

Motifs in Protein Sequences

Motifs are combinations of secondary structures in proteins with a specific **structure** and a specific **function**. They are also called **super-secondary structures**.

Examples: Helix-Turn-Helix, Zinc-finger, Homeobox domain, Hairpin-beta motif, Calcium-binding motif, Beta-alpha-beta motif, Coiled-coil motifs.

Several motifs may combine to form **domains**.

- Serine proteinase domain, Kringle domain, calcium-binding domain, homeobox domain.

Motif Detection Problem

Input:

Set, S , of known (**aligned**) examples of a motif M ,
A new protein sequence, P .

Output:

Does P have a copy of the motif M ?

Example: Zinc Finger Motif

...**Y**K**C**GL**C**ERS**F**VEKS**L**SR**H**ORV**H**KN...
 3 6 19 23

Input:

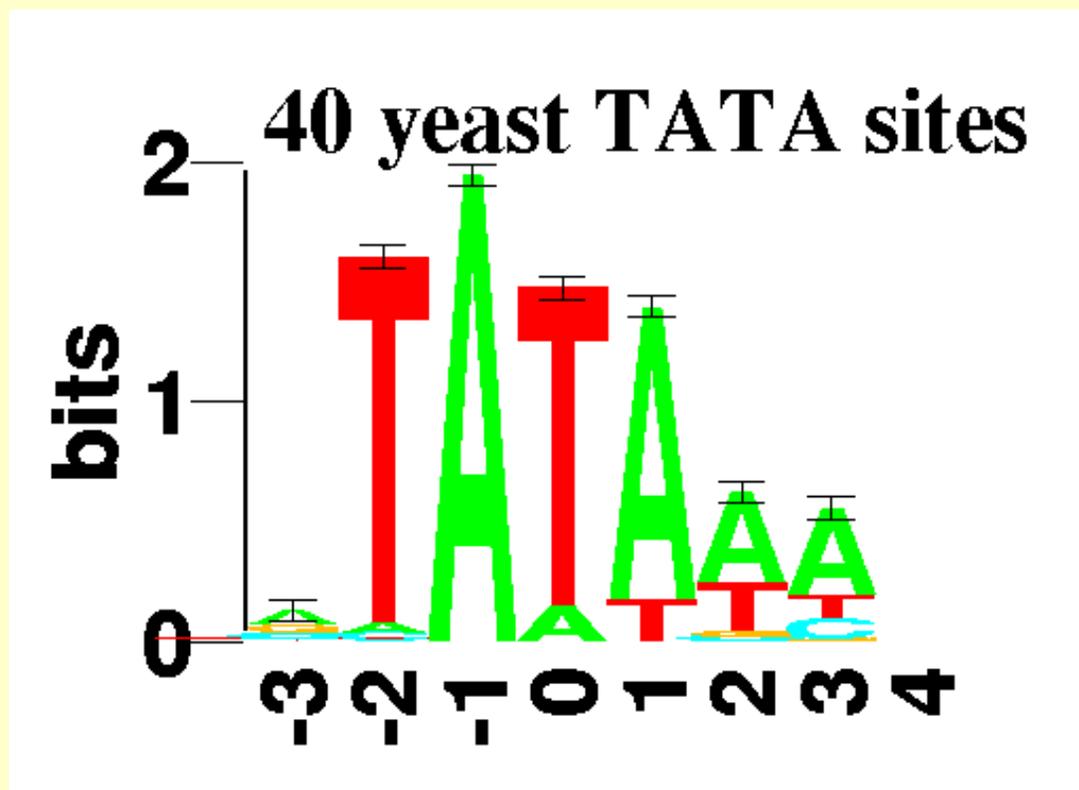
Database, D , of known protein sequences,
A new protein sequence, P .

Output:

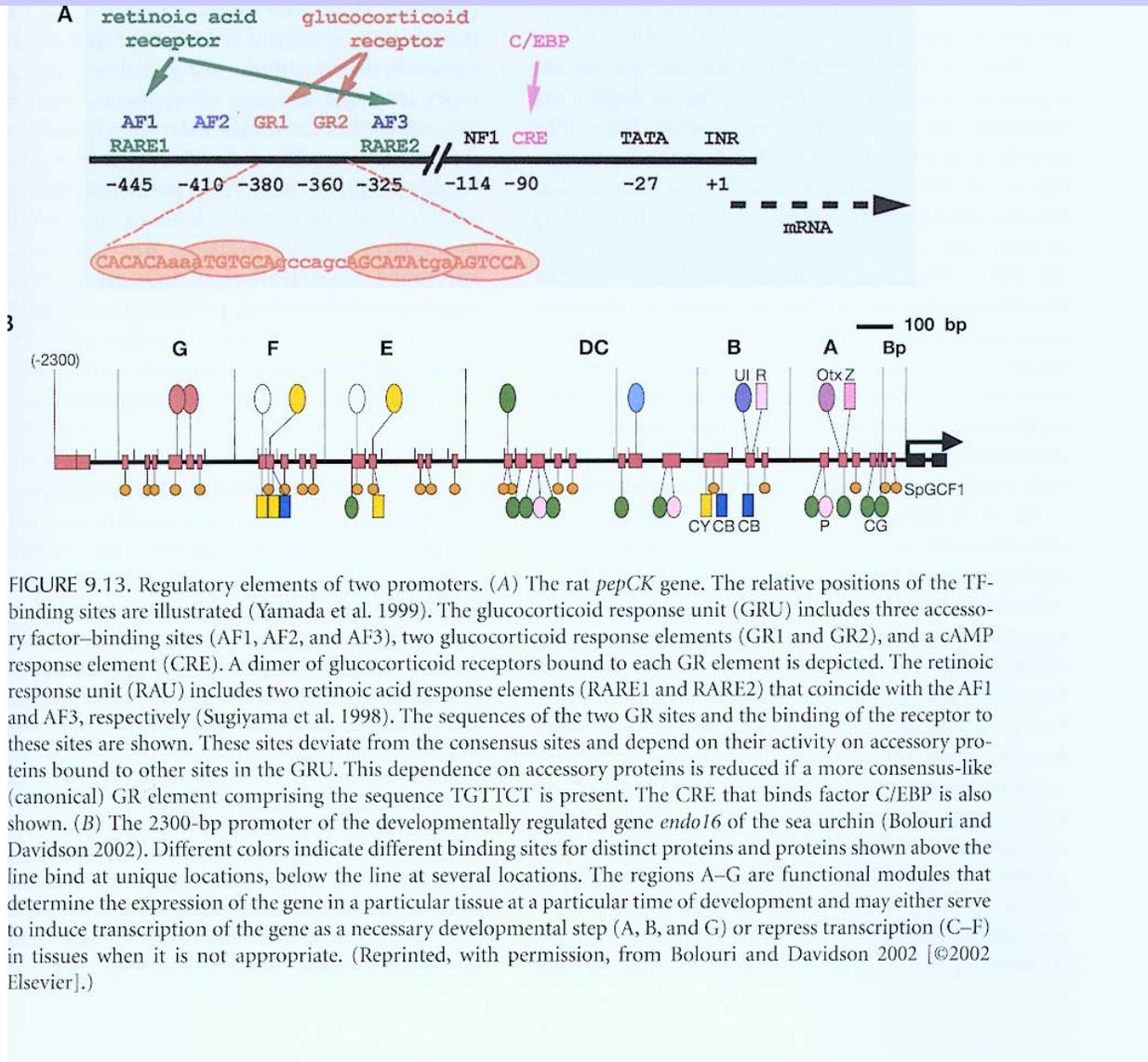
What interesting patterns from D
are present in P ?

Motifs in DNA Sequences

- Given a collection of DNA sequences of promoter regions, locate the transcription factor binding sites (also called regulatory elements)
 - Example:



Motifs in DNA Sequences



Motif Detection (TFBMs)

- See evaluation by Tompa et al.
 - [bio.cs.washington.edu/assessment]
- **Gibbs Sampling Methods:** AlignACE, GLAM, SeSiMCMC, MotifSampler
- **Weight Matrix Methods:** ANN-Spec, Consensus,
- **EM:** Improbizer, MEME
- **Combinatorial & Misc.:** MITRA, oligo/dyad, QuickScore, Weeder, YMF

Motif Detection

- Profile Method
 - If many examples of the motif are known, then
 - **Training**: build a **Profile** and compute a **threshold**
 - **Testing**: **score** against profile
- Gibbs Sampling
- Expectation Method
- HMM
- Combinatorial Pattern Discovery Methods

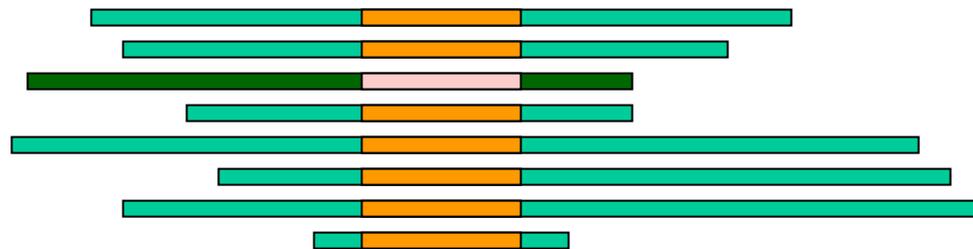
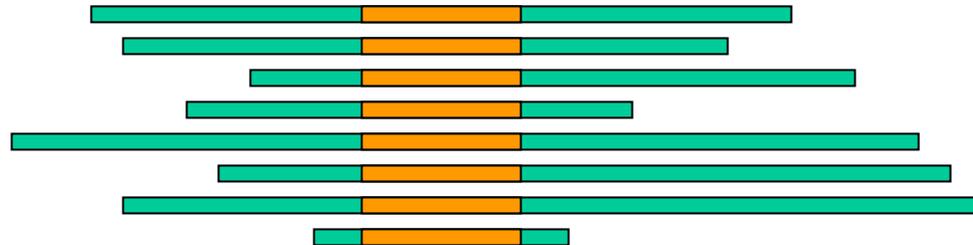
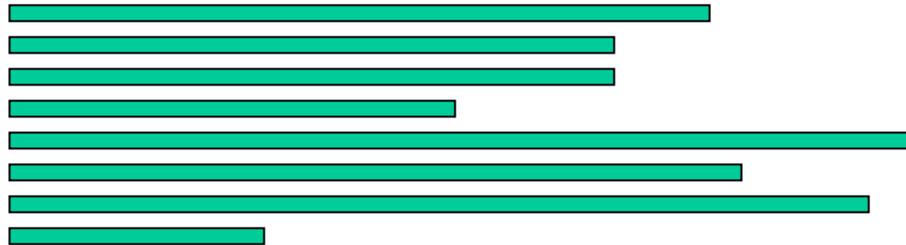
How to evaluate these methods?

- Calculate TP, FP, TN, FN
- Compute **sensitivity** fraction of known sites predicted, **specificity**, and more.
 - **Sensitivity** = $TP / (TP + FN)$
 - **Specificity** = $TN / (TN + FP)$
 - Positive Predictive Value = $TP / (TP + FP)$
 - Performance Coefficient = $TP / (TP + FN + FP)$
 - Correlation Coefficient =

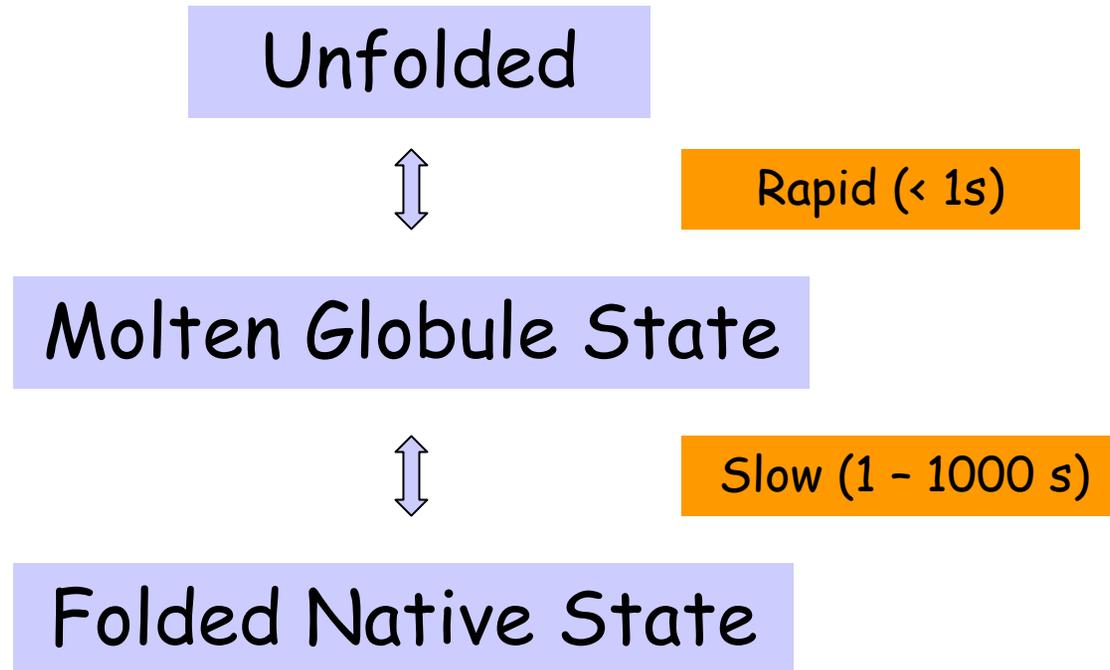
Motif Detection (TFBMs)

- See evaluation by Tompa et al.
 - [bio.cs.washington.edu/assessment]
- **Gibbs Sampling Methods:** AlignACE, GLAM, SeSiMCMC, MotifSampler
- **Weight Matrix Methods:** ANN-Spec, Consensus,
- **EM:** Improbizer, MEME
- **Combinatorial & Misc.:** MITRA, oligo/dyad, QuickScore, Weeder, YMF

Gibbs Sampling for Motif Detection



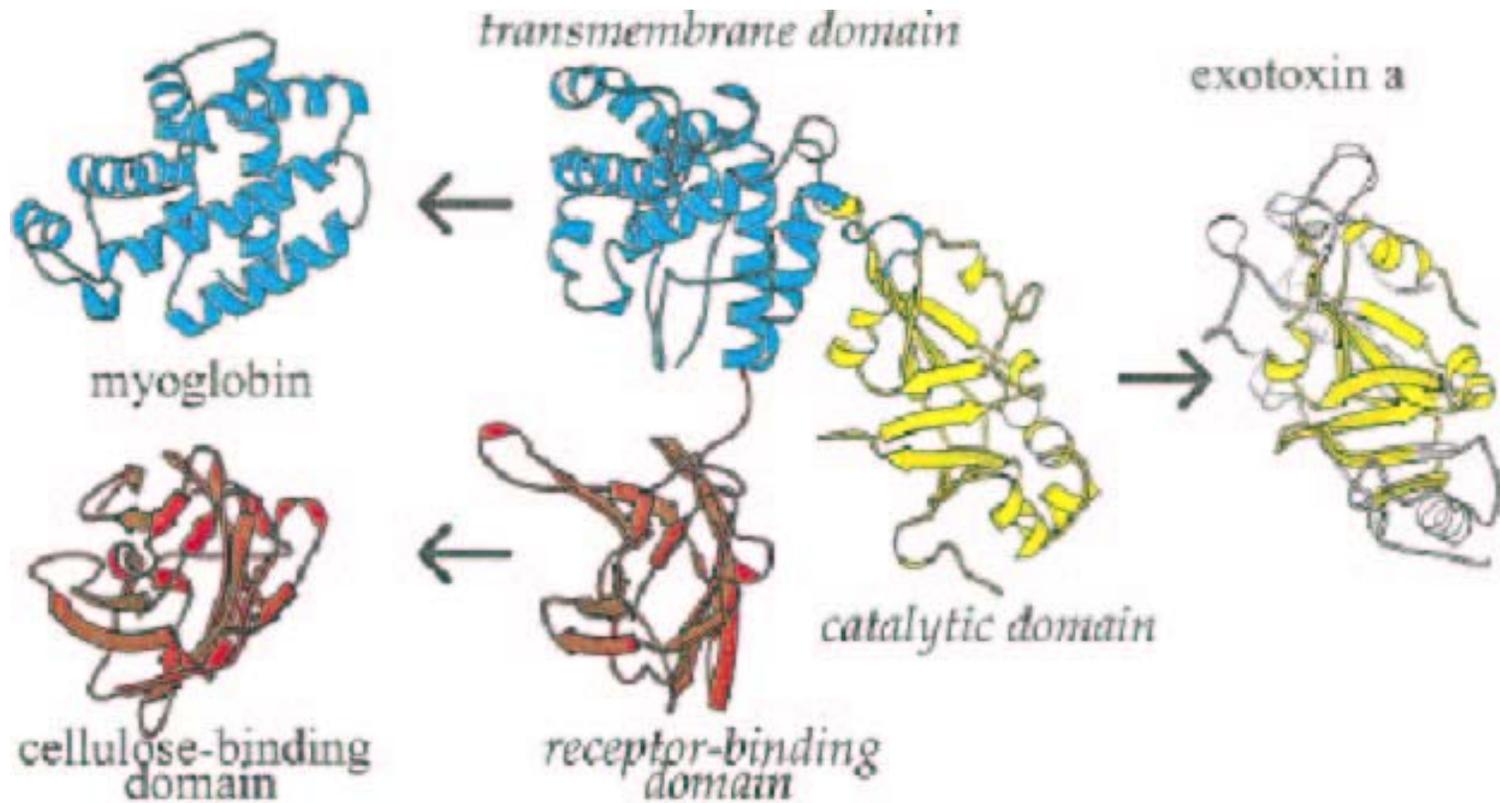
Protein Folding



- How to find minimum energy configuration?

Modular Nature of Protein Structures

Example: Diphtheria Toxin



Protein Structures

- Most proteins have a **hydrophobic core**.
- Within the core, specific **interactions** take place between amino acid side chains.
- Can an amino acid be replaced by some other amino acid?
 - Limited by space and available contacts with nearby amino acids
- Outside the core, proteins are composed of loops and structural elements in contact with water, solvent, other proteins and other structures.

Viewing Protein Structures

- SPDBV
- RASMOL
- CHIME

Structural Classification of Proteins

- Over 1000 protein families known
 - Sequence alignment, motif finding, block finding, similarity search
- **SCOP** (Structural Classification of Proteins)
 - Based on structural & evolutionary relationships.
 - Contains ~ 40,000 domains
 - Classes (groups of folds), Folds (proteins sharing folds), Families (proteins related by function/evolution), Superfamilies (distantly related proteins)

SCOP Family View

The screenshot shows the NCSA Mosaic WWW browser interface. At the top, the document title is "SCOP: Family: Interleukin 8-like" and the URL is "http://scop.mrc-lmb.cam.ac.uk/scop/data/scop.0.004". The main content area is titled "Structural Classification of Proteins" and shows the "Family: Interleukin 8-like". Below this, a "Lineage" tree is shown, followed by a list of "Proteins". The first protein is "Interleukin-8", which is further detailed with its PDB entry names and chain information. Annotations with arrows point to various elements: "scop navigation buttons" at the top, "click here to display protein in 3D-viewer" pointing to a link in the lineage, "click here for sequence and references (NCBI)" pointing to a link in the protein list, "PDB entry names" pointing to the protein list, "click here to fetch image" pointing to a link, and "keyword search facility" pointing to a search box. Two 3D viewers are shown: "RasMol Version 2.4" displaying a protein structure, and "xv 3.00: scratch/xcaa09590.gif" displaying a comparison of Human MIP-1β and Interleukin 8 Dimers.

Figure 2. A typical scop session is shown on a unix workstation. A scop page, of the Interleukin 8-like family, is displayed by the WWW browser program (NCSA Mosaic) (Schatz & Hardin, 1994). Navigating through the tree structure is accomplished by selecting any underlined entry; by clicking on buttons (at the top of each page) and by keyword searching (at the bottom of each page). The static image comparing two proteins in this family was downloaded by clicking on the icon indicated and is displayed by image-viewer program xv. By clicking on one of the green icons, commands were sent to a molecular viewer program (RasMol) written by Roger Sayle (Sayle, 1994), instructing it to automatically display the relevant PDB file and colour the domain in question by secondary structure. Since sending large PDB files over the network can be slow, this feature of scop can be configured to use local copies of PDB files if they are available. Equivalent WWW browsers, image display programs and molecular viewers are also available free for Windows PC and Macintosh platforms.

CATH: Protein Structure Classification

- Semi-automatic classification; ~36K domains
- 4 levels of classification:
 - Class (C), depends on sec. Str. Content
 - α class, β class, α/β class, $\alpha+\beta$ class
 - Architecture (A), orientation of sec. Str.
 - Topology (T), topological connections &
 - Homologous Superfamily (H), similar str and functions.

DALI/FSSP Database

- Completely automated; 3724 domains
- Criteria of compactness & recurrence
- Each domain is assigned a Domain Classification number DC_l_m_n_p representing fold space attractor region (l), globular folding topology (m), functional family (n) and sequence family (p).

Structural Alignment

- What is structural alignment of proteins?
 - 3-d superimposition of the atoms as "best as possible", i.e., to minimize RMSD (root mean square deviation).
 - Can be done using **VAST** and **SARF**
- Structural similarity is common, even among proteins that do not share sequence similarity or evolutionary relationship.

Other databases & tools

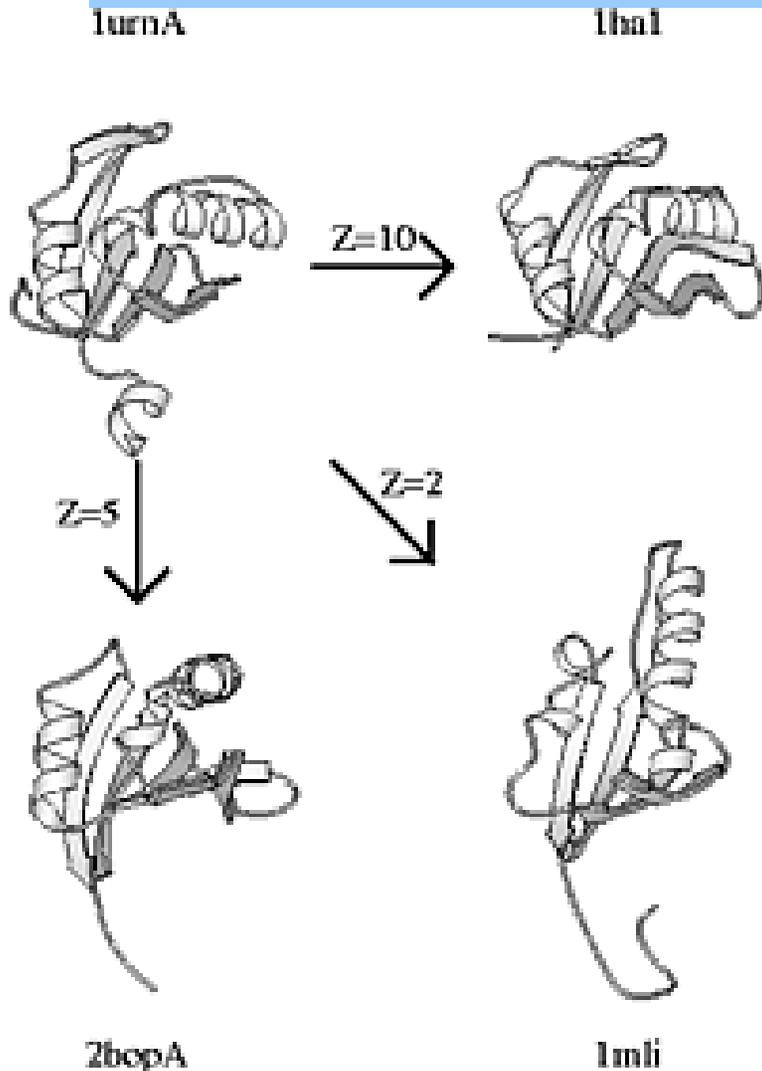
- **MMDB** contains groups of structurally related proteins
- **SARF** structurally similar proteins using secondary structure elements
- **VAST** Structure Neighbors
- **SSAP** uses double dynamic programming to structurally align proteins

5 Fold Space classes



Attractor 1 can be characterized as alpha/beta, attractor 2 as all-beta, attractor 3 as all-alpha, attractor 5 as alpha-beta meander (1mli), and attractor 4 contains antiparallel beta-barrels e.g. OB-fold (1prtF).

Fold Types & Neighbors



Structural neighbours of 1urnA (top left). 1mli (bottom right) has the same topology even though there are shifts in the relative orientation of secondary structure elements.

Sequence Alignment of Fold Neighbors

B

```
1urnA  --RPNHTIYINNLNEKI-----KKDELKKSLHAIFSRFG---QILDILV-SRS---LKM---
Z=10      *      *      *      *      *      *
1ha1    ahLTVKKIFVGGIKEDT-----EEHHLRDYFEOYG---KIEVIEI-MTDrgsGKK---
Z=5      *
2bopA   ----sCFALIS-GTANO-----vKCYRFRVKKNHRHR-----YENCTTtWFT---Vadnga
Z=2      *
1mli    ---mlFHVKMTTVKLpvdmdpakatgIkadeKELAQRIgregTWRHLWR-IAG-----

1urnA   ----RGQAFVIFKEV--SSATNALRSMQGFPFYDKPMRIQYAKTSDIIAKM-----
Z=10     **  ***  *      *      *
1ha1     ----RGFAFVTFDDH--DSVDKIVIO-kyHTVNGHNCEVRKAL-----
Z=5      *      *      *      *      *      *
2bopA   erggQAQILITFGSP--SORODFLKHVPLPP---GMNISGF-----tASLdf-----
Z=2      *      *      **      *      *
1mli     ----HYANYSVFDVpsvEALHDTLMQLpLFPY---MDIEVD-----gLCRHpssihsddr
```

Frequent Fold Types



(141) 1hdcA:1
alpha/beta domain



(85) 1mfaA:3
immunoglobulin fold



(63) 1ceo:2
TIM barrel



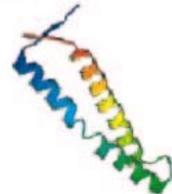
(43) 1befA:1
helical bundle



(36) 2pii:2
alpha/beta-meander



(33) 1vdfA:1
single helix



(27) 1grj:2
coiled coil



(25) 1bbt2:1
beta-meander



(19) 1rro:2
EF-hand



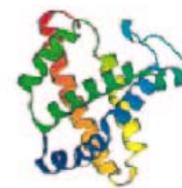
(18) 1oetC:3
HTH-motif



(18) 1ptf:1
OB-fold



(17) 3grs:2
FAD/NAD binding domain



(14) 1mbd:1
globin fold



(13) 1vin:3
cyclin fold



(13) 1aozA:15
blue copper protein



(13) 1lcf:17
periplasmic binding protein



(12) 1eelA:3



(12) 1epaA:1
lipocalin fold



(12) 2arcA:4
beta-roll



(12) 2yhx:3
actin fold

Protein Structure Prediction

- **Holy Grail** of bioinformatics
- **Protein Structure Initiative** to determine a set of protein structures that span protein structure space sufficiently well. **WHY?**
 - Number of folds in natural proteins is limited. Thus a newly discovered proteins should be within modeling distance of some protein in set.
- **CASP**: Critical Assessment of techniques for structure prediction
 - To stimulate work in this difficult field

PSP Methods

- *homology*-based modeling
- methods based on *fold recognition*
 - *Threading* methods
- *ab initio* methods
 - From first principles
 - With the help of databases

ROSETTA

- Best method for PSP
- As proteins fold, a large number of partially folded, low-energy conformations are formed, and that local structures combine to form more global structures with minimum energy.
- Build a database of known structures (I-sites) of short sequences (3-15 residues).
- Monte Carlo simulation assembling possible substructures and computing energy

Threading Methods

- *See p471, Mount*

- `http://www.bioinformaticsonline.org/links/ch_10_t_7.html`

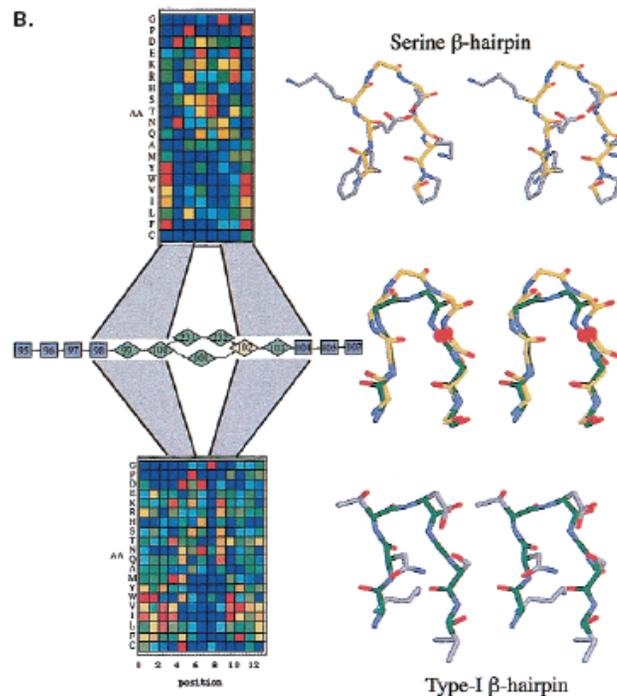
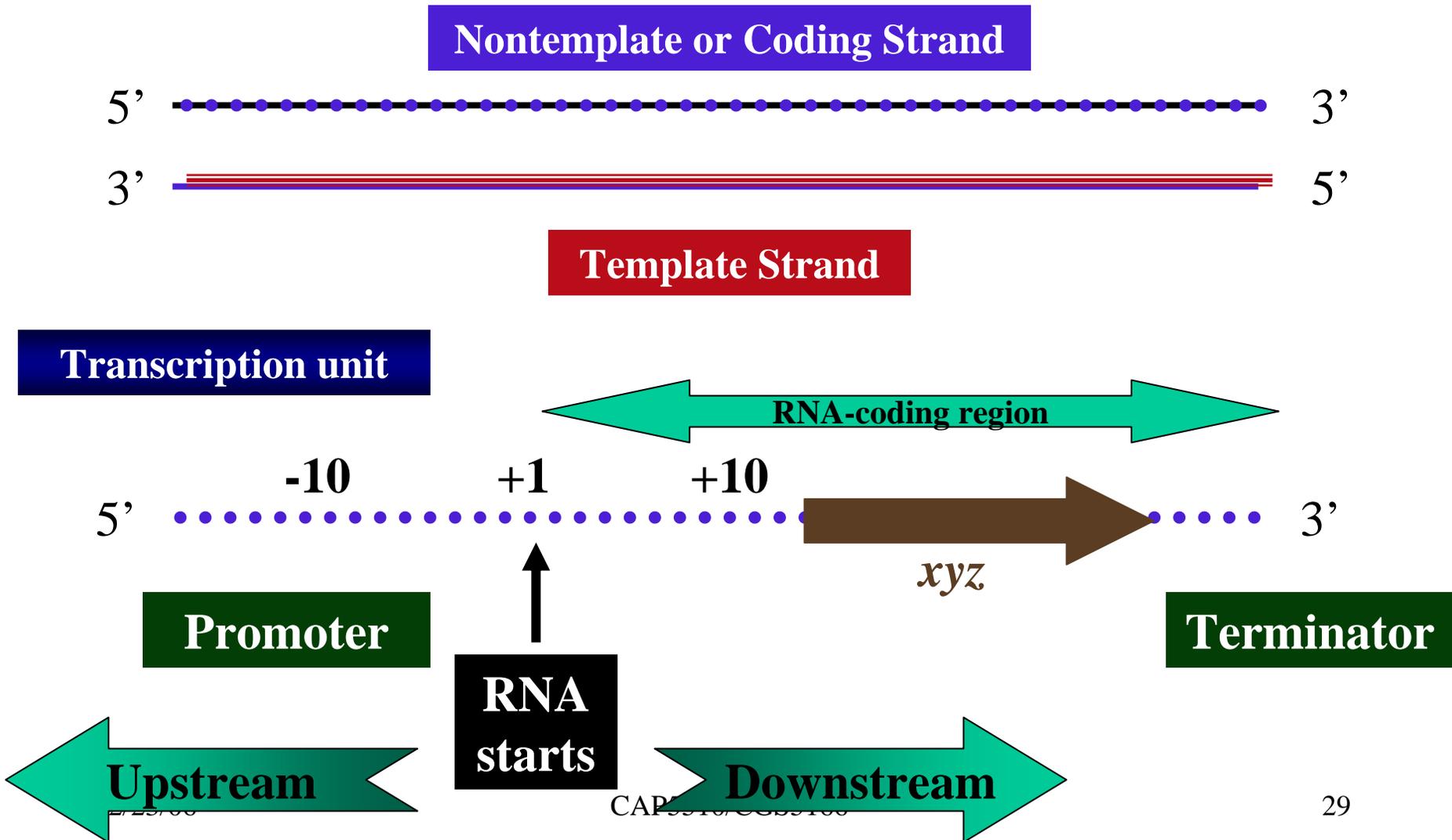


FIGURE 10.30. A hidden Markov model (discrete state-space model) of protein three-dimensional structure. (B) HMM called HMMSTR based on I-sites, 3- to 15-amino-acid patterns that are associated with three-dimensional structural features. The two matrices with colored squares represent alignment of sets of patterns that are found to be associated with a structure, in this case the hairpin turns shown on the right. Each column in the table corresponds to the amino acid variation found for one structural position in one of the turns. (*Blue* side chains) Conserved nonpolar residues; (*green*) conserved polar residues; (*red*) conserved proline; and (*orange*) conserved glycine. The two hairpins are aligned structurally in the middle structure on the right and the observed variation in the corresponding amino acid positions is represented by the HMM between the matrices on the left. The HMM represents an alignment of the two hairpin structural motifs in three-dimensional space and an alignment of the sequences. A short mismatch in the turn is represented by splitting the model into two branches. The shaped icons represent states, each of which represents a structure and a sequence position. Each state contains probability distributions about the sequence and structural attributes of a single position in the motif, including the probability of observing a particular amino acid, secondary structure, Φ - Ψ backbone angles, and structural context, e.g., location of β strand in a β sheet. Rectangles are predominantly β -strand states, and diamonds are predominantly turns. The color of the icon indicates a sequence preference as follows: (*blue*) hydrophobic; (*green*) polar; and (*yellow*) glycine. Numbers in icons are arbitrary identification numbers for the HMM states. There is a transition probability of moving from each state in the model to the next, as in HMMs that represent *msa*'s. This model is a small component of the main HMMSTR model that represents a merging of the entire I-sites library. Three different models, designated λ^P , λ^C , and λ^R , are included in HMMSTR, which differ in details as to how the alignment of the I-sites was obtained to design the branching patterns (topology) of the model and which structural data were used to train the model. HMMSTR may be used for a variety of different predictions, including secondary structure prediction, structural context prediction, and Φ - Ψ dihedral angle prediction. Predictions are made by aligning the model with a sequence, finding if there is a high-scoring alignment, and deciphering the highest-scoring path through the model. The HMMSTR program may be downloaded or used on a server that can be readily located by a Web search. (B, reprinted, with permission, from Bystroff et al. 2000 [©2000 Elsevier].)

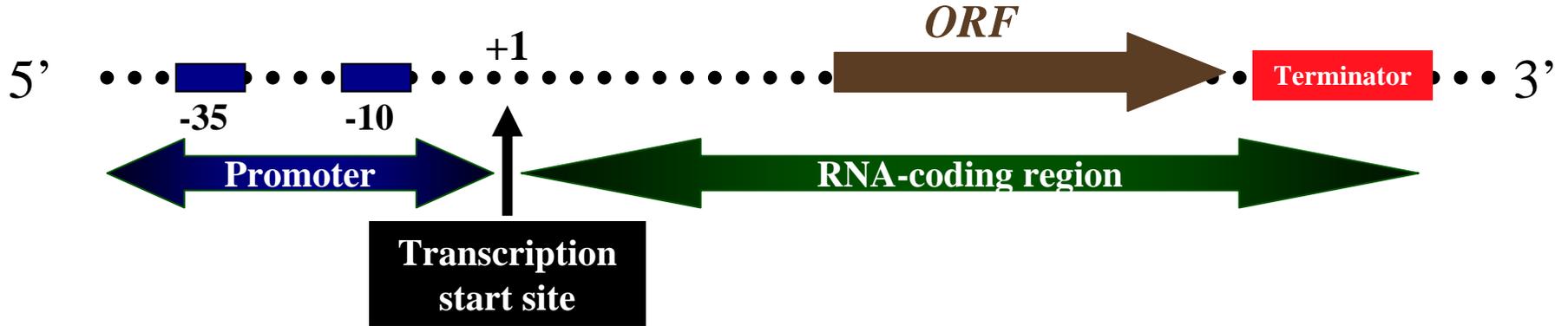
Nomenclature

RNA Polymerization occurs 5' to 3'

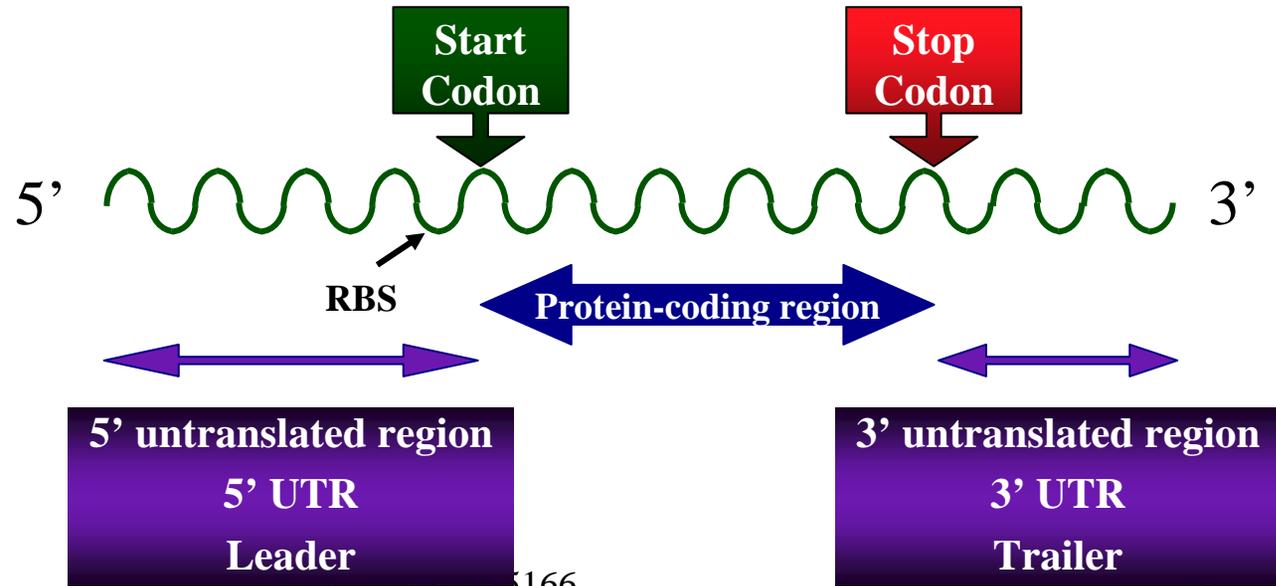


Transcriptional unit and single gene mature mRNA

Transcriptional unit



mRNA



RBS
Ribosome
binding site

2/23/06

CGS510/CGS5166

30

Messenger RNA or mRNA

Initiation Codon

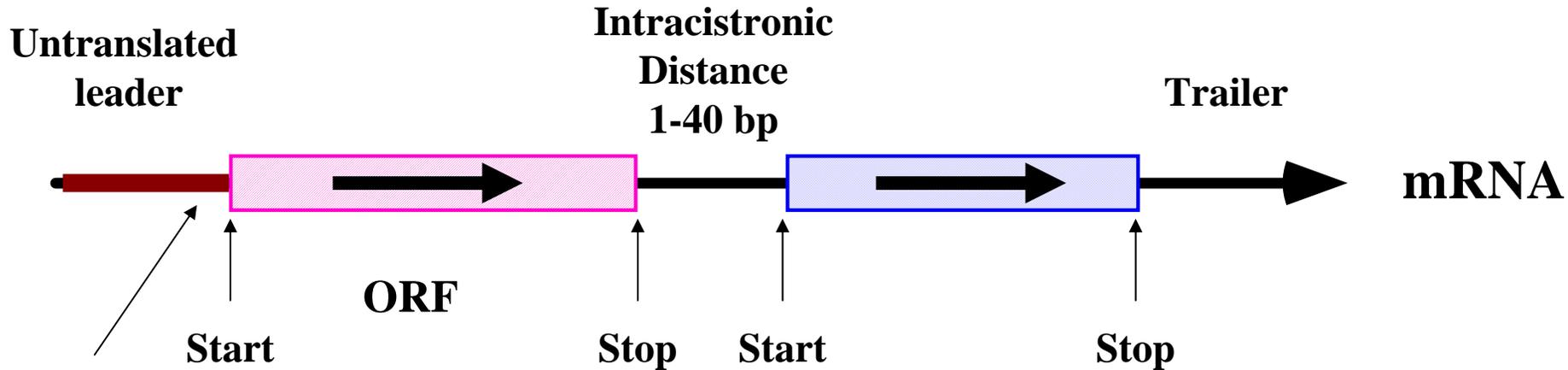
AUG **Methionine**

Termination Codons

Others:

GUG **Valine**
UUG **Leucine**
AUU **Isoleucine**

UAA **Ochre**
UAG **Amber**
UGA **Opal**



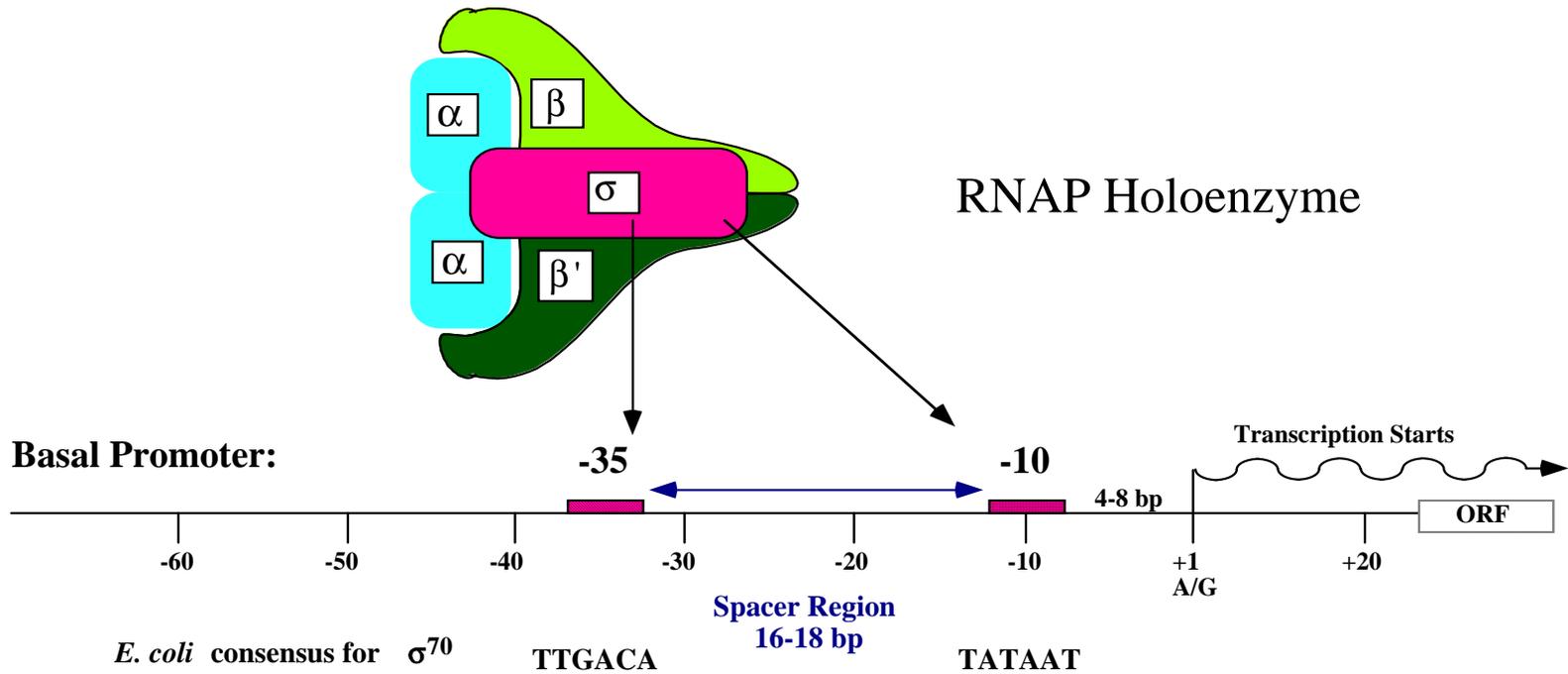
RBS
Ribosome Binding Site
Shine-Dalgarno Sequence

7 bp upstream of start codon
5'--AGGAGG--3'

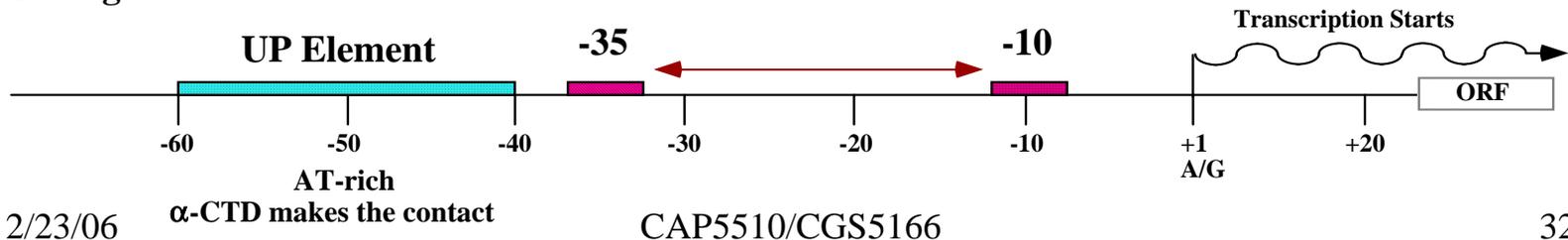
Coding region
Open Reading Frame (ORF)

Reading frame is one of three possible ways of reading a nucleotide sequence as a series of triplets.

Transcriptional machinery: RNA Polymerase and DNA



Stronger Promoter:



2/23/06

32

Prokaryotic Gene Characteristics

76 ■ CHAPTER 9

DNA PATTERNS IN THE *E. coli* *lexA* GENE

GENE SEQUENCE	PATTERN
1 GAATTCGATAAATCTCTGGTTTTATTTGTGCAGTTTATGGTT	CTGNNNNNNNNNNCAG
TT	TTGACA
41 CCAAATCGCCTTTTGTGCTATATACTCACAGCATAACTG	CTGNNNNNNNNNNCAG
CAA -35 -10 TATACT >	TATAAT, > mRNA start
81 TATAATCACCCAGGGGGCGAATGAAAGCGTTAACGGCCA	CTGNNNNNNNNNNCAG
+10 GGGGG Ribosomal binding site	GGAGG
121 GGCAACAAGAGGTGTTTGTATCTCATCCGTGATCACATCAG	
161 CCAGACAGGTATGCGCCGACGCGTGCAGAAATCGCCAG	ATG
201 CGTTTGGGGTTCGGTTCCCAAACGCGCTGAAGAATCATC	
241 TGAAGGCGCTGGCACGCAAGGCGTTATTGAAATTTGTTT	
281 CGCGCATCACGCGGGATTTCGTCTGTGTGCAGGAAGAGGAA	
321 GAAGGGTTGCCGCTGGTAGGTCGTGTGGCTGCCGGTGAAC	
361 CACTTCTGGCGCAACAGCATATTGAAGGTCATTATCAGGT	OPEN READING FRAME
401 CGATCCTTCCTTATTCAGCCGAATGCTGATTTCTGCTG	
441 CGCGTCAGCGGGATGTCGATGAAAGATATCGGCATTATGG	
481 ATGGTGAAGTGTGCTGGCAGTGCATAAACTCAGGATGTACG	
521 TAACGGTCAGGTCGTTGTCGACGTATTGATGACGAAGTT	
561 TCCCTTTCAGCCCTTCTTAAAAACAGGGCAATTAAGTCAAC	
601 TGTTCGCAGAAATAGCGAGTTTAAACCAATTTGTCGTTGA	
641 CCTTCGTCAGCAGAGCTTCCACCATTAAGGGCTGGCCGTT	TAA
681 GGGGTTATTCGCAACGGCGACTGGCTGTAACATATCTCTG	
721 AGACCGCGATGCGCCCTGGCGTCCGCGTTTGTITTTTCATC	
761 TCTCTTCATCAGGCTTGTCTGCATGGCATTCCCTCACTTCA	
801 TCTGATAAAGCACTCTGGCATCTCGCCTTACCCATGATTT	
841 TCTCCAAATATCACCGTTTCCGTTGCTGGGACTGGTTCGATAC	
881 GGCGTAAATGGTTCATCTTGATAGCCCGGTTTATTTGGGC	
921 GGCGTGGCGGTTGGCGCAACGGCGGACCAAGCT	

Shown are matches to approximate consensus binding sites for LexA repressor (CTGNNNNNNNNNNCAG), the -10 and -35 promoter regions relative to the start of the mRNA (TTGACA and TATAAT), the ribosomal binding site on the mRNA (GGAGG), and the open reading frame (ATG...TAA). Only the second two of the predicted LexA binding sites actually bind the repressor.

FIGURE 9.6. The promoter and open reading frame of the *E. coli* *lexA* gene.

Motif Detection (TFBMs)

- See evaluation by Tompa et al.
 - [bio.cs.washington.edu/assessment]
- **Gibbs Sampling Methods:** AlignACE, GLAM, SeSiMCMC, MotifSampler
- **Weight Matrix Methods:** ANN-Spec, Consensus,
- **EM:** Improbizer, MEME
- **Combinatorial & Misc.:** MITRA, oligo/dyad, QuickScore, Weeder, YMF

Start and Stop Codon Distribution

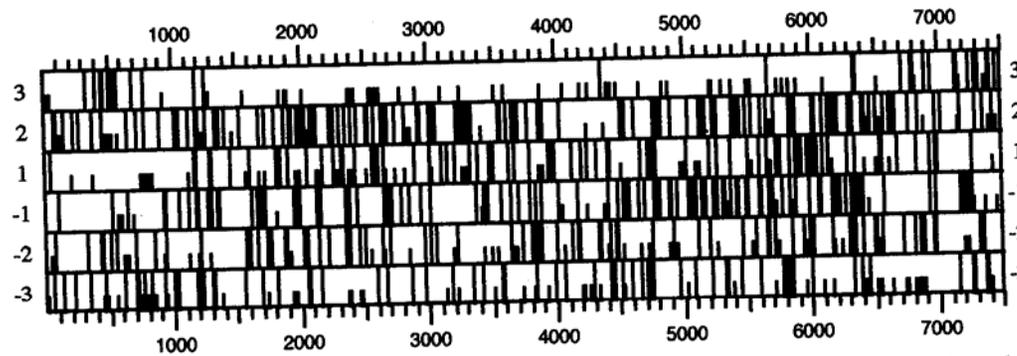
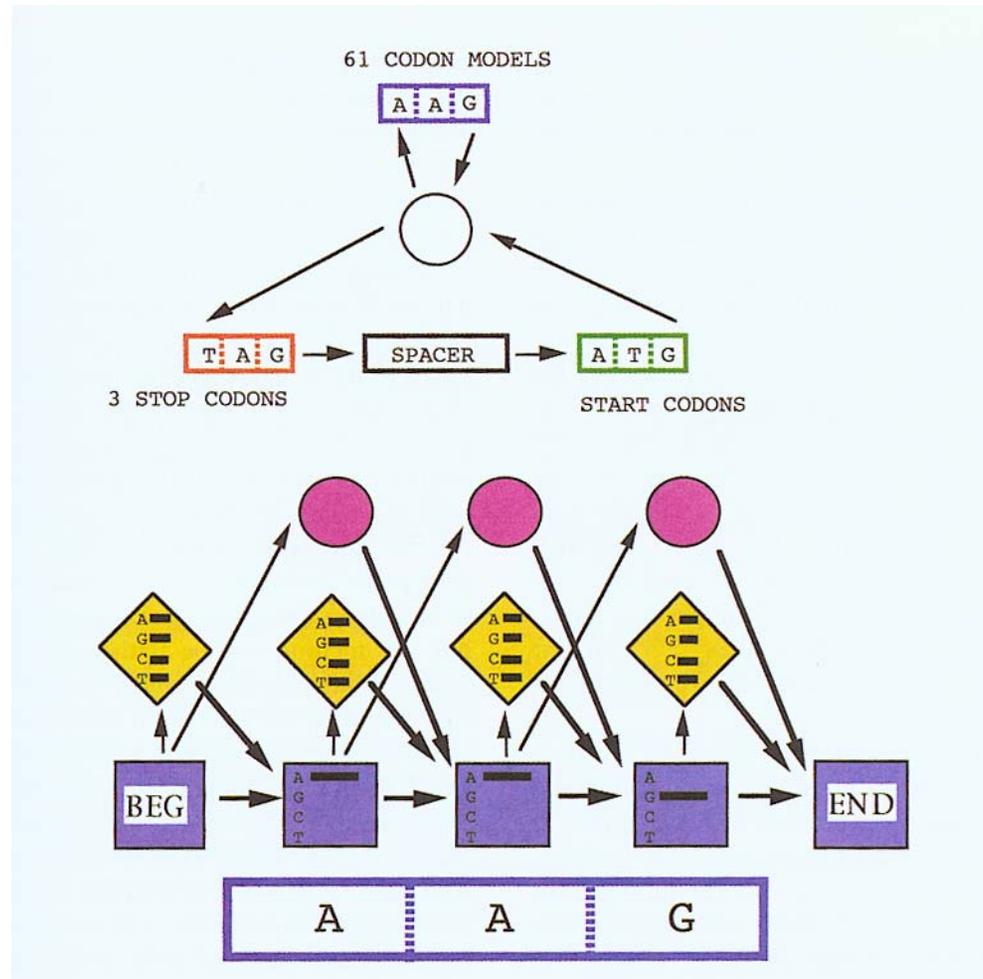


FIGURE 9.1. ORF map of a portion of the *E. coli lac* operon using the DNA STRIDER program (Marck 1988). Shown are AUG and termination codons as one-half and full vertical bars, respectively, in all six possible reading frames. The *lacZ* gene is visible as an ORF that runs from positions 1284 to 4355 in frame 3.

Genetic Code

		Second letter					
		U	C	A	G		
First letter	U	UUU UUC	UCU UCC UCA UCG	UAU UAC	UGU UGC	U	C
		UUA UUG		UAA UAG		UGA UGG	A
	C	CUU CUC CUA CUG	CCU CCC CCA CCG	CAU CAC	CGU CGC CGA CGG	U	C
				CAA CAG			A
A	AUU AUC AUA	ACU ACC ACA ACG	AAU AAC	AGU AGC	U	C	
	AUG		AAA AAG		AGA AGG	A	G
G	GUU GUC GUA GUG	GCU GCC GCA GCG	GAU GAC	GGU GGC GGA GGG	U	C	
			GAA GAG			A	G

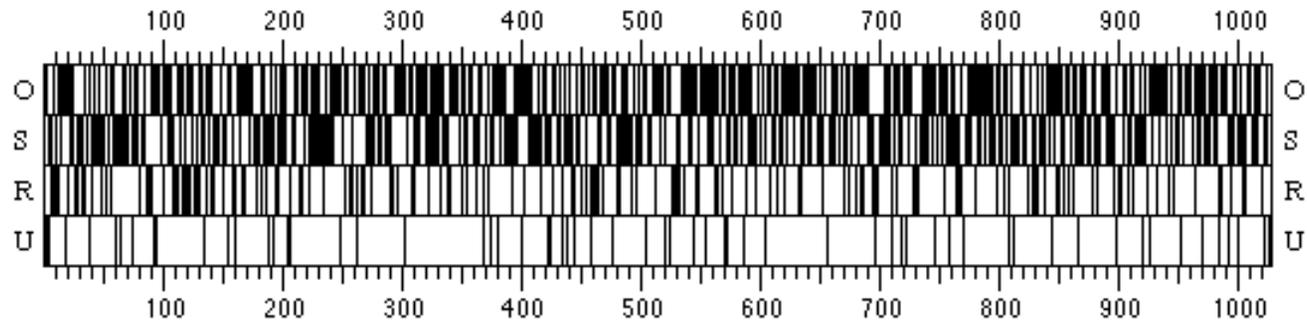
Recognizing Codons



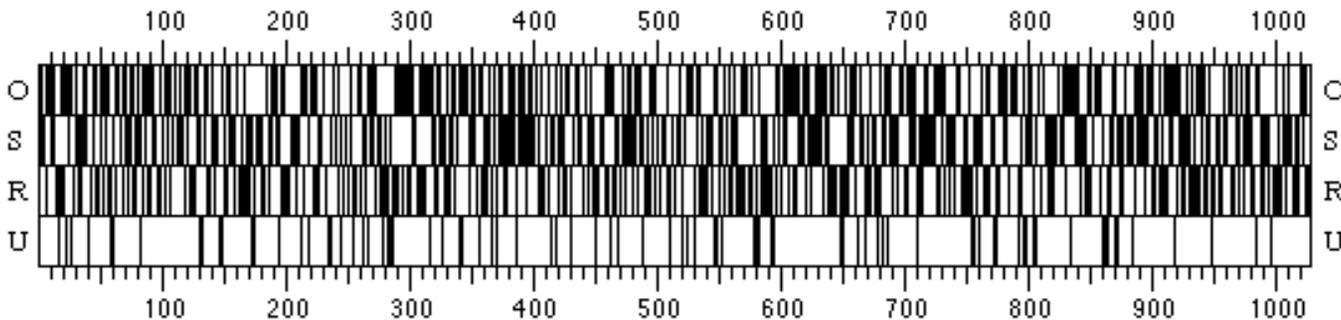
Codon Bias

- Some codons preferred over others.

O = optimal
S = suboptimal
R = rare
U = unfavorable



Frame Shift 1

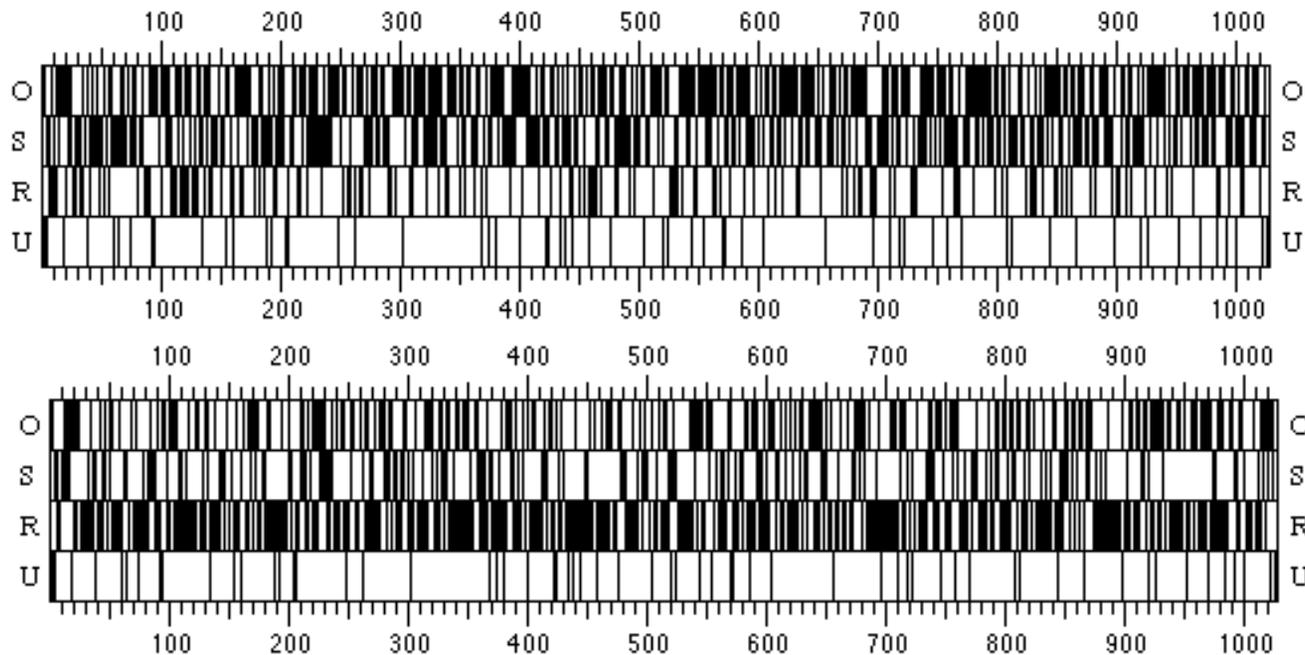


Frame Shift 2

Codon Bias

- Codon biases specific to organisms

O = optimal
S = suboptimal
R = rare
U = unfavorable



Same Frames;
Different labeling
of codon types
(i.e., from yeast)

Eukaryotic Gene Prediction

- Complicated by introns & alternative splicing
- Exons/introns have different GC content.
- Many other measures distinguish exons/introns
- Software:
 - **GENEPARSER** Snyder & Stormo (NN)
 - **GENIE** Kulp, Haussler, Reese, Eckman (HMM)
 - **GENSCAN** Burge, Karlin (Decision Trees)
 - **XGRAIL** Xu, Einstein, Mural, Shah, Uberbacher (NN)
 - **PROCRUSTES** Gelfand (Formal Languages)
 - **MZEF** Zhang

Introns/Exons in *C. elegans*

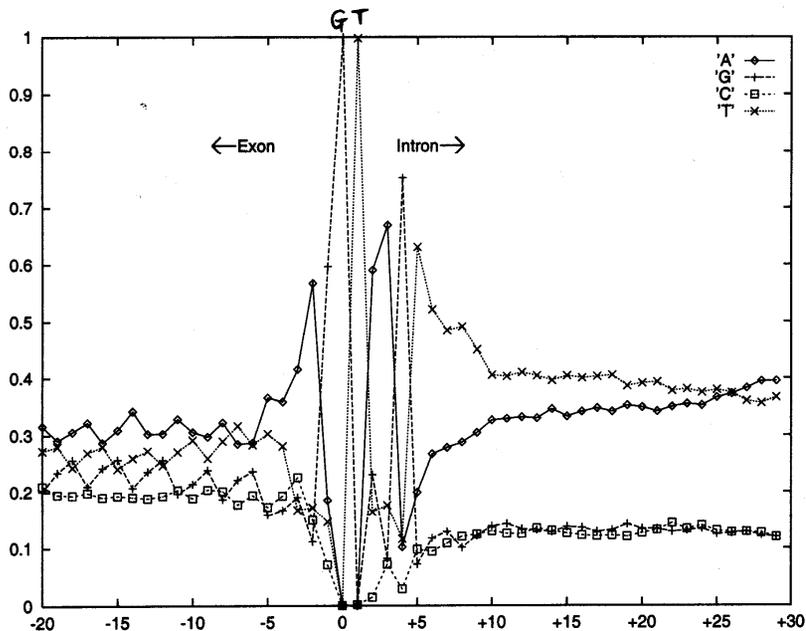


Figure 2: Profile of the same 5' collection but around a larger window.

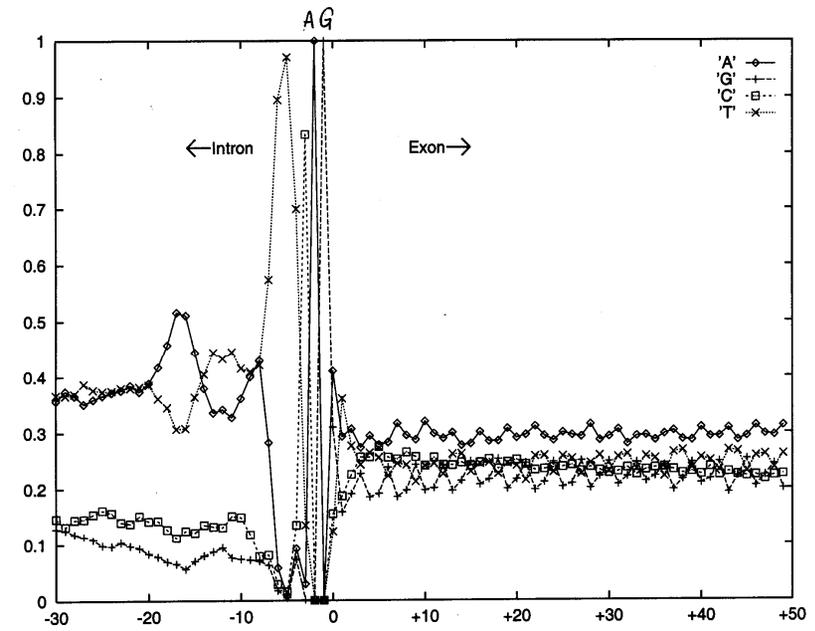
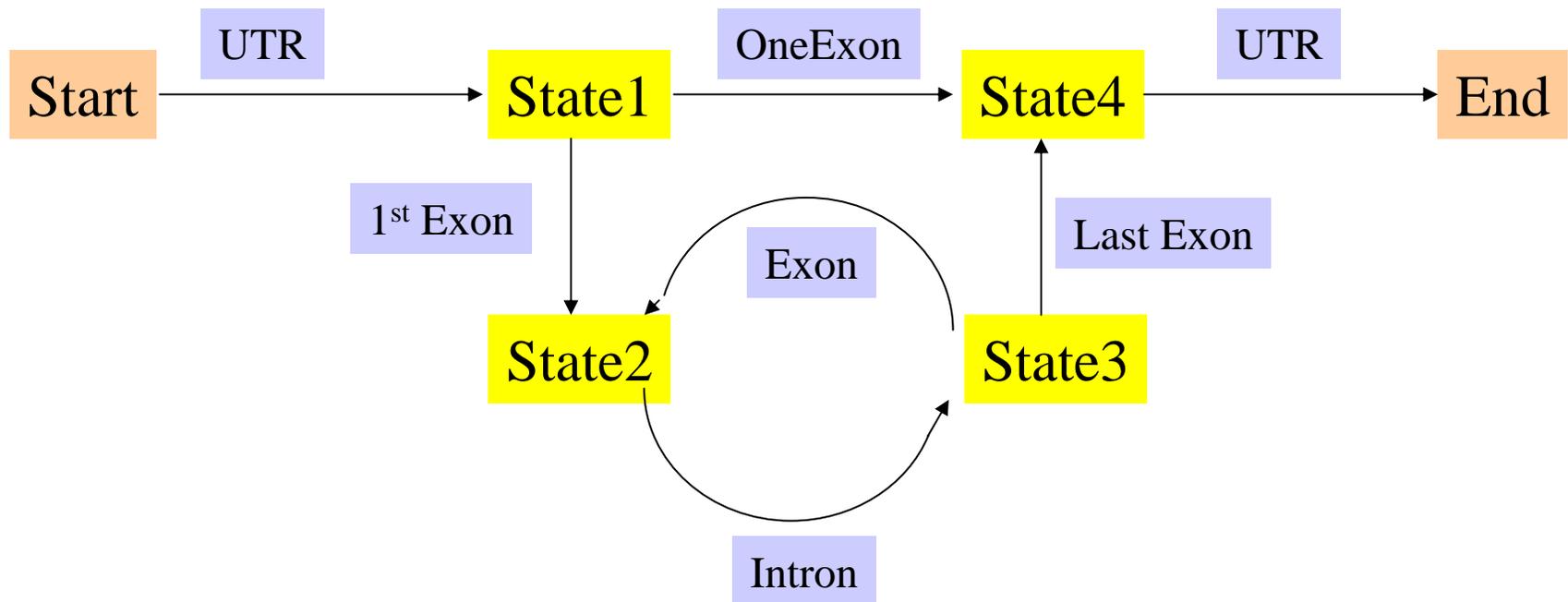


Figure 4: Profile of 8,192 sequences of length 80 around the 3' site. The first position in the exon is labeled 0.

- 8192 Introns in *C. elegans*: [GT...AG]
- Vary in lengths from 30 to over 600; Complexity varies

HMM structure for Gene Finding



Motifs in Protein Sequences

Motifs are combinations of secondary structures in proteins with a specific **structure** and a specific **function**. They are also called **super-secondary structures**.

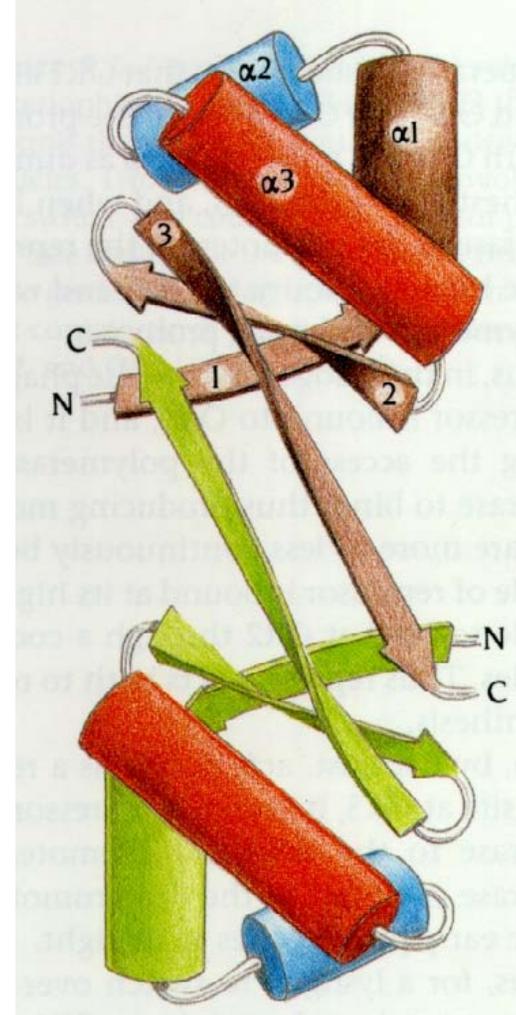
Examples: Helix-Turn-Helix, Zinc-finger, Homeobox domain, Hairpin-beta motif, Calcium-binding motif, Beta-alpha-beta motif, Coiled-coil motifs.

Several motifs may combine to form **domains**.

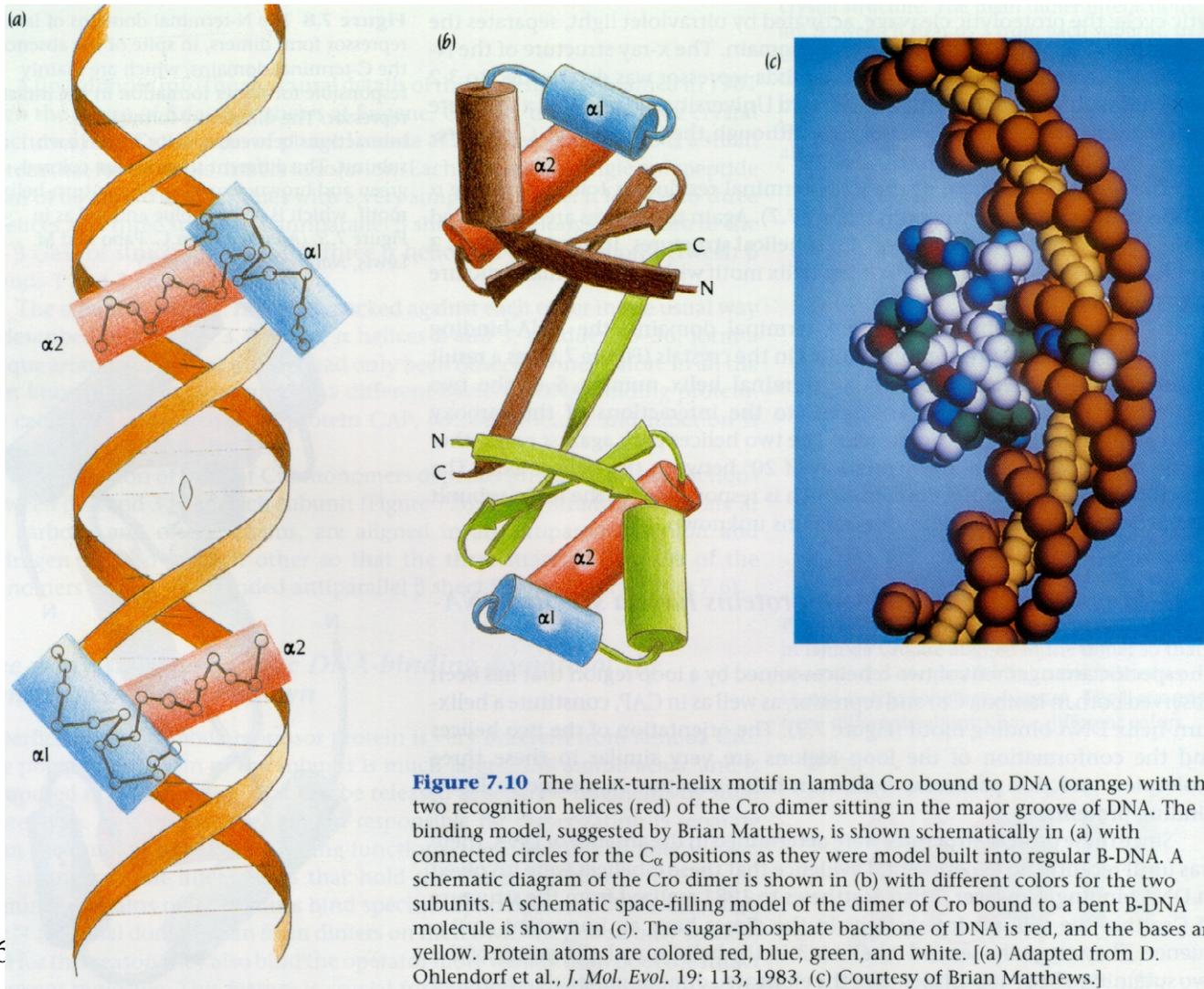
- Serine proteinase domain, Kringle domain, calcium-binding domain, homeobox domain.

Helix-Turn-Helix Motifs

- Structure
 - 3-helix complex
 - Length: 22 amino acids
 - Turn angle
- Function
 - Gene regulation by binding to DNA



DNA Binding at HTH Motif



HTH Motifs: Examples

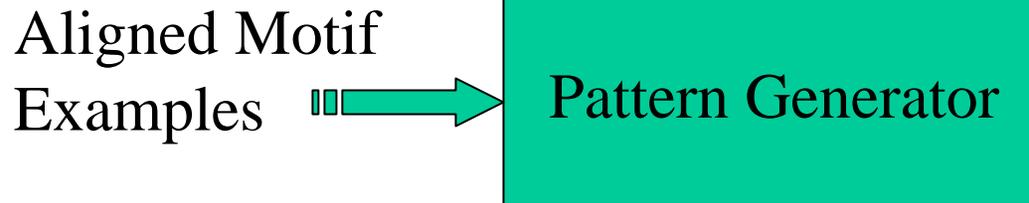
<i>Loc</i>	<i>Protein Name</i>	<i>Helix 2</i>									<i>Turn</i>				<i>Helix 3</i>								
		-1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
14	Cro	F	G	Q	E	K	T	A	K	D	L	G	V	Y	Q	S	A	I	N	K	A	I	H
16	434 Cro	M	T	Q	T	E	L	A	T	K	A	G	V	K	Q	Q	S	I	Q	L	I	E	A
11	P22 Cro	G	T	Q	R	A	V	A	K	A	L	G	I	S	D	A	A	V	S	Q	W	K	E
31	Rep	L	S	Q	E	S	V	A	D	K	M	G	M	G	Q	S	G	V	G	A	L	F	N
16	434 Rep	L	N	Q	A	E	L	A	Q	K	V	G	T	T	Q	Q	S	I	E	Q	L	E	N
19	P22 Rep	I	R	Q	A	A	L	G	K	M	V	G	V	S	N	V	A	I	S	Q	W	E	R
24	CII	L	G	T	E	K	T	A	E	A	V	G	V	D	K	S	Q	I	S	R	W	K	R
4	LacR	V	T	L	Y	D	V	A	E	Y	A	G	V	S	Y	Q	T	V	S	R	V	V	N
167	CAP	I	T	R	Q	E	I	G	Q	I	V	G	C	S	R	E	T	V	G	R	I	L	K
66	TrpR	M	S	Q	R	E	L	K	N	E	L	G	A	G	I	A	T	I	T	R	G	S	N
22	BlaA Pv	L	N	F	T	K	A	A	L	E	L	Y	V	T	Q	G	A	V	S	Q	Q	V	R
23	TrpI Ps	N	S	V	S	Q	A	A	E	Q	L	H	V	T	H	G	A	V	S	R	Q	L	K

Basis for New Algorithm

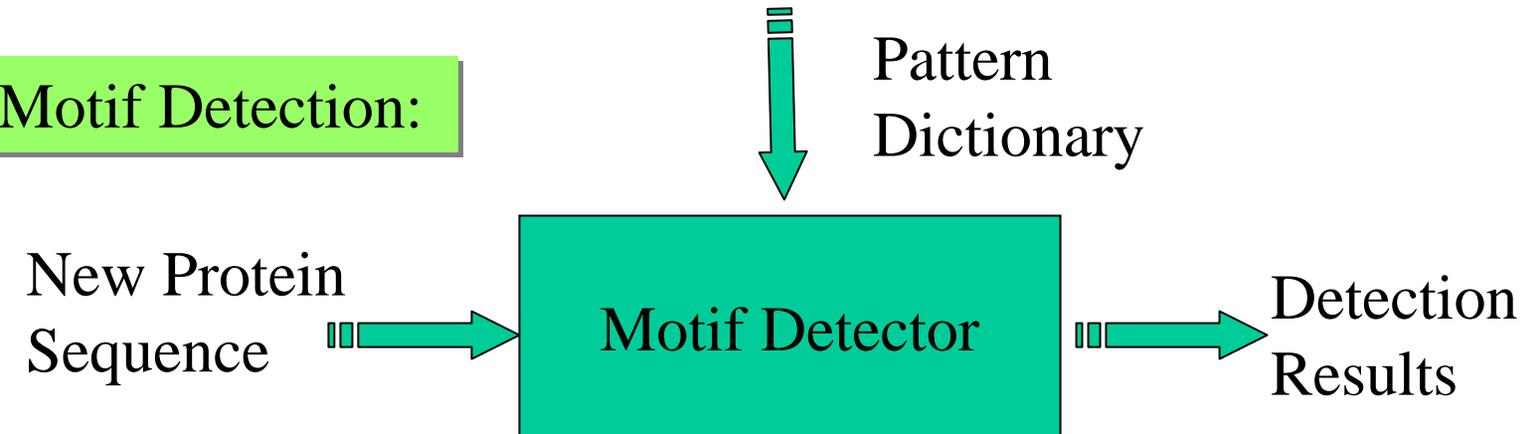
- Combinations of residues in specific locations (may not be contiguous) contribute towards stabilizing a structure.
- Some **reinforcing** combinations are relatively rare.

New Motif Detection Algorithm

Pattern Generation:



Motif Detection:



Patterns

Loc	Protein Name	Helix 2									Turn				Helix 3								
		-1	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
14	Cro	F	G	Q	E	K	T	A	K	D	L	G	V	Y	Q	S	A	I	N	K	A	I	H
16	434 Cro	M	T	Q	T	E	L	A	T	K	A	G	V	K	Q	Q	S	I	Q	L	I	E	A
11	P22 Cro	G	T	Q	R	A	V	A	K	A	L	G	I	S	D	A	A	V	S	Q	W	K	E
31	Rep	L	S	Q	E	S	V	A	D	K	M	G	M	G	Q	S	G	V	G	A	L	F	N
16	434 Rep	L	N	Q	A	E	L	A	Q	K	V	G	T	T	Q	Q	S	I	E	Q	L	E	N
19	P22 Rep	I	R	Q	A	A	L	G	K	M	V	G	V	S	N	V	A	I	S	Q	W	E	R
24	CII	L	G	T	E	K	T	A	E	A	V	G	V	D	K	S	Q	I	S	R	W	K	R
4	LacR	V	T	L	Y	D	V	A	E	Y	A	G	V	S	Y	Q	T	V	S	R	V	V	N
167	CAP	I	T	R	Q	E	I	G	Q	I	V	G	C	S	R	E	T	V	G	R	I	L	K
66	TrpR	M	S	Q	R	E	L	K	N	E	L	G	A	G	I	A	T	I	T	R	G	S	N
22	BlaA Pv	L	N	F	T	K	A	A	L	E	L	Y	V	T	Q	G	A	V	S	Q	Q	V	R
23	TrpI Ps	N	S	V	S	Q	A	A	E	Q	L	H	V	T	H	G	A	V	S	R	Q	L	K

- Q1 G9 N20
- A5 G9 V10 I15

Pattern Mining Algorithm

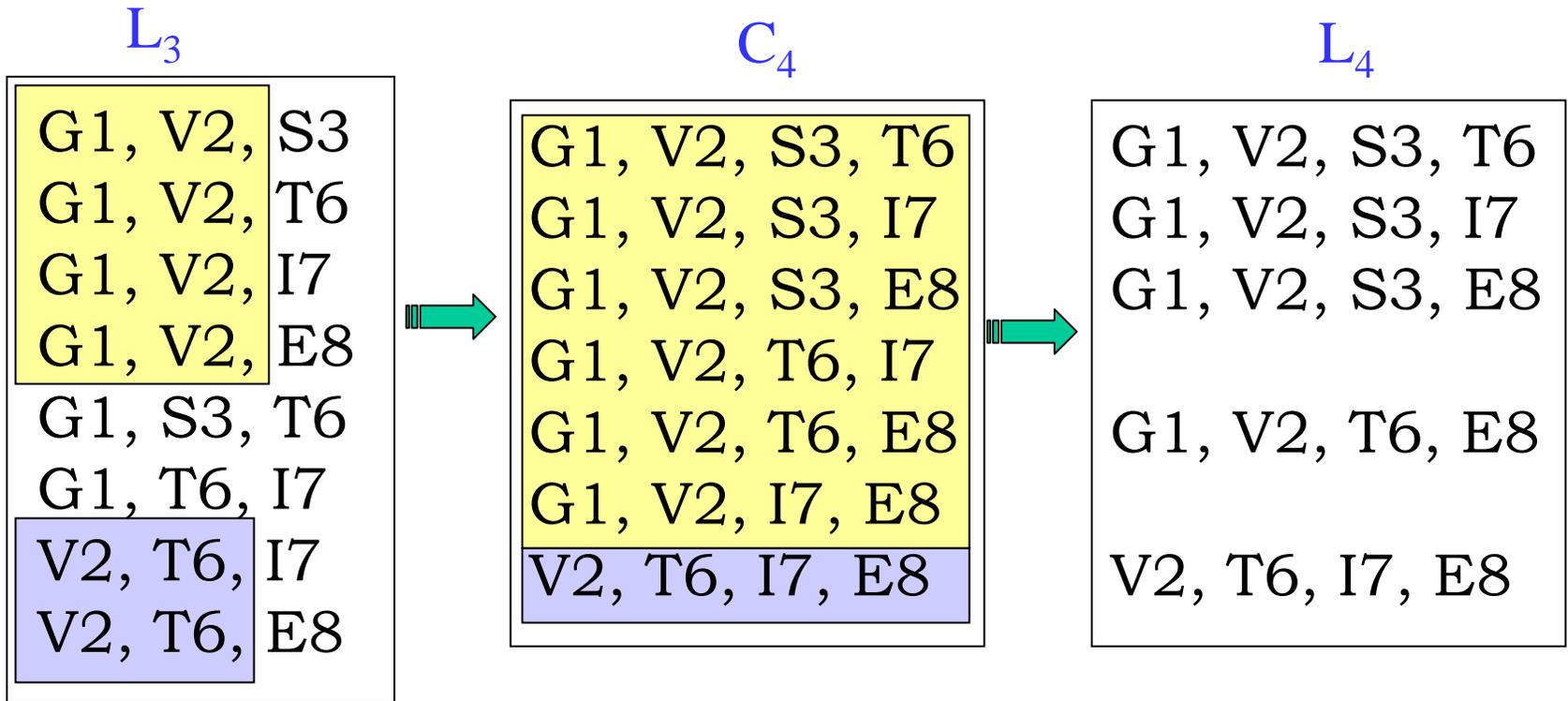
Algorithm **Pattern-Mining**

Input: Motif length **m**, support threshold **T**,
list of aligned motifs **M**.

Output: Dictionary **L** of frequent patterns.

1. $L_1 :=$ All frequent patterns of length 1
2. **for** $i = 2$ **to** **m** **do**
3. $C_i :=$ **Candidates**(L_{i-1})
4. $L_i :=$ Frequent candidates from C_i
5. **if** ($|L_i| \leq 1$) **then**
6. **return** **L** as the union of all L_j , $j \leq i$.

Candidates Function



Motif Detection Algorithm

Algorithm **Motif-Detection**

Input : Motif length **m**, threshold score **T**, pattern dictionary **L**, and input protein sequence **P**[1..n].

Output : Information about motif(s) detected.

1. **for** each location **i do**
2. **S** := **MatchScore**(**P**[**i**..**i+m-1**], **L**).
3. **if** (**S** > **T**) **then**
4. Report it as a possible motif

Experimental Results: GYM 2.0

<i>Motif</i>	<i>Protein Family</i>	<i>Number Tested</i>	<i>GYM = DE Agree</i>	<i>Number Annotated</i>	<i>GYM = Annot.</i>
<i>HTH Motif (22)</i>	Master	88	88 (100 %)	13	13
	Sigma	314	284 + 23 (98 %)	96	82
	Negates	93	86 (92 %)	0	0
	LysR	130	127 (98 %)	95	93
	AraC	68	57 (84 %)	41	34
	Rreg	116	99 (85 %)	57	46
	Total	675	653 + 23 (94 %)	289	255 (88 %)

Experiments

- Basic Implementation (Y. Gao)
- Improved implementation & comprehensive testing (K. Mathee, GN).
- Implementation for homeobox domain detection (X. Wang).
- Statistical methods to determine thresholds (C. Bu).
- Use of substitution matrix (C. Bu).
- Study of patterns causing errors (N. Xu).
- Negative training set (N. Xu).
- NN implementation & testing (J. Liu & X. He).
- HMM implementation & testing (J. Liu & X. He).