# CAP 5510: Introduction to Bioinformatics

## Giri Narasimhan

ECS 254; Phone: x3748

*giri@cis.fiu.edu*

www.cis.fiu.edu/~giri/teach/BioinfS07.html

# Pattern Discovery

# Patterns

- Nature **stumbles** upon **recipes** to accomplish tasks.
- With high probability, such recipes are reused.
- This causes the recipe to be conserved through evolution.
- Such recipes give rise to **patterns**.

# Why Pattern Discovery?

- ❑ **Modern Biomedical Research**
  - ● Generates a "ton of data".
  - ● Use analytical tools to find patterns in data.
- ❑ **Pattern Discovery** facilitates this process!
  - ● Pattern Discovery in sequences
  - ● Pattern Discovery in structures
  - ● Pattern Discovery in quantitative data
- ❑ Patterns help to detect members of a class
- ❑ Patterns help to characterize classes

# Sequence Patterns: Examples

❑ Protein active sites and functional domains

  ● For e.g., Zinc-finger motifs & Helix-turn-helix motifs

❑ Protein family signatures

❑ Signals in DNA e.g., protein binding sites

❑ MicroRNA and Anti-sense RNA

# Example 1: Protein Motifs

❑ DNA-binding motifs

   ● Helix-turn-Helix

❑ Motifs in $Cys_2His_2$-Zinc-binding proteins

Example: Zinc Finger Motif
…**Y**KC**C**GL**C**ERS**F**VEKSA**L**SR**H**ORV**H**KN…
3    6                    19        23

❑ Motifs in proteins that bind to [4Fe-4S]-complex

Example: Ferredoxin subfamily
…**C**xx**C**xx**C**xxx**CP**…

# How to Represent Patterns

- ❏ Consensus sequence
- ❏ Alignments
- ❏ LOGO format
- ❏ Frequency Matrices
- ❏ Weight Matrices (Profiles, PSSMs, PWMs)

# Pattern Representations

☐ *Consensus sequences*

**[Pribnow, 1975]**
```
TACGAT
TATAAT
TATAAT
GATACT
TATGAT
TATGTT
------
TATAAT  Consensus
```

Needs Alignment

```
TATRNT  Consensus
        w/ IUPAC
```
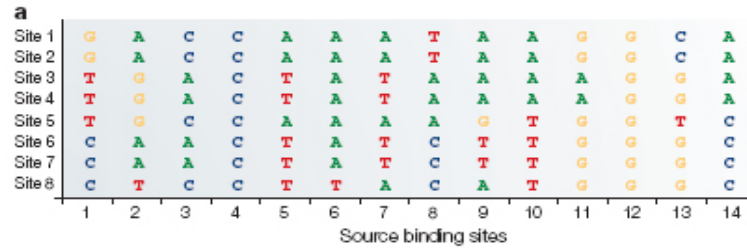
```
TATAAT  Multi-level
G CGC   Consensus
   T
```

# Pattern Representations

❑ Consensus sequences

❑ Weight Matrices (Profiles, PSSMs)

- Frequency Counts
- Relative Frequency Measures
- Normalized Measures
- Log-transformed Measures
- Information content
- "Logo" technique
- HMMs

# Pattern Representation: Weight Matrix

**Alignment**

**Consensus**

**Frequencies**

**Profile/ PSSM/PWM**

**Scoring a sequence against a profile**

**Visualizing a profile**

[Wasserman, Sandelin, Nat Genet, 2004]

# Formulae

❑ Prob of char **b** in position **i**:

$$p(b,i) = \frac{f_{b,i}}{N}$$

Frequency

\# Sequences

❑ Corrected prob:

$$P(b,i) = \frac{f_{b,i} + s(b)}{N + \sum_{a \in \mathbf{A}} s(a)}$$

PseudoCount

❑ Weight matrix entry:

$$W_{b,i} = \log_2 \frac{P(b,i)}{BP(b)}$$

Background Frequency

❑ Information content of position of **i**:

$$D_i = 2 + \sum_b P(b,i) \log_2 P(b,i)$$

[Wasserman, Sandelin, Nat Genet, 2004]

# Statistical Evaluation Fundamentals

❑ Probability of finding a sequence **w** in some position of a DNA/protein sequence (assuming independence at each position)

$$\Pr(w) = \prod_{i=1}^{m} \Pr(w_i)$$

❑ Pr(w$_i$) = BP(b) [Background Frequency]

# Statistical Evaluation

□ **Z-score** of a motif with a certain frequency: ➡️

$$z(w) = \frac{Obs(w) - Exp(w)}{\sqrt{Var(w)}}$$

□ **Information Content** or Relative Entropy of an alignment or profile: ➡️

$$IC(M) = \sum_{i=1}^{4} \sum_{j=1}^{m} m_{i,j} \log \frac{m_{i,j}}{b_i}$$

□ **Maximum a Posteriori** (MAP) Score: ➡️

$$MAP(M) = -\sum_{i=1}^{4} \sum_{j=1}^{m} n_{i,j} \log \frac{m_{i,j}}{b_i}$$

□ **Model Vs Background** Score: ➡️

$$L(w) = \frac{\Pr(w \mid M)}{\Pr(w \mid Bg)} = \prod_{j=1}^{m} \frac{m_{i,j}}{b_i}$$

# Pattern Discovery 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

❑ **Profile Method**
- 🔴 If many examples of the motif are known, then
  - ➢ **Training**: build a **Profile** and compute a **threshold**
  - ➢ **Testing**: **score** against profile

❑ **Combinatorial Pattern Discovery Methods**

❑ **Gibbs Sampling**

❑ **Expectation Method**

❑ **HMM**

# 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+FN)
- Positive Predictive Value = TP/(TP+FP)
- Performance Coefficient = TP/(TP+FN+FP)
- Correlation Coefficient =

# 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**KC**GL**CERS**F**VEKSAL**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?

# Supervised Pattern Discovery

❑ <u>Input</u>:   Alignment of known motifs, and

Query sequence

<u>Output</u>: Is the query sequence a motif?

● Profile Method [Gribskov et al., 1996]
  ➢ Build a profile from the alignment and score query sequence against the profile to decide if it "fits the profile".
  ➢ Need to pick a threshold score.
● Enumerative/Combinatorial Methods

# Profile HMMs

PROFILE METHOD, [M. Gribskov et al., '90]

| Location in Seq. | Sequence 1 2 3 4 5 6 | | Protein Name |
|---|---|---|---|
| 14 | G V S A S A | | Ka RbtR |
| 32 | G V S E M T | | Ec DeoR |
| 33 | G V S P G T | | Ec RpoD |
| 76 | G A G I A T | | Ec TrpR |
| 178 | G C S R E T | | Ec CAP |
| 205 | C L S P S R | | Ec AraC |
| 210 | C L S P S R | | St AraC |
| 13 | G V N K E T | | Br MerR |

START → STATE 1 → STATE 2 → STATE 3 → STATE 4 → STATE 5 → STATE 6 → END

# Combinatorial Method: GYM

Pattern Generation:

Aligned Motif Examples → Pattern Generator

Pattern Dictionary

Motif Detection:

New Protein Sequence → Motif Detector → Detection Results

[Narasimhan, Bu, Wang, Xu, Yang, Mathee, J Comput Biol, 2002]

# Helix-Turn-Helix Motifs

- Structure
  - 3-helix complex
  - Length: 22 amino acids
  - Turn angle

- Function
  - Gene regulation by binding to DNA

Branden & Tooze

**Figure 7.10** The helix-turn-helix motif in lambda Cro bound to DNA (orange) with the two recognition helices (red) of the Cro dimer sitting in the major groove of DNA. The binding model, suggested by Brian Matthews, is shown schematically in (a) with connected circles for the $C_\alpha$ positions as they were model built into regular B-DNA. A schematic diagram of the Cro dimer is shown in (b) with different colors for the two subunits. A schematic space-filling model of the dimer of Cro bound to a bent B-DNA molecule is shown in (c). The sugar-phosphate backbone of DNA is red, and the bases are yellow. Protein atoms are colored red, blue, green, and white. [(a) Adapted from D. Ohlendorf et al., *J. Mol. Evol.* 19: 113, 1983. (c) Courtesy of Brian Matthews.]

# HTH Motifs: Examples

| 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 |

# Combinatorial Method: GYM

❑ **Combinations of residues** in specific locations (may not be contiguous) contribute towards stabilizing a structure.

❑ Some reinforcing combinations are relatively rare.

❑ GYM algorithm is inspired by the APriori algorithm [Agrawal et al., 1996]

[Narasimhan, Bu, Wang, Xu, Yang,  Mathee, J Comput Biol, 2002]

# 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**
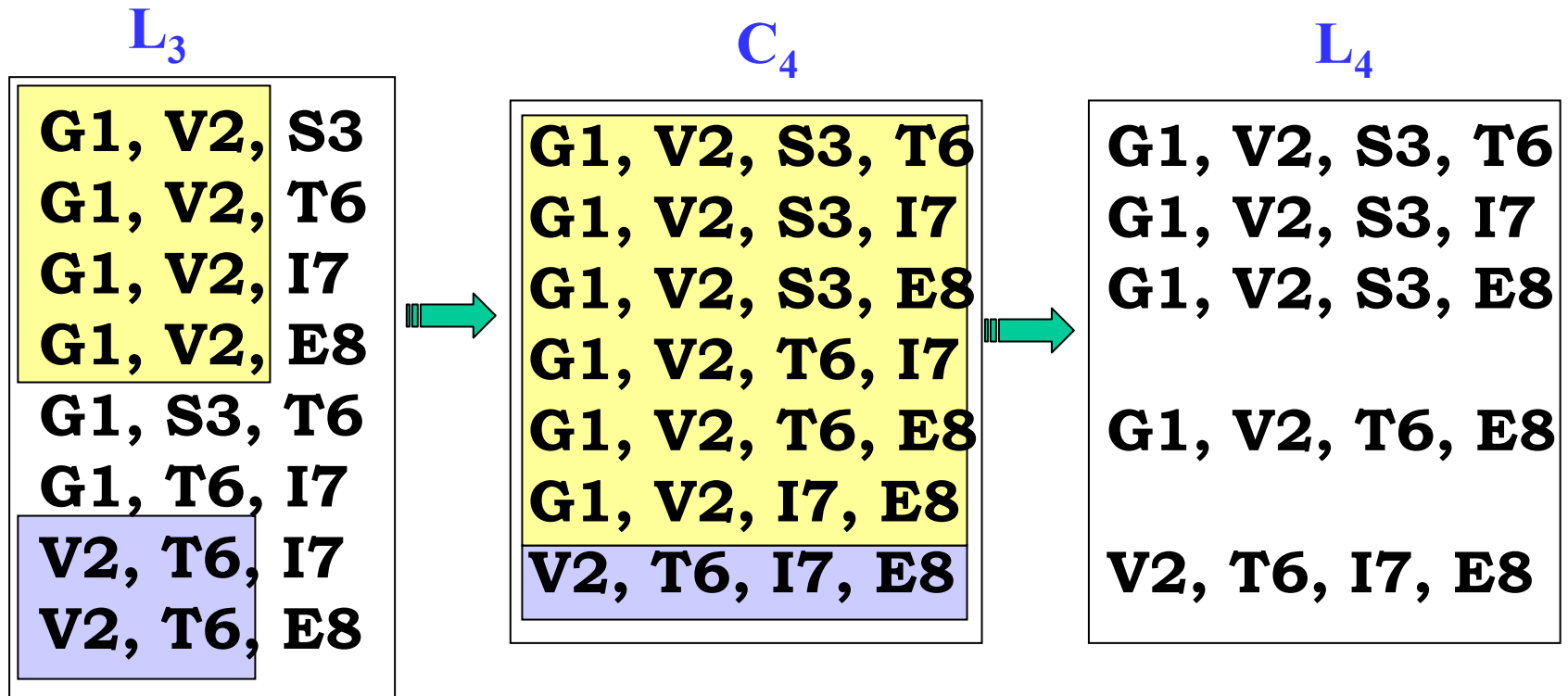
# Pattern Mining Algorithm

**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. $\quad$ $C_i$ **:= Candidates**$(L_{i-1})$
4. $\quad$ $L_i$ := Frequent candidates from $C_i$
5. $\quad$ **if** $(|L_i| <= 1)$ **then**
6. $\quad\quad$ **return** $L$ as the union of all $L_j$ , $j <= i$.

# **Candidates** Function

**L₃**

| | |
|---|---|
| **G1, V2,** | **S3** |
| **G1, V2,** | **T6** |
| **G1, V2,** | **I7** |
| **G1, V2,** | **E8** |
| **G1, S3,** | **T6** |
| **G1, T6,** | **I7** |
| **V2, T6,** | **I7** |
| **V2, T6,** | **E8** |

**C₄**

**G1, V2, S3, T6**
**G1, V2, S3, I7**
**G1, V2, S3, E8**
**G1, V2, T6, I7**
**G1, V2, T6, E8**
**G1, V2, I7, E8**
**V2, T6, I7, E8**

**L₄**

**G1, V2, S3, T6**
**G1, V2, S3, I7**
**G1, V2, S3, E8**

**G1, V2, T6, E8**

**V2, T6, I7, E8**

# Motif Detection Algorithm

**Algorithm Motif-Detection**

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

**Output** :      Detected motif(s).

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

| Motif | Protein Family | Number Tested | GYM = DE Agree | Number Annotated | GYM = Annot. |
|-------|---------------|---------------|----------------|------------------|--------------|
| HTH Motif (22) | 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 %) |

# Unaligned Pattern Discovery

**TEIRESIAS**:
The algorithm is similar to that used in GYM for aligned Pattern discovery.

| Protein Sequence Database | → | TEIRESIAS | → | **Seqlet Dictionary** |
|---|---|---|---|---|

**Seqlet Dictionary**

**A..GV**

**L..H...H**

**Y.C..C...F**

**V..G..G.G.T.L**
.
.
.

**Rigoutsos & Floratos**, Bioinformatics, '98

# TEIRESIAS: Key Features

- ❑ Starts with a set of __seed__ patterns (Enumeration step)
- ❑ __Convolution__ operator applied to all pairs of patterns:

$$A..GV.S \oplus V.S.GR = A..GV.S.GR$$

- ❑ __Order of Evaluation__ carefully chosen so that long patterns get longer first
- ❑ Finds **all maximal patterns.**
- ❑ **Combinatorial explosion** avoided by generating only relevant maximal patterns.

# SPLASH

- Structural Pattern Localization Analysis by Sequential Histogram (SPLASH)
- Not limited to fixed alphabet size
- Patterns are modeled by a homology metric and thus allow mismatches
- Early pruning of inconsistent seed patterns, leading to increased efficiency.
- Easily parallelized with availability of extra resources.

Califano, Bioinformatics, '00; Califano et al., J Comput Biol, '00

# Precomputed Sequence Patterns

- ❑ PROSITE
- ❑ BLOCKS and PRINTS
- ❑ *eMOTIF*
- ❑ SPAT
- ❑ PRODOM
- ❑ Pfam

# Motif Detection Tools

- ❑ PROSITE (Database of protein families & domains)
  - 🔴 Try PDOC00040. Also Try PS00041
- ❑ PRINTS Sample Output
- ❑ BLOCKS (multiply aligned ungapped segments for highly conserved regions of proteins; automatically created) Sample Output
- ❑ Pfam (Protein families database of alignments & HMMs)
  - 🔴 Multiple Alignment, domain architectures, species distribution, links: Try
- ❑ MoST
- ❑ PROBE
- ❑ ProDom
- ❑ DIP

# Protein Information Sites

❑SwissPROT & GenBank

❑InterPRO is a database of protein families, domains and functional sites in which identifiable features found in known proteins can be applied to unknown protein sequences. See sample.

❑PIR Sample Protein page

# Modular Nature of Proteins

❑ Proteins are collections of "modular" domains. For example,

**Coagulation Factor XII**

—[ F2 ]—[ E ]—[ F2 ]—[ E ]——[ K ]—[ Catalytic Domain ]—

—[ F2 ]—[ E ]——[ K ]——[ K ]—[ Catalytic Domain ]—

PLAT

# Domain Architecture Tools

❑ CDART
- 🔴 Protein `AAH24495`;  Domain Architecture;
- 🔴 It's domain relatives;
- 🔴 Multiple alignment for 2$^{nd}$ domain

❑ SMART

# Predicting Specialized Structures

- ❑ COILS – Predicts coiled coil motifs
- ❑ TMPred – predicts transmembrane regions
- ❑ SignalP – predicts signal peptides
- ❑ SEG – predicts nonglobular regions

# Patterns in DNA Sequences

❏ Signals in DNA sequence control events

- Start and end of genes
- Start and end of introns
- Transcription factor binding sites (regulatory elements)
- Ribosome binding sites

❏ Detection of these patterns are useful for

- Understanding gene structure
- Understanding gene regulation

# Motifs in DNA Sequences

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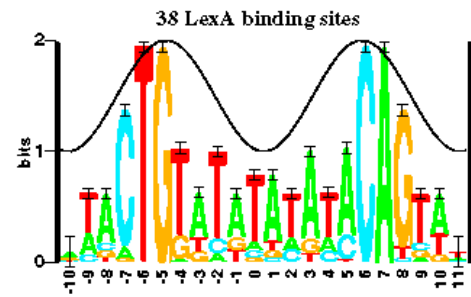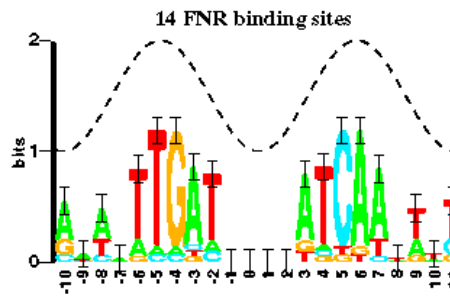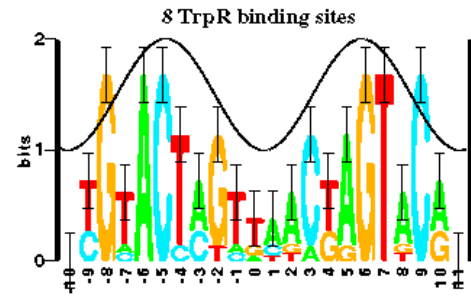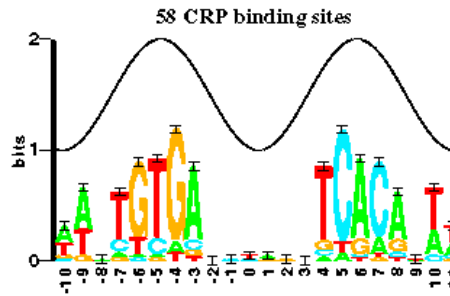
● Example:

http://www.lecb.ncifcrf.gov/~toms/sequencelogo.html

# Motifs



(a) CAP-DNA Complex

Helix-Turn-Helix

(b) CAP recognition site DNA Logo

(c) CAP Helix-Turn-Helix Logo

Sidechain-Base Interactions

# Motifs in DNA Sequences



Fig. 1. Some aligned sequences and their sequence logo. At the top of the figure are listed the 12 DNA sequences from the $P_L$ and $P_R$ control regions in bacteriophage lambda. These are bound by both the cI and cro proteins [16]. Each even numbered sequence is the complement of the preceding odd numbered sequence. The sequence logo, described in detail in the text, is at the bottom of the figure. The cosine wave is positioned to indicate that a minor groove faces the center of each symmetrical protein. Data which support this assignment are given in reference [17].
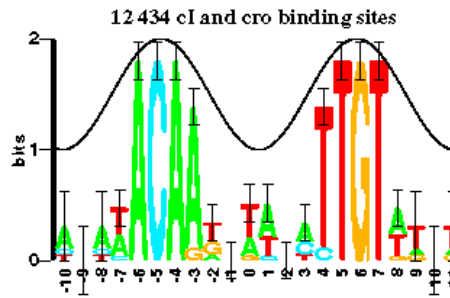
# More Motifs in *E. Coli* DNA Sequences

# E. coli Ribosome binding sites

http://www.lecb.ncifcrf.gov/~toms/sequencelogo.html

This figure shows two "sequence logos" which represent sequence conservation at the 5'(donor) and 3'(acceptor) ends of human introns. The region between the black vertical bars is removed during mRNA splicing. The logos graphically demonstrate that most of the pattern for locating the intron ends resides on the intron. This allows more codon choices in the protein-coding exons. The logos also show a common pattern "CAG|GT", which suggests that the mechanisms that recognize the two ends of the intron had a common ancestor. See R. M. Stephens and T. D. Schneider, "Features of spliceosome evolution and function inferred from an analysis of the information at human splice sites", J. Mol. Biol., 228, 1124-1136, (1992)

# Other Motifs in DNA Sequences: Human Splice Junctions
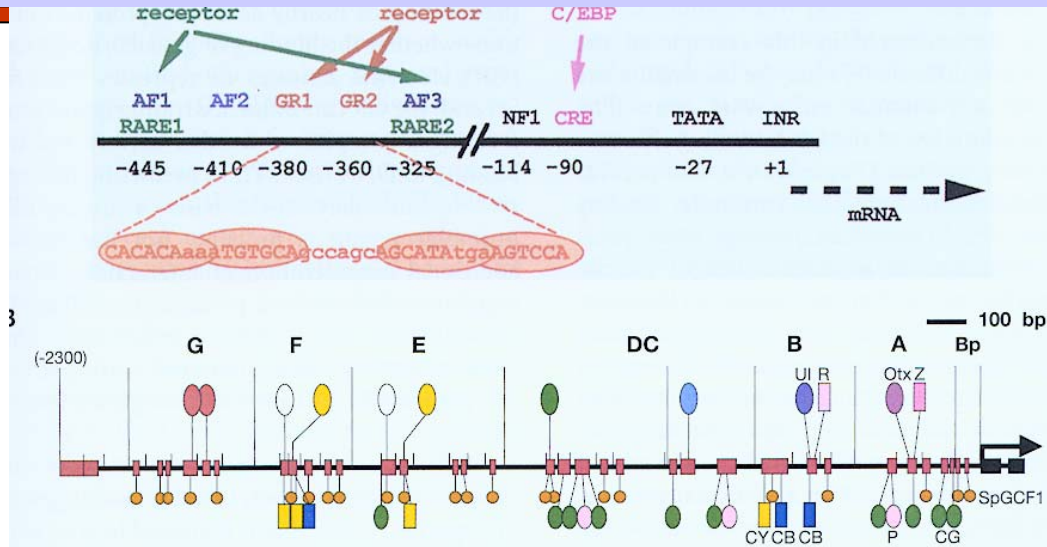
# Motifs in DNA Sequences



FIGURE 9.13. Regulatory elements of two promoters. (A) The rat *pepCK* gene. The relative positions of the TF-binding sites are illustrated (Yamada et al. 1999). The glucocorticoid response unit (GRU) includes three accessory factor–binding sites (AF1, AF2, and AF3), two glucocorticoid response elements (GR1 and GR2), and a cAMP response element (CRE). A dimer of glucocorticoid receptors bound to each GR element is depicted. The retinoic response unit (RAU) includes two retinoic acid response elements (RARE1 and RARE2) that coincide with the AF1 and AF3, respectively (Sugiyama et al. 1998). The sequences of the two GR sites and the binding of the receptor to these sites are shown. These sites deviate from the consensus sites and depend on their activity on accessory proteins bound to other sites in the GRU. This dependence on accessory proteins is reduced if a more consensus-like (canonical) GR element comprising the sequence TGTTCT is present. The CRE that binds factor C/EBP is also shown. (B) The 2300-bp promoter of the developmentally regulated gene *endo16* of the sea urchin (Bolouri and Davidson 2002). Different colors indicate different binding sites for distinct proteins and proteins shown above the line bind at unique locations, below the line at several locations. The regions A–G are functional modules that determine the expression of the gene in a particular tissue at a particular time of development and may either serve to induce transcription of the gene as a necessary developmental step (A, B, and G) or repress transcription (C–F) in tissues when it is not appropriate. (Reprinted, with permission, from Bolouri and Davidson 2002 [©2002 Elsevier].)

# 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, <u>Weeder</u>, YMF