

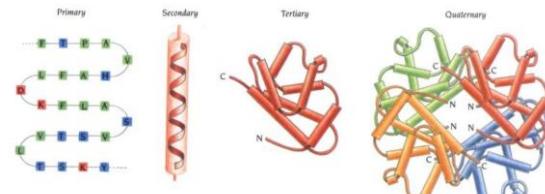
Protein Structure Analysis

<http://binf.gmu.edu/vaisman/binf731/>

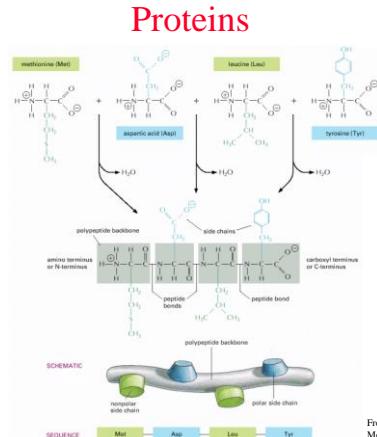
Iosif Vaisman

2020

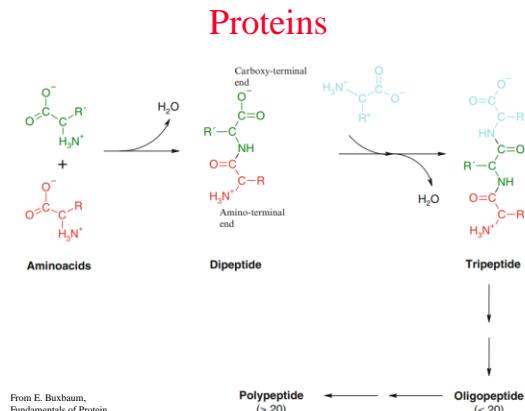
1



2

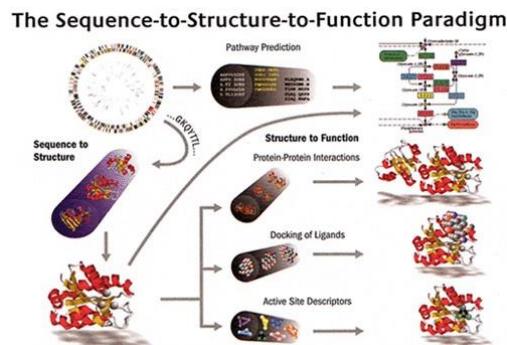


3



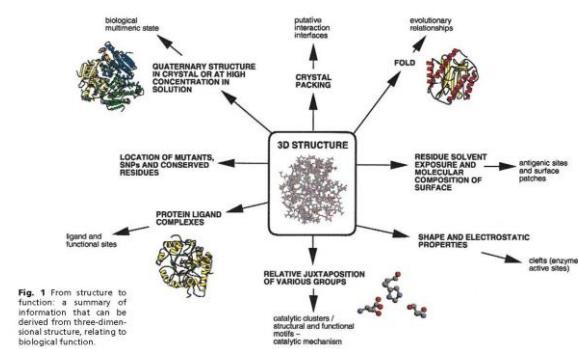
4

Proteins: Sequence, Structure, Function



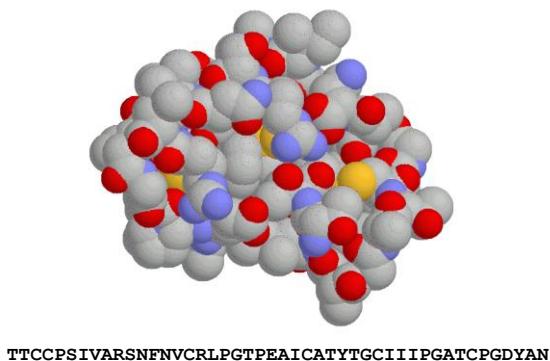
5

Proteins: Sequence, Structure, Function



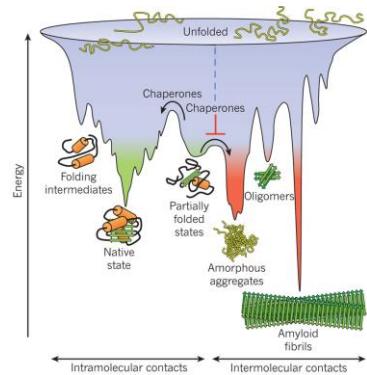
6

Proteins



7

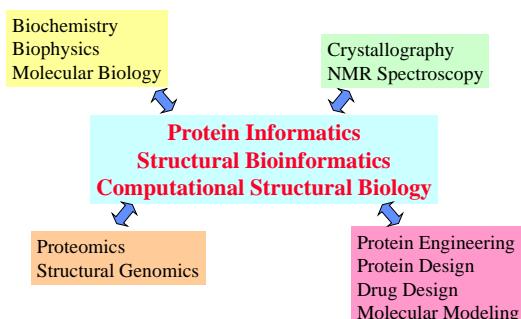
Proteins



Hartl F.U. et al., Nature, 2011

8

Protein Science



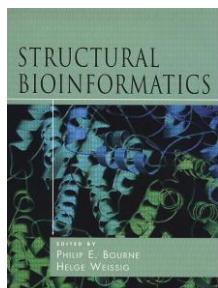
9

Protein Structure and...

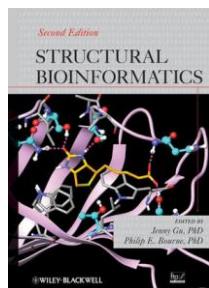
Business
Law
Ethics
Medicine
...

10

Recommended book



Philip Bourne, Helge Weissig (Eds)
Structural bioinformatics
Hoboken, N.J. : Wiley-Liss, 2003.



Jenny Gu, Philip Bourne (Eds)
Structural bioinformatics
Hoboken, N.J. : Wiley-Liss, 2009.

Protein Informatics

SEQUENCE
↓
STRUCTURE
↓
DYNAMICS
↓
FUNCTION

11

12

Information

General

knowledge or intelligence communicated, received or gained

Information theory

indication of the number of possible choices

```
Th_ qui_k br_wn _ox ju_ps ov__ th_ laz_ d_g
Ae_h uz_ ko_ wm so_g oqr_it ypu_vn tr_e oj_
```

Information

```
Th_ qui_k br_wn _ox ju_ps ov__ th_ laz_ d_g
Ae_h uz_ ko_ wm so_g oqr_it ypu_vn tr_e oj_
```

```
The quick brown fox jumps over the lazy dog
Aedh uzh kox wm sobg oqrfit ypulvn tree ojc
```

13

14

Information and uncertainty

Information is a decrease in uncertainty

$$\log_2(M) = -\log_2(M^{-1}) = -\log_2(P)$$

Shannon's formula for uncertainty

$$H = - \sum_{i=1}^M P_i \log_2 P_i$$

only information essential to understand what is transmitted

Communication

Fundamental problem of communication:

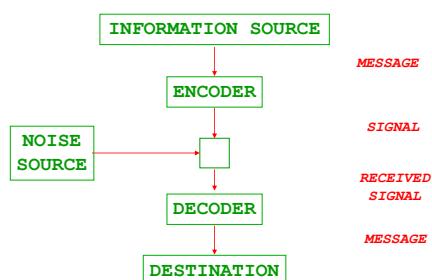
reproducing at one point either exactly or approximately a message selected at another point

The Mathematical Theory of Communication
Claude Shannon and Warren Weaver

15

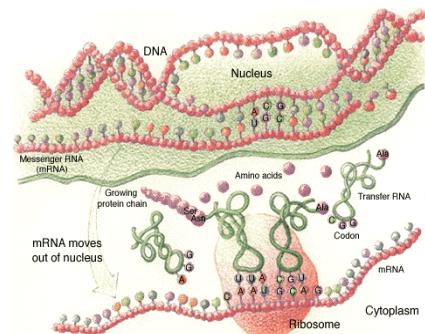
16

Communication system



Adopted from C.E. Shannon,
The Mathematical Theory of Communication, 1949

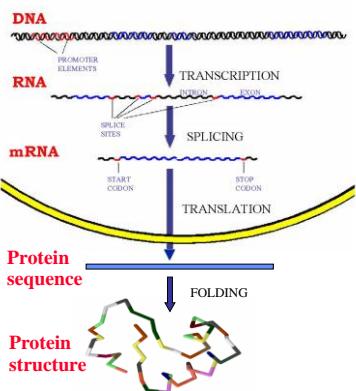
Cell Informatics



17

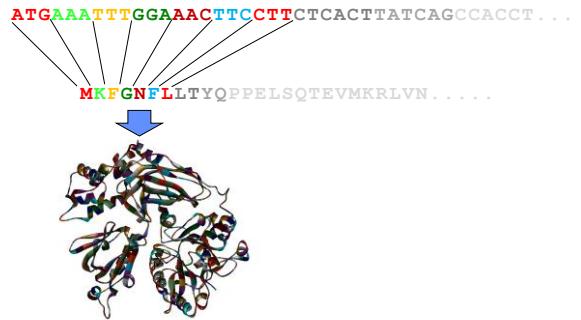
18

Cell Informatics



19

DNA Sequence – Protein Sequence – Protein Structure



20

Communication system duality

"This duality can be pursued further and is related to the duality between past and future and the notions of control and knowledge. Thus we may have knowledge of the past but cannot control it; we may control the future but have no knowledge of it."

C. E. Shannon (1959)

21

Error correcting codes

a	b	c	d	e
a				
b	X			X
c		X		
d			X	
e				X

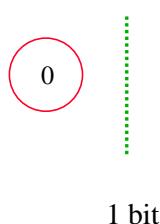
Code words ac, ba, be, db, ed in the permutation space of $[a..e] \times [a..e]$

Hamming metric

The sum of bit changes necessary to move from one point in the permutation space to another point in the permutation space

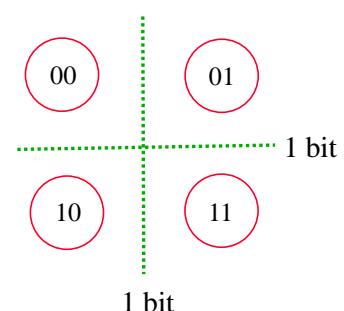
0000 and 0111 are separated by Hamming distance of 3:
0000 - 0001 - 0011 - 0111

Information Theory



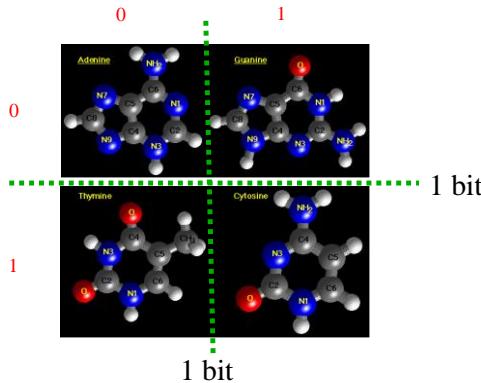
23

Information Theory



24

Nucleotide permutation space



25

Standard genetic code

TTT	F Phe	TCT	S Ser	TAT	Y Tyr	TGT	C Cys
TTC	F Phe	TCC	S Ser	TAC	Y Tyr	TGC	C Cys
TTA	L Leu	TCA	S Ser	TAA	* Ter	TGA	* Ter(Sec)
TTG	L <u>Leu</u>	TCG	S Ser	TAG	* Ter(Pyl)	TGG	W Trp
<hr/>							
CTT	L Leu	CCT	P Pro	CAT	H His	CGT	R Arg
CTC	L Leu	CCC	P Pro	CAC	H His	CGC	R Arg
CTA	L Leu	CCA	P Pro	CAA	Q Gln	CGA	R Arg
CTG	L <u>Leu</u>	CCG	P Pro	CAG	Q Gln	CGG	R Arg
<hr/>							
ATT	I Ile	ACT	T Thr	AAT	N Asn	AGT	S Ser
ATC	I Ile	ACC	T Thr	AAC	N Asn	AGC	S Ser
ATA	I Ile	ACA	T Thr	AAA	K Lys	AGA	R Arg
ATG	M <u>Met</u>	ACG	T Thr	AAG	K Lys	AGG	R Arg
<hr/>							
GTT	V Val	GCT	A Ala	GAT	D Asp	GGT	G Gly
GTC	V Val	GCC	A Ala	GAC	D Asp	GGC	G Gly
GTA	V Val	GCA	A Ala	GAA	E Glu	GGA	G Gly
GTG	V Val	GCG	A Ala	GAG	E Glu	GGG	G Gly

26

Standard genetic code

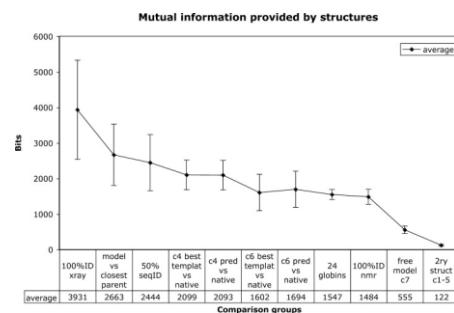
Frameshift Errors

ATGAAATTGGAACTCCTCTCACTTATCAGCCACCTGAGCTATCTCAGACCGAAGTGATGAAGCGATTGGTTAATCT

5'3'Frame1 MKFGNFLLTQPPELSQTEVMKRLVNT
5'3'Frame2 -NLETSFSLISHLSYLRPK--SDWLI
5'3'Frame3 EIWKLPSPHSAT-AISDRSDEAIG-S
3'5'Frame1 RLTNRFITSV-DSSGG--VRRKFPNFT
3'5'Frame2 D-PIASSLRSEIAQVADK-EGSFQIS
3'5'Frame3 INOSLHHFGLR-LRWLISEKEVSKFH

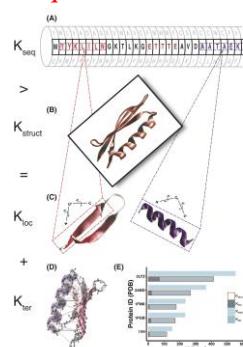
27

Information theory in protein structure evaluation



B.C. 1-3000

Information content in protein sequences and structures



A. Possenti et al. 201

Amino acid naming conventions

- Amino acids with a unique first letter: **Cys**, **His**, **Ile**, **Met**, **Ser**, **Val**
 - Where several amino acids start with the same letter, common amino acids are given preference: **Ala**, **Gly**, **Leu**, **Pro**, **Thr**
 - Letters other than the first letter are used for **Asn** (asparagi**N**), **Arg** (**a**rginine), **Tyr** (**t**yrosine)
 - Similar sounding names: **Asp** (asp**a**r**Dic** acid), **Glu** (glu**e**mate), **Gln** (**Q**amine), **Phe** (**F**enylalanine)
 - The remaining amino acids have letters that do not occur in their name: **Lys** (**K** close to L), **Trp** (**W** reminds of double ring), **Sec** (**U**), **Pyl** (**O**)
 - **X** is used as placeholder, meaning “any amino acid”. **B** is used for “Asp or Asn”, **Z** for “Gln or Glu”, **J** for “Ile or Leu”. The **-** is used to denote gaps in a protein sequence, e.g., in sequence alignments. **h** is used to denote hydrophobic amino acids (do not confuse with H for His!)

E Ruxbaum 2015

29

30