

## Single-letter and three-letter symbols for amino acids

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

## The genetic code Remember U in mRNA and T in DNA

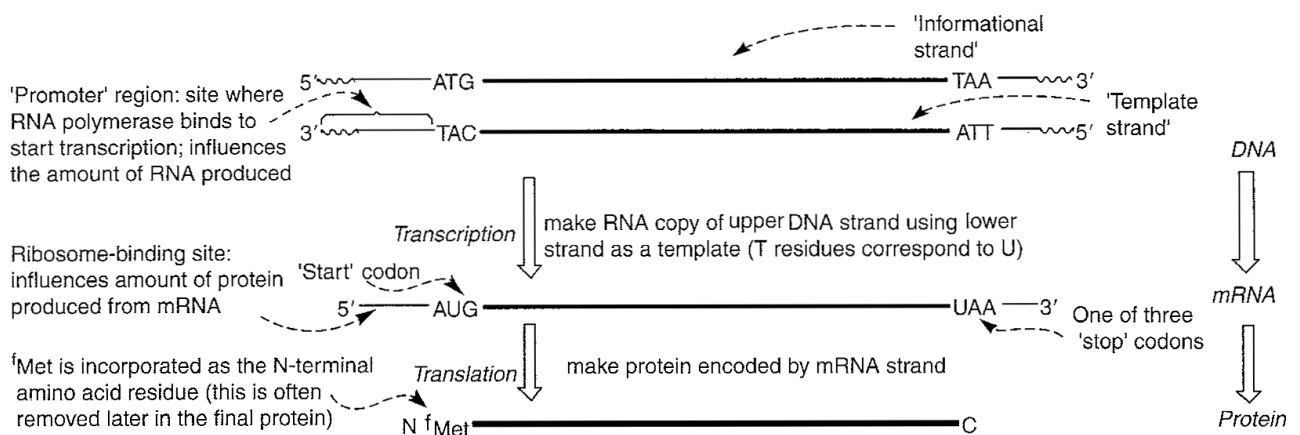
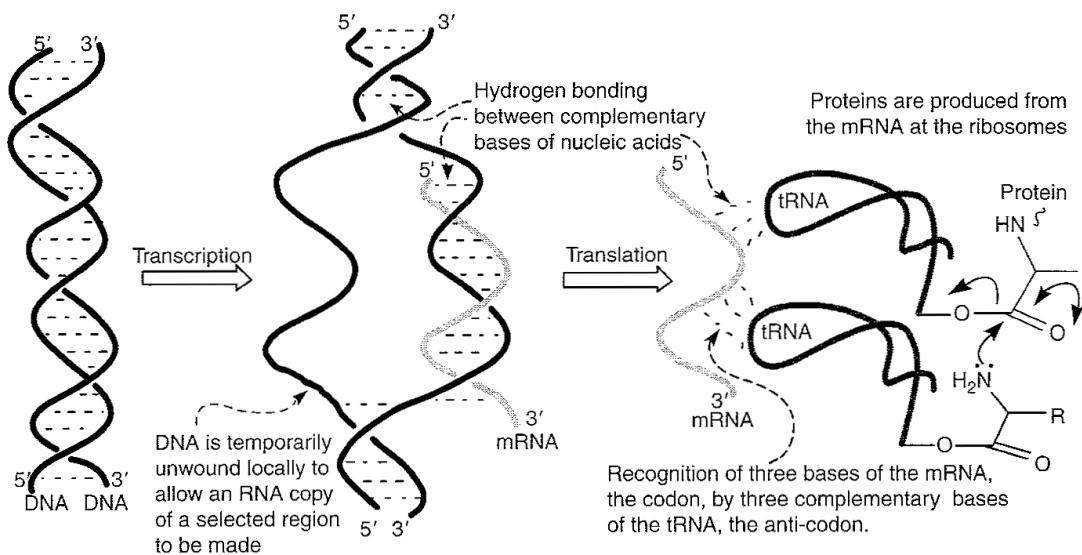
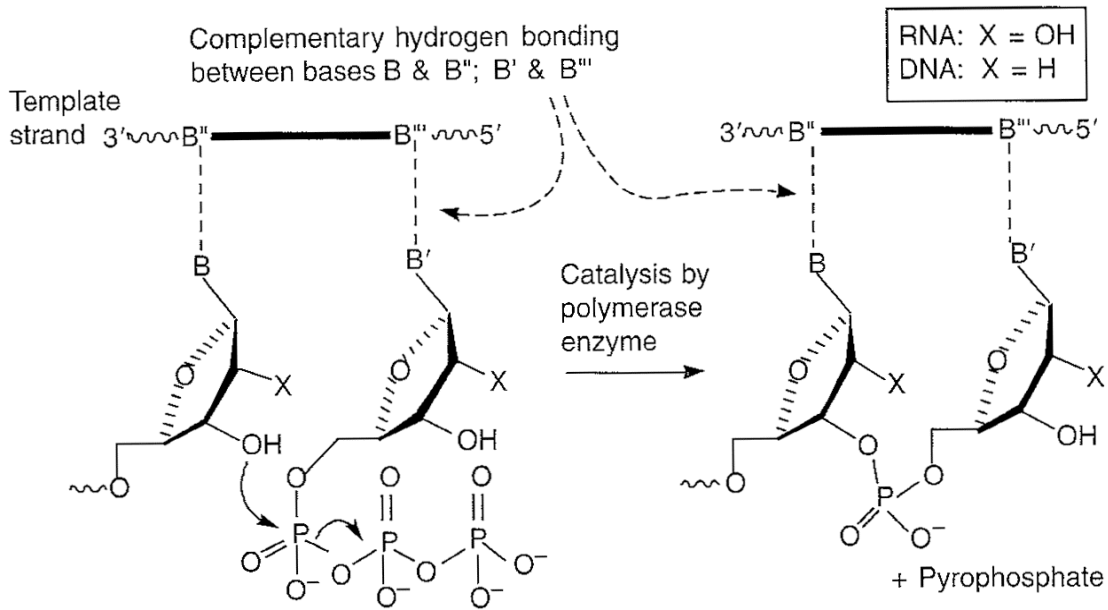
5' base	Middle base				3' base
	U	C	A	G	
<b>U</b>	UUU Phe	UCU Ser	UAU Tyr	UGU Cys	<b>U</b>
	UUC Phe	UCC Ser	UAC Tyr	UGC Cys	<b>C</b>
	UUA Leu	UCA Ser	UAA Stop*	UGA Stop*	<b>A</b>
	UUG Leu	UCG Ser	UAG Stop*	UGG Trp	<b>G</b>
<b>C</b>	CUU Leu	CCU Pro	CAU His	CGU Arg	<b>U</b>
	CUC Leu	CCC Pro	CAC His	CGC Arg	<b>C</b>
	CUA Leu	CCA Pro	CAA Gln	CGA Arg	<b>A</b>
	CUG Leu	CCG Pro	CAG Gln	CGG Arg	<b>G</b>
<b>A</b>	AUU Ile	ACU Thr	AAU Asn	AGU Ser	<b>U</b>
	AUC Ile	ACC Thr	AAC Asn	AGC Ser	<b>C</b>
	AUA Ile	ACA Thr	AAA Lys	AGA Arg	<b>A</b>
	AUG Met <sup>†</sup>	ACG Thr	AAG Lys	AGG Arg	<b>G</b>
<b>G</b>	GUU Val	GCU Ala	GAU Asp	GGU Gly	<b>U</b>
	GUC Val	GCC Ala	GAC Asp	GGC Gly	<b>C</b>
	GUA Val	GCA Ala	GAA Glu	GGA Gly	<b>A</b>
	GUG Val	GCG Ala	GAG Glu	GGG Gly	<b>G</b>

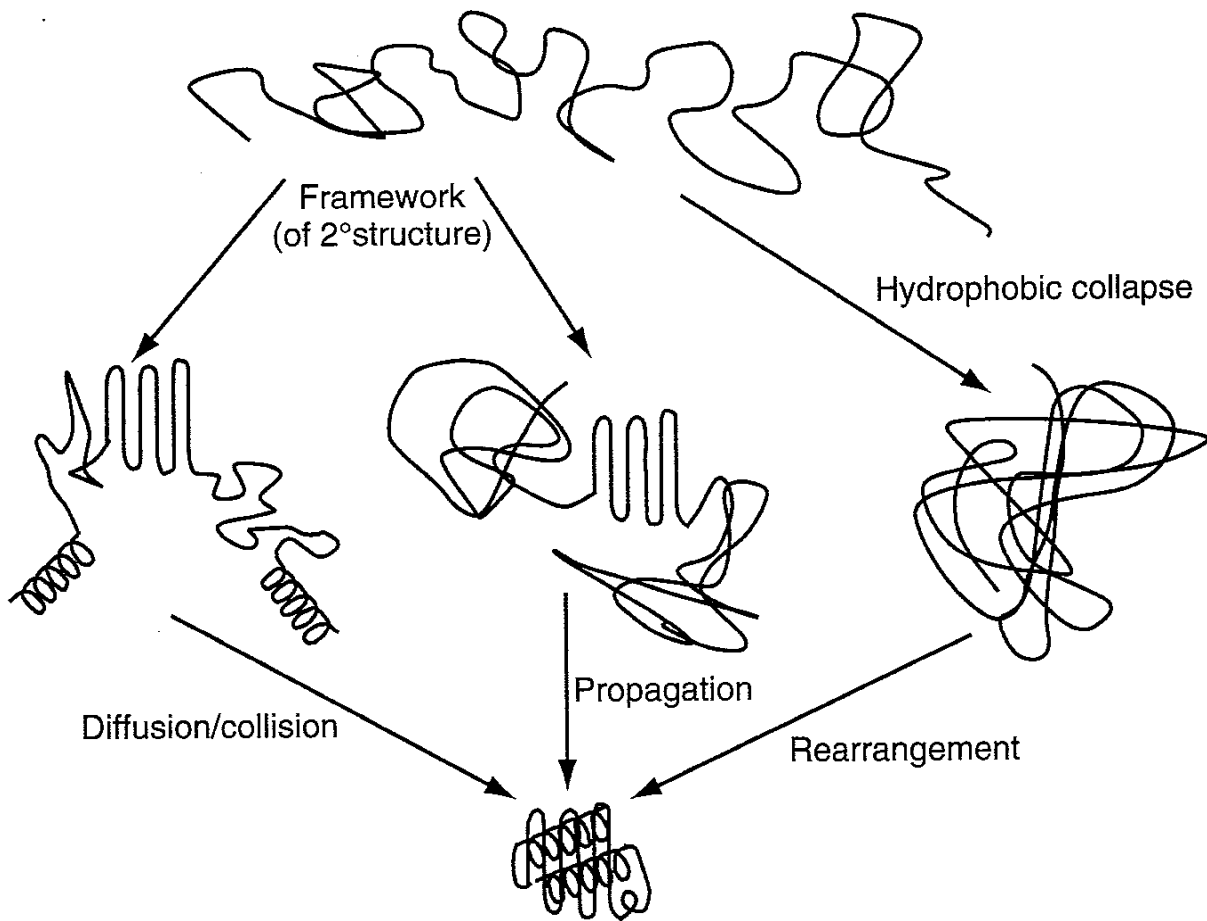
\*Stop codons have no amino acids assigned to them.

<sup>†</sup>The AUG codon is the usual initiation codon as well as that for methionine residues elsewhere. The code is almost universal but differences have been found in mitochondrial DNA from some organisms

### Genome sizes

Organism	Number of Base pairs	Number of Genes	Comment
$\phi$ X-174	5386	10	virus infecting <i>E. coli</i>
Human mitochondrion	16 569	37	subcellular organelle
Epstein-Barr virus (EBV)	172 282	80	cause of mononucleosis
<i>Mycoplasma pneumoniae</i>	816 394	680	cause of cyclic pneumonia epidemics
<i>Rickettsia prowazekii</i>	1 111 523	878	bacterium, cause of epidemic typhus
<i>Treponema pallidum</i>	1 138 011	1 039	bacterium, cause of syphilis
<i>Borrelia burgdorferi</i>	1 471 725	1 738	bacterium, cause of Lyme disease
<i>Aquifex aeolicus</i>	1 551 335	1 749	bacterium from hot spring
<i>Thermoplasma acidophilum</i>	1 564 905	1 509	archaeal prokaryote, lacks cell wall
<i>Campylobacter jejuni</i>	1 641 481	1 708	frequent cause of food poisoning
<i>Helicobacter pylori</i>	1 667 867	1 589	chief cause of stomach ulcers
<i>Methanococcus jannaschii</i>	1 664 970	1 783	archaeal prokaryote, thermophile
<i>Hemophilus influenzae</i>	1 830 138	1 738	bacterium, cause of middle ear infections
<i>Thermotoga maritima</i>	1 860 725	1 879	marine bacterium
<i>Archaeoglobus fulgidus</i>	2 178 400	2 437	another archaeon
<i>Deinococcus radiodurans</i>	3 284 156	3 187	radiation-resistant bacterium
<i>Synechocystis</i>	3 573 470	4 003	cyanobacterium, 'blue-green alga'
<i>Vibrio cholerae</i>	4 033 460	3 890	cause of cholera
<i>Mycobacterium tuberculosis</i>	4 411 529	4 275	cause of tuberculosis
<i>Bacillus subtilis</i>	4 214 814	4 779	popular in molecular biology
<i>Escherichia coli</i>	4 639 221	4 406	molecular biologists' all-time favourite
<i>Pseudomonas aeruginosa</i>	6 264 403	5 570	largest prokaryote sequenced as yet
<i>Saccharomyces cerevisiae</i>	$12.1 \times 10^6$	5 885	yeast, first eukaryotic genome sequenced
<i>Caenorhabditis elegans</i>	$95.5 \times 10^6$	19 099	the worm
<i>Arabidopsis thaliana</i>	$1.17 \times 10^8$	25 498	flowering plant (angiosperm)
<i>Drosophila melanogaster</i>	$1.8 \times 10^8$	13 601	the fruit fly
<i>Fugu rubripes</i>	$3.9 \times 10^8$	30 000	puffer fish (fugu fish)
Human	$3.2 \times 10^9$	34 000?	
Wheat	$16 \times 10^9$	30 000	
Salamander	$10^{11}$	?	
<i>Psilotum nudum</i>	$10^{11}$	?	whisk fern – a simple plant





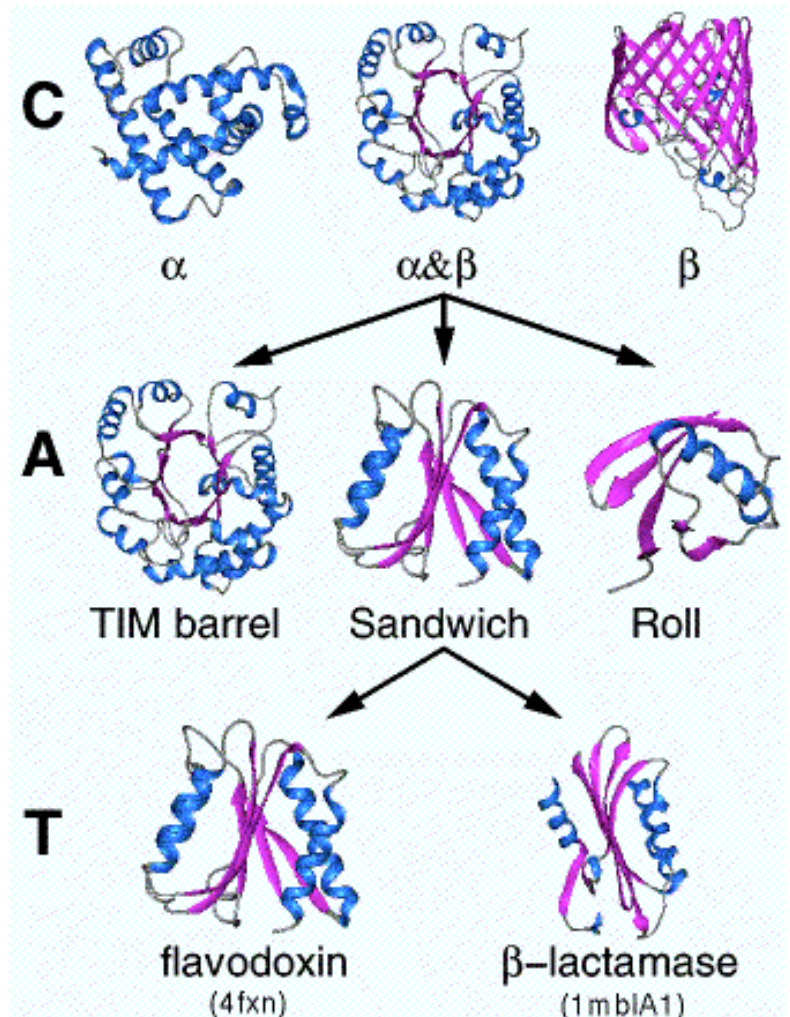
A database for 3D shapes of protein domains

**C** Class – Secondary Structure

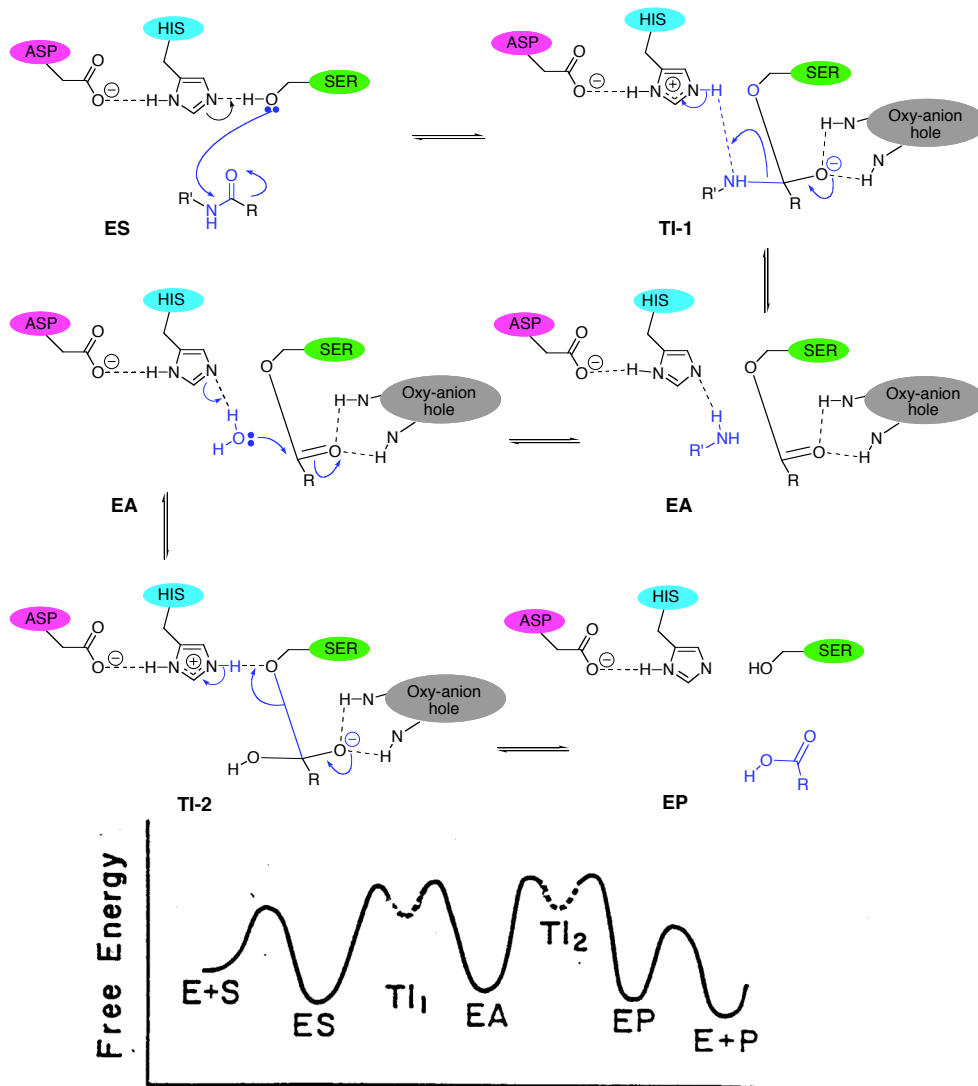
**A** Architecture – Arrangements of 2ndary Structure but ignores connections

**T** Topology – Clustering of A

**H** Homologous Superfamily – Combines Sequence and Structure information to create groups.



## Serine Hydrolases



*A representation of the expected free energy diagram for serine proteinase catalysis. From evolutionary principles the free energies of all the transition states are expected to be similar, and the energies of all the intermediates are anticipated to be similar*

### Additional Reading:

- Ashburner *et al*; Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nature Genetics* **2000**, 25, 25-29
- Ursing *et al*; EXProt: a database for proteins with experimentally verified function. *Nucleic Acids Research* **2002**, 30, 50-51.
- Orengo *et al* CATH- A Hierarchic Classification of Protein Domain Structures. *Structure* **1997**, 5, 1093-1108.
- Pearl *et al* Assigning genomic sequences to CATH *Nucleic Acids Research*. **2000**, 28, 277-282