
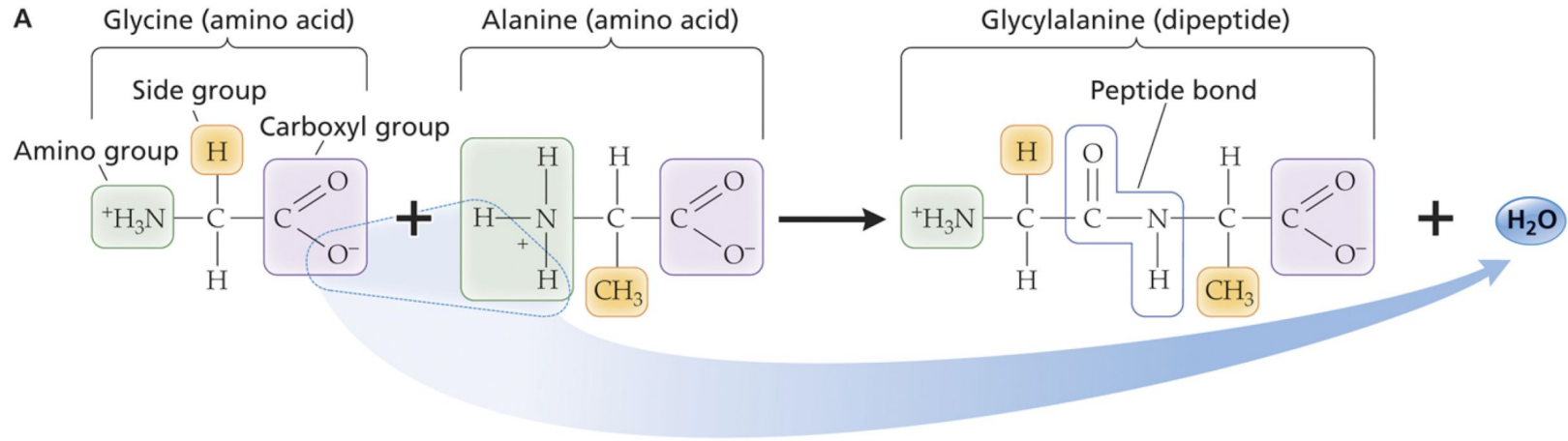


Cysteine
(Cys) **C**



Proline
(Pro) **P**

Protein polymeric structure



Python string representation of polypeptides

- Alphabet size = 20
- By convention written from left to right (N to C)

Amino Acid	3-Letter Code	1-Letter Code
Alanine	Ala	A
Cysteine	Cys	C
Aspartic acid or aspartate	Asp	D
Glutamic acid or glutamate	Glu	E
Phenylalanine	Phe	F
Glycine	Gly	G
Histidine	His	H
Isoleucine	Ile	I
Lysine	Lys	K
Leucine	Leu	L
Methionine	Met	M
Asparagine	Asn	N
Proline	Pro	P
Glutamine	Gln	Q
Arginine	Arg	R
Serine	Ser	S
Threonine	Thr	T
Valine	Val	V
Tryptophan	Trp	W
Tyrosine	Tyr	Y

Standard Genetic Code

		Second position								
		T		C		A		G		
		Code	Amino acid	Code	Amino acid	Code	Amino acid	Code	Amino acid	
First position	T	T T T	phe	T C T	ser	T A T	tyr	T G T	cys	T
		T T C		T C C		T A C		T G C		C
		T T A	leu	T C A		T A A	STOP	T G A	STOP	A
		T T G		T C G		T A G	STOP	T G G	trp	G
	C	C T T	leu	C C T	pro	C A T	his	C G T	arg	T
		C T C		C C C		C A C		C G C		C
		C T A		C C A		C A A	gln	C G A		A
		C T G		C C G		C A G		C G G		G
	A	A T T	ile	A C T	thr	A A T	asn	A G T	ser	T
		A T C		A C C		A A C		A G C		C
		A T A		A C A		A A A	lys	A G A	arg	A
		A T G	met	A C G		A A G		A G G		G
	G	G T T	val	G C T	ala	G A T	asp	G G T	gly	T
		G T C		G C C		G A C		G G C		C
		G T A		G C A		G A A	glu	G G A		A
		G T G		G C G		G A G		G G G		G

Third position									
----------------	--	--	--	--	--	--	--	--	--

Image from Momand and McCurdy, 2016

Implement a translation function

- Input: string representing a DNA or RNA sequence
- Output: string representing a Protein sequence