

# Transcription Factor Binding Site motifs

Wstęp do biologii obliczeniowej

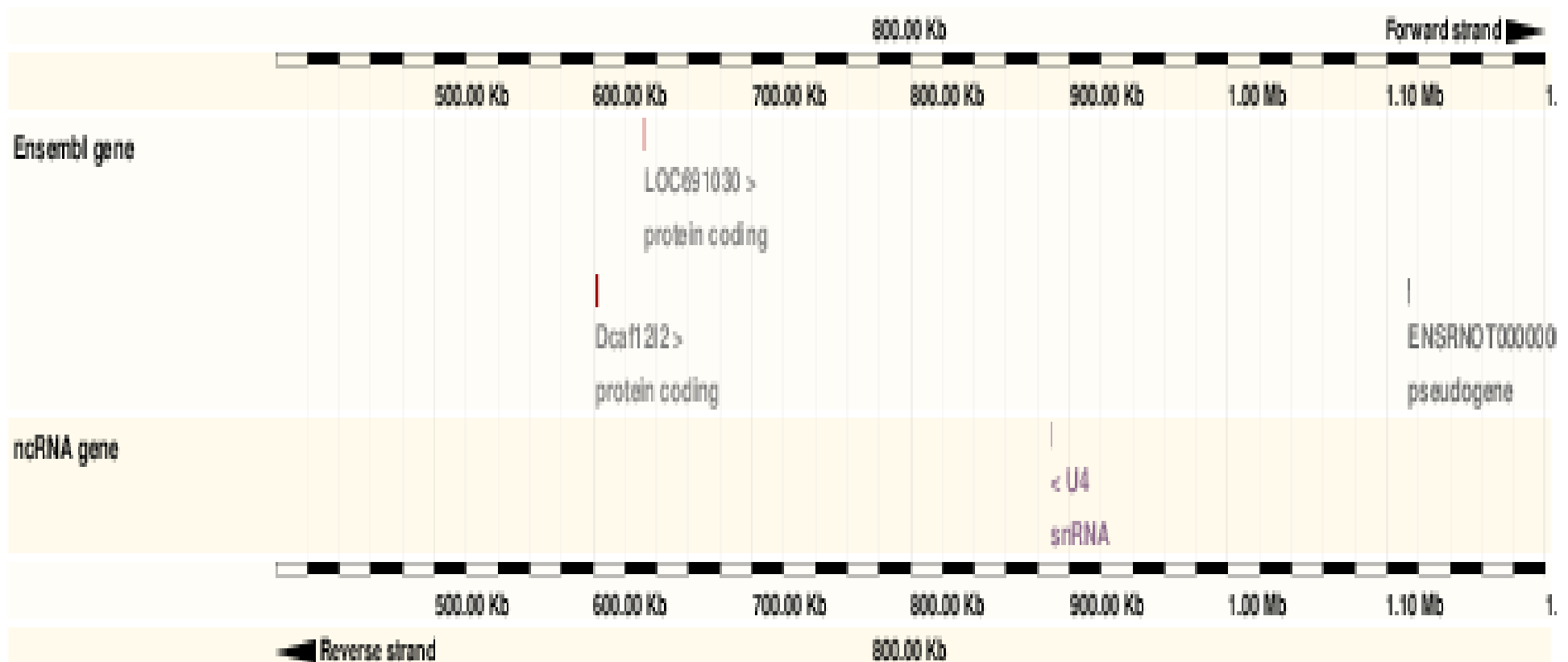
Lecture 11

May 12<sup>th</sup> 2020

# We know about genes, what's with the rest of the genome?

- We know that genomes have “genes” which are the “important” bits of inherited information
- We know that protein coding sequences are “genes”, as they are very important for cellular function
- Non-coding sequences include introns and intergenic sequences, and comprise large proportions of the genome

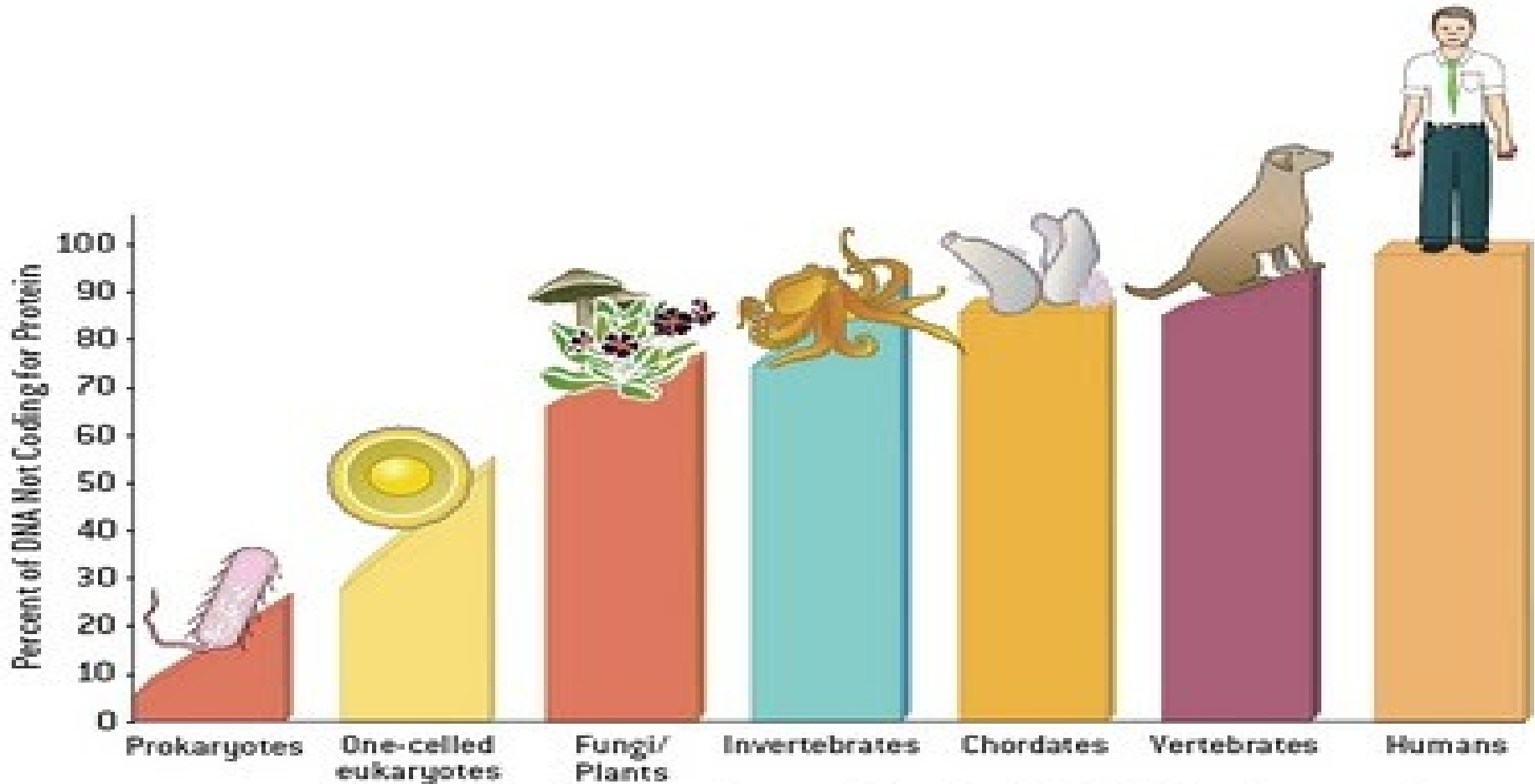
# Random genomic locus from Human genome...



There are currently 90 tracks turned off.

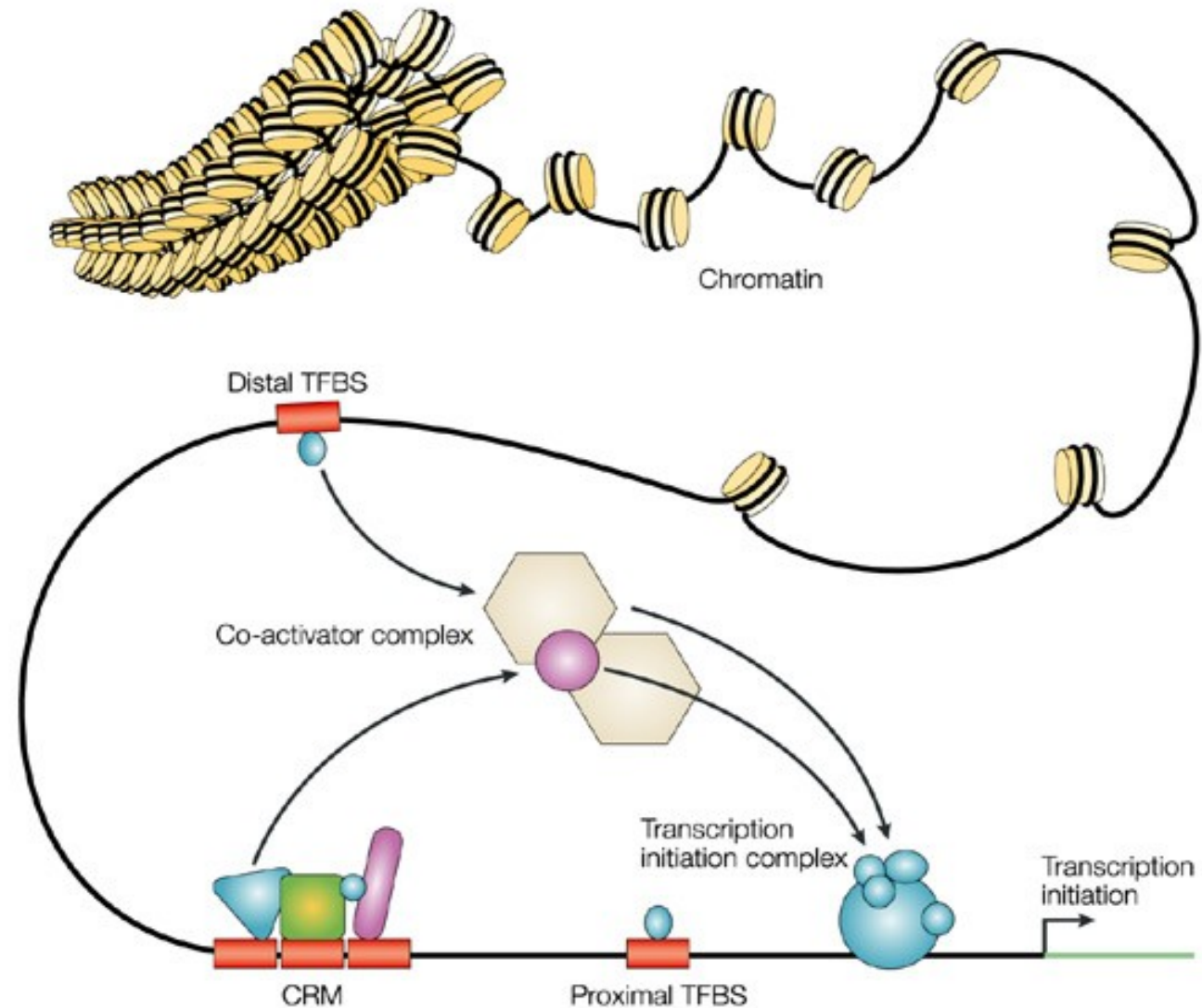
Ensembl Rattus norvegicus version 62.34d (RGSC3.4) Chromosome X: 400,004 - 1,200,007

# Is this “junk” DNA?

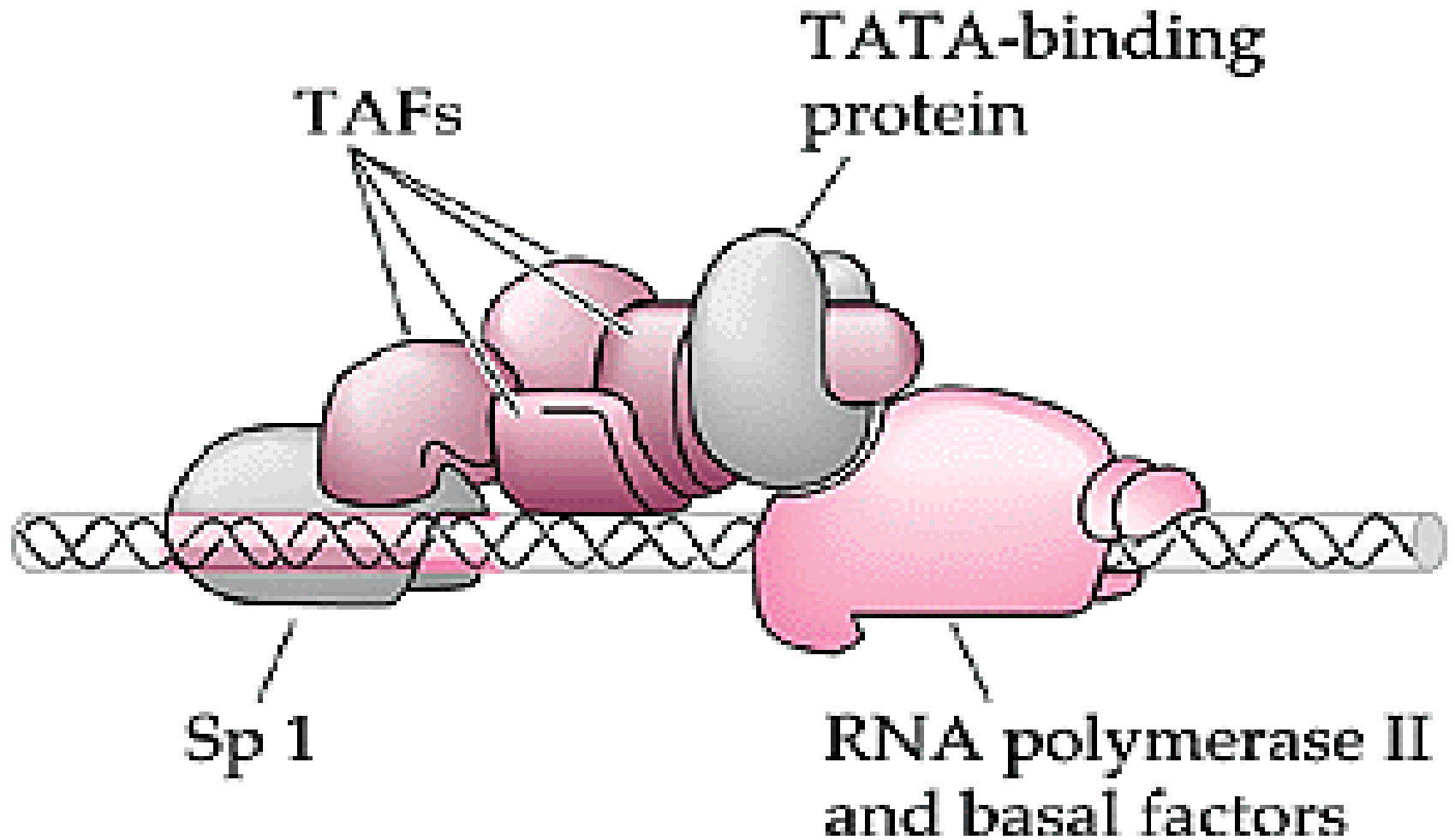


**NONPROTEIN-CODING SEQUENCES** make up only a small fraction of the DNA of prokaryotes. Among eukaryotes, as their complexity increases, generally so, too, does the proportion of their DNA that does not code for protein. The noncoding sequences have been considered junk, but perhaps it actually helps to explain organisms' complexity.

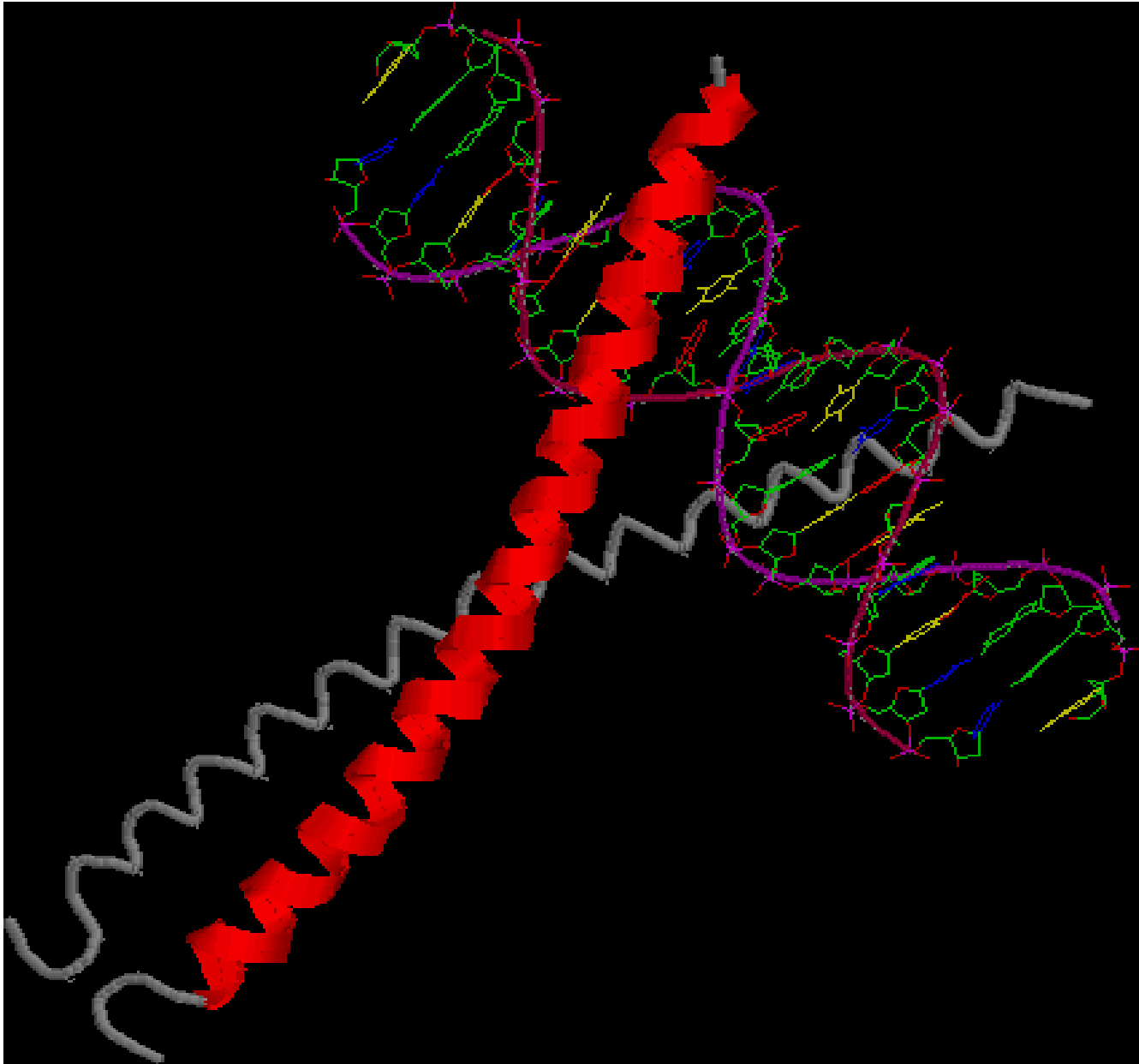
# Let us look at transcription initiation



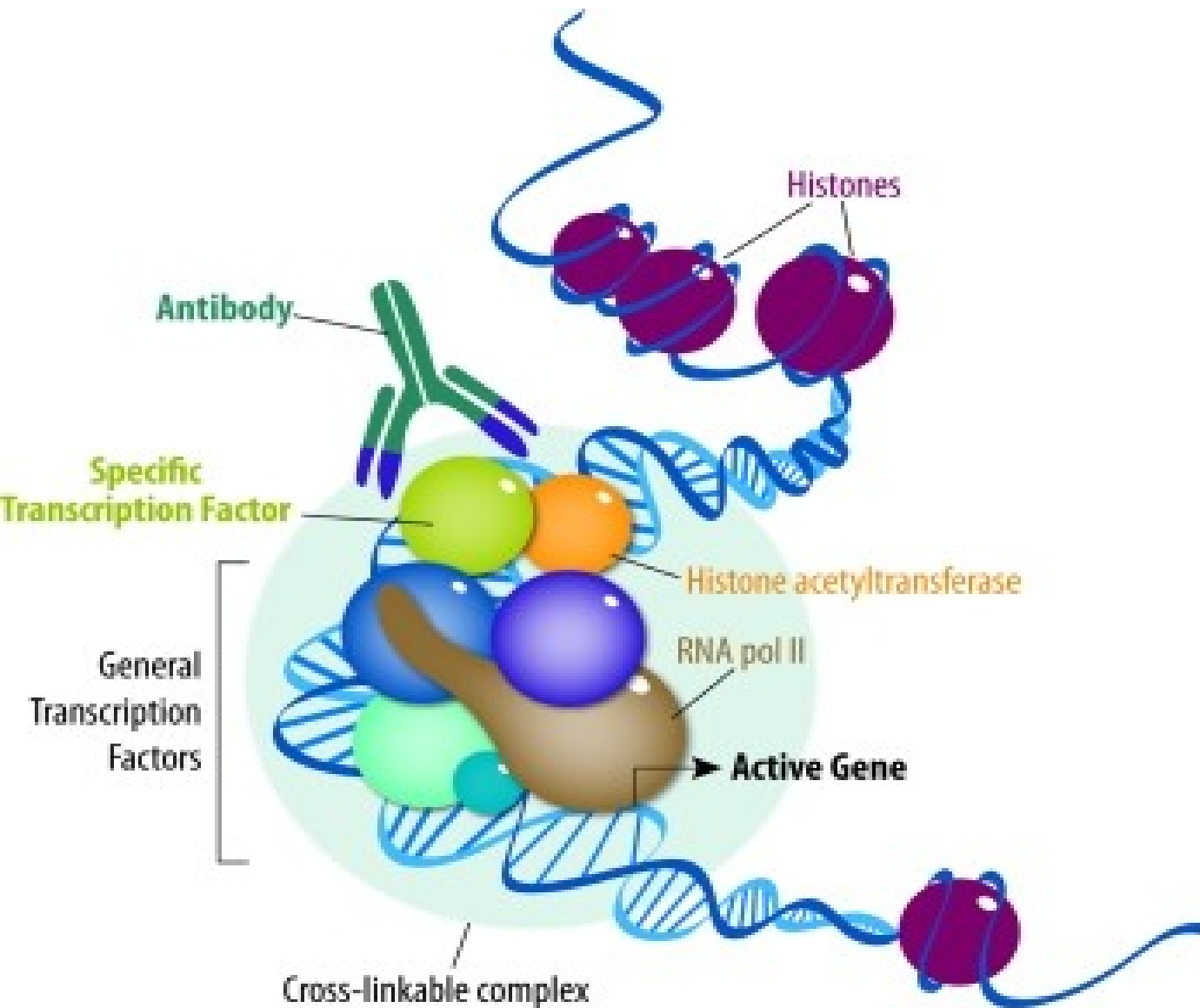
# How does a promoter look like



# Transcription factor binding



# Discovering binding sites





# TFBS representation

- Consensus words
- Regular expressions (IUPAC code)
- Position specific frequency Matrices (probability)
- Position specific Weight Matrices (log-odds)
- Higher Order Markov Models
- Bayesian Networks

# Position specific matrices

- We can represent a sequence alignment by position specific matrix, which retains only the counts of symbols and not the correlations
- We can also compute Position specific probability matrix, normalizing all columns to 1

```

cacagtCGCGTGt
actatCGCGTGtt
actgttCGCGTGc
tctcatCGCGTG
aggaatCGCGTGc
gagtgtCGCGTG
aggggtCGCGTG
  ggatCGCGTGtcc
    atgTGCCTGaagg
      ttaggTGCCTG
        TGCCTGccacctc
          gagtTGCCTGc
            aaggTGCCTGc
              aggtTGCCTGc
                tagcgTGCCTGc
                  tgattTGCCTG
                    atggaTGCCTGc
                      ggcTGCCTGacc
                        cttatTGCCTGc
                          tgggttAGCGTGc
                            ttacttAGCGTGc
                              gatccgGGCGTGa
                                tgggttaGTCGTG
                                  gtgtgaAGTTTGc

```

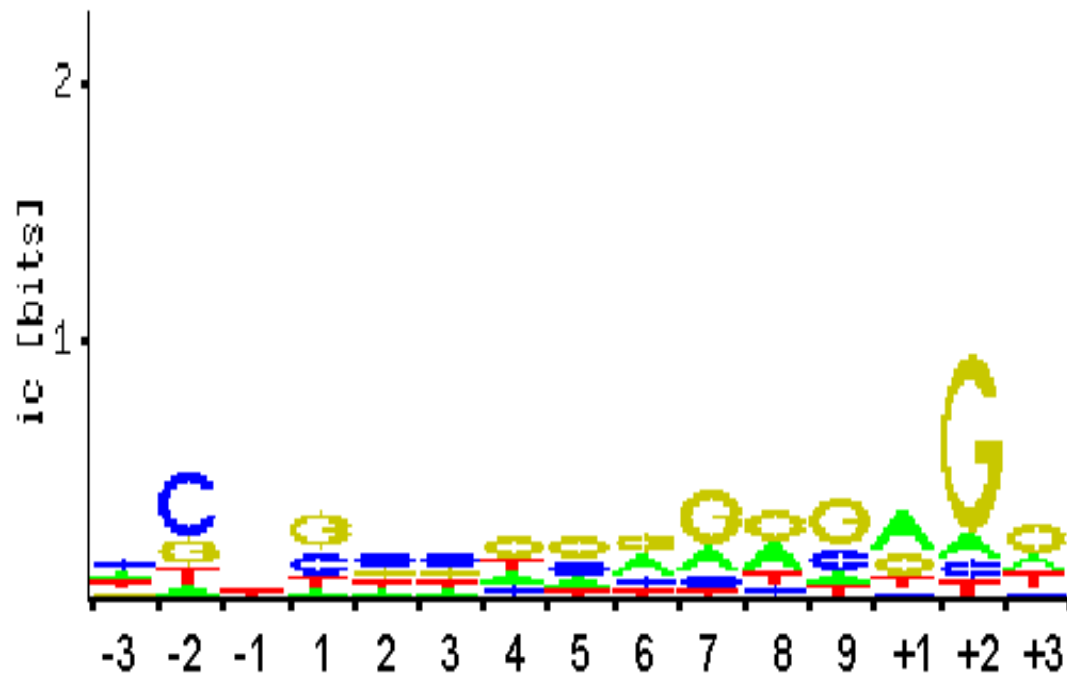
A[	3	0	0	0	0	0	]
C[	8	0	23	0	0	0	]
G[	2	23	0	23	0	24	]
T[	11	1	1	1	24	0	]



# Log-odds and PWMs

- PSPMs can be improved by including background information and converting to an additive (log) scale, this is usually called a position weight matrix (PWM)
- We can score new sequences against a PWM, by summing up scores, it corresponds to calculating log-odds of observing a sequence from PWM against random model
- It is also important to correct for missing values, by adding pseudo-counts

# We can find important positions



- Information content is defined as the entropy of the probability distribution
- Information content of a position in a PWM is the relative entropy of the motif compared to the background

<b>A:</b>	21	7	21	9	13	12	18	15	34	23	35	12	46	10	20
<b>C:</b>	39	58	21	32	40	37	14	27	18	14	10	23	7	7	10
<b>G:</b>	19	22	22	46	24	30	47	44	37	55	42	53	31	77	49
<b>T:</b>	20	12	33	12	21	19	19	12	9	6	10	9	14	5	19

# Predicting Binding sites in genomes

- Using PWMs, we can score new sequences
- If we want to predict binding sites, we need to decide on the threshold on the score to find predicted sites (above the threshold sequences)
- This is in fact a problem in space of all words of certain length  $L$
- We can estimate false positive/false negative rates of any threshold
- But the complexity grows exponentially with  $L$

# Motif databases

- There are two major eucaryotic TFBS motif databases: JASPAR and TRANSFAC
- TRANSFAC is commercial (BIOBASE GMBH)
- JASPAR is an open alternative
- They collect binding sites and matrices
- Provide tools for scanning sequences and motif comparisons

# Comparing motifs

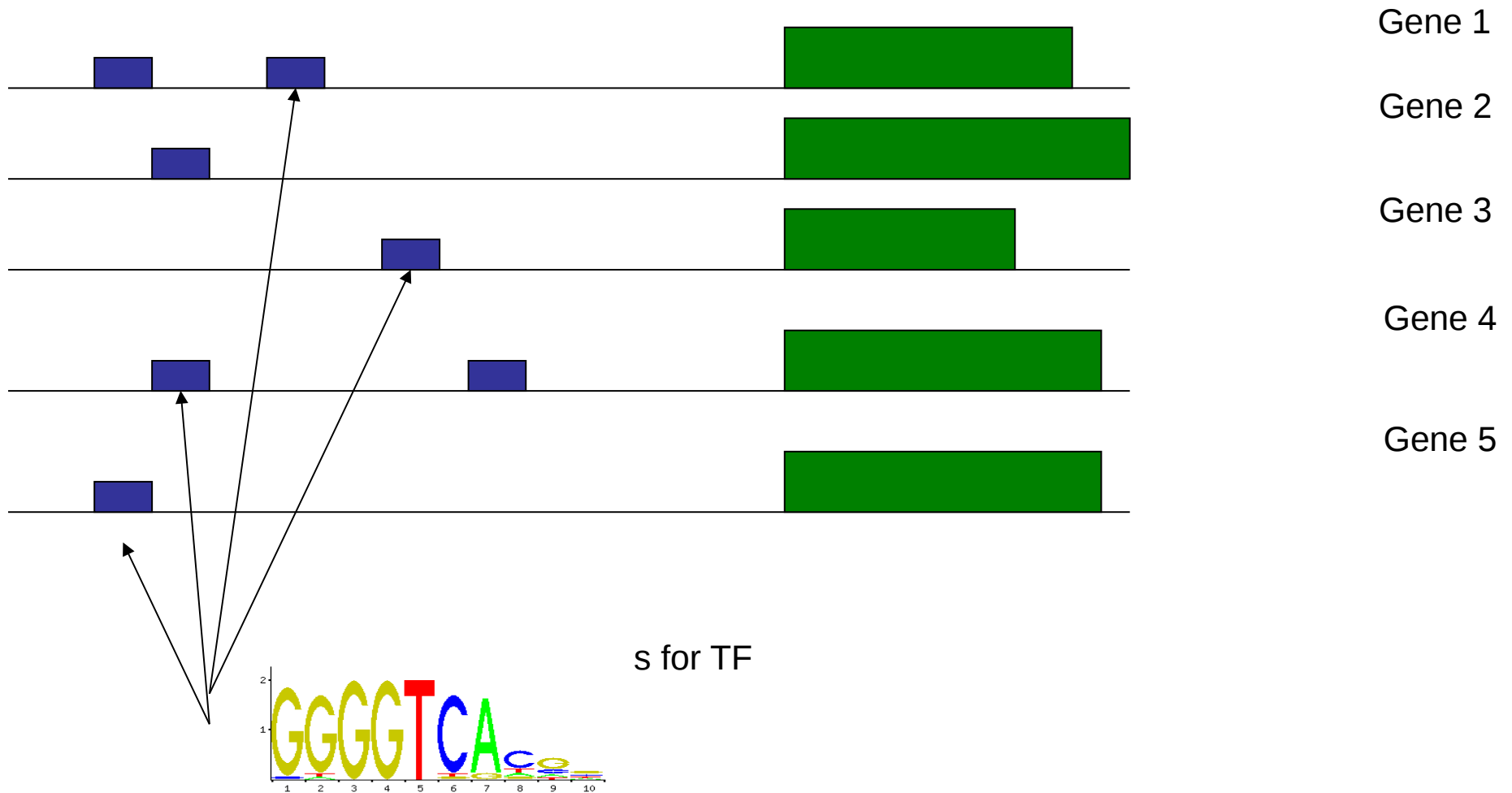
- It is important to be able to compare different motifs
- There are several approaches:
  - alignment based methods (like STAMP: gap-less alignment of PSSMs)
  - probability based (like MOSTA: probability of co-occurrence of the motifs in random sequences)
  - Distribution based: comparing motif score distributions across some sequences of interest (MEMOFinder)

# Motif finding – problem statement

- We are presented with a number of experimentally defined regions (several hundred bases in size) bound by the factor of interest
- What is the most likely sequence motif bound by this factor? (motif of a given length 5-20bp)
- Hint: what are the positions of the supposed binding sites of that motif in input sequences?



# Motif finding - illustration



# Terminology

- Information content of the motif

$$\sum_i \sum_j p_{ij} \log_2 \frac{p_{ij}}{b_j}$$

- log-odds of a motif occurrence

$$PWM(w) = \sum_j \log_2 \frac{p_{w[j]}}{b_j}$$

# Background models

- Different models
  - Uniform nucleotide distributions (weak)
  - GC-content (most commonly used )
  - Higher order Markov models (strongest, but complex and dependent on source sequences)
- Different types of sources
  - Whole genome
  - Non-coding sequences
  - Promoter-related
  - Input sequences

# Consensus method

CYCLE 1

sequence 1      sequence 2      sequence 3  
A C T G A      T A G C G      C T T G C

	A	C	T	G
A	1	0	0	0
C	0	1	0	0
G	0	0	0	1
T	0	0	1	0

$I_{seq} = 5.5$

CYCLE 2

	A	C	T	G
A	1	1	0	0
C	0	1	0	1
G	0	0	1	1
T	1	0	1	0

$I_{seq} = 2.8$

	A	C	T	G
A	2	0	0	0
C	0	1	1	0
G	0	1	0	2
T	0	0	1	0

$I_{seq} = 4.2$

	A	C	T	G
A	1	0	0	0
C	1	1	0	0
G	0	0	0	2
T	0	1	2	0

$I_{seq} = 4.2$

	A	C	T	G
A	1	0	0	0
C	0	1	0	1
G	0	0	1	1
T	1	1	1	0

$I_{seq} = 2.8$

CYCLE 3

	A	C	T	G
A	2	0	0	0
C	1	1	1	0
G	0	1	0	3
T	0	1	2	0

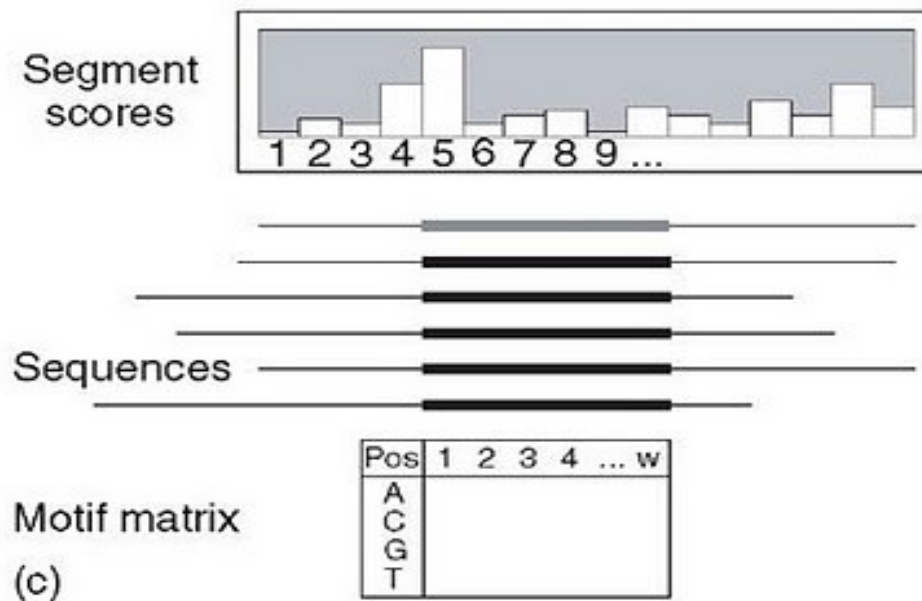
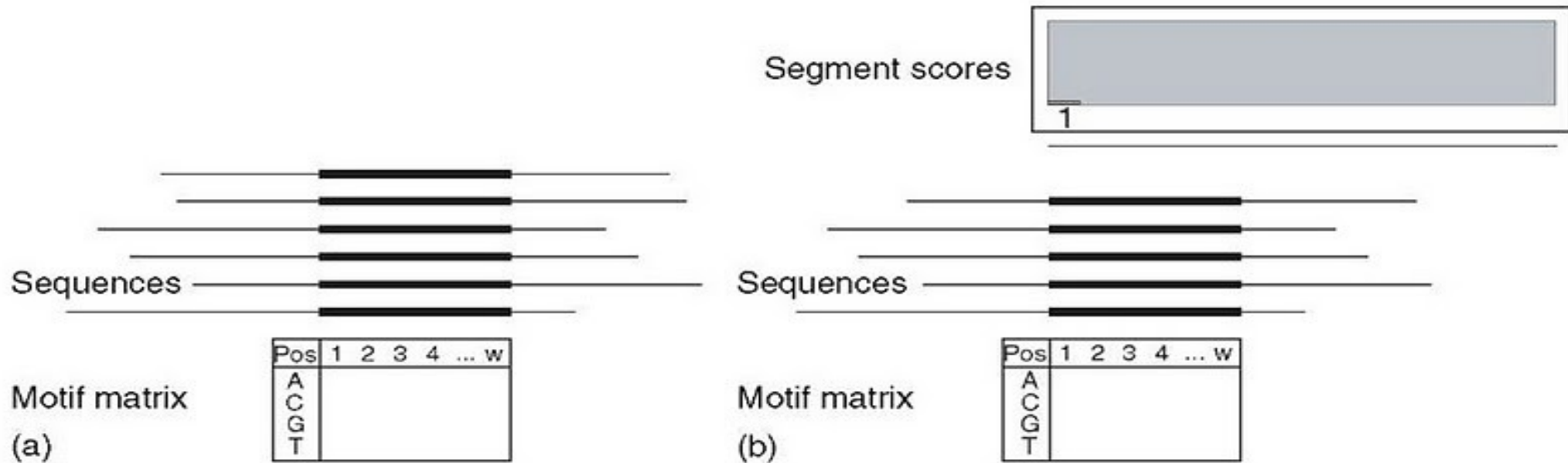
$I_{seq} = 3.2$

	A	C	T	G
A	2	0	0	0
C	0	1	1	1
G	0	1	1	2
T	1	1	1	0

$I_{seq} = 2.1$

- Greedy algorithm
- Dependant on sequence ordering
- Take only a few “most informant alignments

# Gibbs sampling



# Motif Expectation Maximization (MEME)

- $S$  unaligned set of sequences (training data)  $S_1, S_2, \dots, S_i, \dots, S_n$  each of length  $L$
- $W$  width of motif
- $Z$  matrix of probabilities that the motif starts in position  $j$  in  $S_i$
- $\rho$  matrix representing the probability of character  $c$  in column  $k$  (the character  $c$  will be A, C, G, or T for DNA sequences or one of the 20 protein characters)
- $\epsilon$  epsilon value

1. EM ( $S, W$ ) {
2.     choose starting point and initial value for  $\rho$
3.     do {
4.         re-estimate  $Z$  from  $\rho$  //the estimation step
5.         re-estimate  $\rho$  from  $Z$  //the maximization step
6.     } until (change in  $\rho < \epsilon$ )
7.     return  $\rho, Z$
8. }