# Finding Regulatory Motifs in DNA Sequences

03

#### Saeedeh Akbari

Department of Computer Engineering

Sharif University of Technology

Fall 2023

Adapted with modifications from lecture notes prepared by Phillip Compeau Bioinformatics Algorithms: An Active LearningApproach

### Outline

#### 03

#### **™** Motif finding problem

- □ Brute Force Motif Finding
- **Search** Trees
- Real Branch-and-Bound Motif Search
- Ranch-and-Bound Median String Search

### What is a DNA Motif?



- ❖ DNA motifs are short, recurring patterns that are presumed to have a biological function.
  - As regulatory motifs may vary at some positions
  - Every plant cell keeps track of day and night independently of other cells
    - e.g. three plant genes, called LCY, CCA1, and TOC1, are the clock's master timekeepers
    - Such regulatory genes, and the regulatory proteins that they encode, are often controlled by external factors (e.g., nutrient availability or sunlight) in order to allow organisms to adjust their gene expression

### Identifying Motifs: Complications



- We do not know the motif sequence
- We do not know where it is located relative to the genes start
- Motifs can differ slightly from one gene to the next
- Mow to discern it from "random" motifs?

### Two Methods



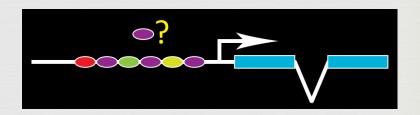
### Pattern Matching Finding known motifs

Does protein X bind upstream of my genes?

e.g. Patser, Pscan, Mast..

### Pattern Discovery Finding unknown motifs

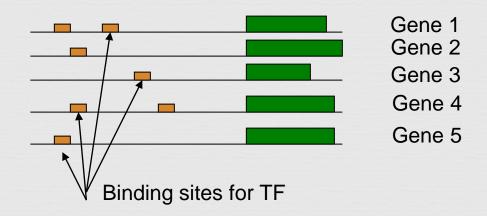
• What motifs are upstream of my genes?



e.g. MEME, Weeder, MDScan ...



- Suppose a transcription factor (TF) controls five different genes



#### 03

- Now suppose we are given the promoter regions of the five genes G1, G2, ... G5
- Can we find the binding sites of TF, without knowing about them *a priori*?
  - ☑ Binding sites are similar to each other, but not necessarily identical
- This is the motif finding problem

#### Identifying Motifs: Complications



- We do not know the motif sequence

  May know its length
- We do not know where it is located relative to the genes start
- Motifs can differ slightly from one gene to the next Non-essential bases could mutate...
- Mow to discern functional motifs from random ones?

### Outline

#### 03

- ™ Motif finding problem
- **™** Implanting Patterns in Random Text

- **Search** Trees
- Branch-and-Bound Motif Search
- Ranch-and-Bound Median String Search

### Random Sample



atgaccgggatactgataccgtatttggcctaggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg tgagtatccctgggatgacttttgggaacactatagtgctctcccgatttttgaatatgtaggatcattcgccagggtccga gctgagaattggatgaccttgtaagtgttttccacgcaatcgcgaaccaacgcggacccaaaggcaagaccgataaaggaga tcccttttgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatggcccacttagtccacttataggtcaatcatgttcttgtgaatggatttttaactgagggcatagaccgcttggcgcacccaaattcagtgtgggcgagcgcaa cggttttggcccttgttagaggcccccgtactgatggaaactttcaattatgagagagctaatctatcgcgtgctgttcataacttgagttggtttcgaaaatgctctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgtattggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatttcaacgtatgccgaaccgaaagggaag ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttctgggtactgatagca

### Implanting Motif AAAAAAAGGGGGGGG



atgaccgggatactgat<mark>AAAAAAAGGGGGGG</mark>ggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg acccctattttttgagcagatttagtgacctggaaaaaaatttgagtacaaaacttttccgaata<mark>AAAAAAAAGGGGGGG</mark>a tgagtatccctgggatgacttAAAAAAAAGGGGGGGtgctctcccgatttttgaatatgtaggatcattcgccagggtccga gctgagaattggatg<mark>AAAAAAAAGGGGGGG</mark>tccacgcaatcgcgaaccaacgcggacccaaaggcaagaccgataaaggaga tcccttttgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAAAAAAAAGGGGGGGCcttatag  $gtca at cat gttctt gt gaat g gat tt {\color{red} AAAAAAAAAGGGGGGG} gacc gctt g gc g cacc caa at t cag t g t g g g c g cacc caa at t cag t g t g g g c g cacc caa at t cag t g t g g g c g cacc caa at t cag t g t g g g c g cacc caa at t cag t g t g g g c g cacc caa at t cag t g t g g g c g cacc caa at t cag t g t g cacc g cacc caa at t cag t g t g cacc g cacc caa at t cag t g t g cacc g cacc caa at t cag t g cacc g cacc caa at t cag t g cacc g cacc caa at t cag t g cacc g cacc$ cggttttggcccttgttagaggcccccgtAAAAAAAAGGGGGGGCcaattatgagagagctaatctatcgcgtgcgtgttcat  $a acttg agtt \color{red} \color{blue} \color{blue}$ ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatAAAAAAAAGGGGGGGGaccgaaagggaag ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttAAAAAAAAGGGGGGGQa

## Implanting Motif **AAAAAAGGGGGGG**with Four Mutations



atgaccgggatactgat<mark>AgAAgAAAGGttGGG</mark>ggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg acccctattttttgagcagatttagtgacctggaaaaaaatttgagtacaaaacttttccgaata<mark>cAAtAAAAcGGcGGG</mark>a tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgatttttgaatatgtaggatcattcgccagggtccga gctgagaattggatg<mark>cAAAAAAAGGGattG</mark>tccacgcaatcgcgaaccaacgcggacccaaaggcaagaccgataaaggaga tcccttttgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaat| AtAAtAAAGGaaGGG cttatag  $gtcaatcatgttcttgtgaatggattt \color{red} \color{blue} \color$  $a acttg agtt \color{red} \color{blue} \color{blue}$ ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcat<mark>ActAAAAAGGaGcGG</mark>accgaaagggaag ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa

### Why Finding (15,4) Motif is Difficult?





### Challenge Problem



#### S Find a motif in a sample of

- 20 "random" sequences (e.g. 600 nt long)
- each sequence containing an implanted pattern of length 15,
- each pattern appearing with 4 mismatches as (15,4)-motif.

### Outline

#### 03

- **What is Motif**
- Motif finding problem
- **™** The Gold Bug Problem
- Rrute Force Motif Finding

- Branch-and-Bound Motif Search

### A Motif Finding Analogy

03

The Motif Finding Problem is similar to the problem posed by Edgar Allan Poe (1809 – 1849) in his *Gold Bug* story



### The Gold Bug Problem



#### **Given** a secret message:

```
53++!305))6*;4826)4+.)4+);806*;48!8`60))85;]8*:+*8!83(88)5*!;
46(;88*96*?;8)*+(;485);5*!2:*+(;4956*2(5*-4)8`8*; 4069285);)6
!8)4++;1(+9;48081;8:8+1;48!85;4)485!528806*81(+9;48;(88;4(+?348)4+;161;:188;+?;
```

## ©Decipher the message encrypted in the fragment

### Hints for The Gold Bug Problem

#### 03

#### Additional hints:

- The encrypted message is in English
- Each symbol correspond to one letter in the English alphabet
- No punctuation marks are encoded

#### The Gold Bug Problem: Symbol Counts



- Naive approach to solving the problem:
  - © Count the frequency of each symbol in the encrypted message
  - Find the frequency of each letter in the alphabet in the English language
  - Compare the frequencies of the previous steps, try to find a correlation and map the symbols to a letter in the alphabet

#### Symbol Frequencies in the Gold Bug Message



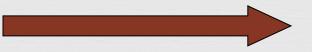
#### **∝**Gold Bug Message:

Symbol	8	•	4	)	+	*	5	6	(	!	1	0	2	9	3	:	?	•	-	]	
Frequency	34	25	19	16	15	14	12	11	9	8	7	6	5	5	4	4	3	2	1	1	1

#### **∝**English Language:

e ta o i n s r h l d c u m f p g w y b v k x j q z

Most frequent



Least frequent

#### The Gold Bug Message Decoding: First Attempt



⊗ By simply mapping the most frequent symbols to the most frequent letters of the alphabet:

sfiilfcsoorntaeuroaikoaiotecrntaeleyrcooestvenpinelefheeosnlt arhteenmrnwteonihtaesotsnlupnihtamsrnuhsnbaoeyentacrmuesotorl eoaiitdhimtaecedtepeidtaelestaoaeslsueecrnedhimtaetheetahiwfa taeoaitdrdtpdeetiwt

The result does not make sense

#### The Gold Bug Problem: *l*-tuple count



#### A better approach:

- Examine frequencies of *l*-tuples, combinations of 2 symbols, 3 symbols, etc.
- ";48" is the most frequent 3-tuple in English and text
- Make inferences of unknown symbols by examining other frequent *l*-tuples

#### The Gold Bug Problem: the ;48 clue



Mapping "the" to ";48" and substituting all occurrences of the symbols:

```
53++!305))6*the26)h+.)h+)te06*the!e`60))e5t]e*:+*e!e3(ee)5
*!t
h6(tee*96*?te)*+(the5)t5*!2:*+(th956*2(5*h)e`e*th0692e5)t)
6!e
)h++t1(+9the0e1te:e+1the!e5th)he5!52ee06*e1(+9thet(eeth(+?3hthe)h+t161t:1eet+?t
```

#### The Gold Bug Message Decoding: Second Attempt



#### Make inferences:

53++!305))6\*the26)h+.)h+)te06\*the!e`60))e5t]e\*:+\*e!e3(ee)5\*!th6(tee\*96\*?te)\*+(the5)t5\*!2:\*+(th956\*2(5\*h)e`e\*th0692e5)t)6!e)h++t1(+9the0e1te:e+1the!e5th)he5!52ee06\*e1(+9thet(eeth(+?3hthe)h+t161t:1eet+?t

- "thet(ee" most likely means "the tree"

  "Infer" (" = "r"
- "th(+?3h" becomes "thr+?3h"

  Can we guess "+" and "?"?

#### The Gold Bug Problem: The Solution



After figuring out all the mappings, the final message is:

AGOODGLASSINTHEBISHOPSHOSTELINTHEDEVILSSEATWENYONEDEGRE ESANDTHIRTEENMINUTESNORTHEASTANDBYNORTHMAINBRANCHSEVENT HLIMBEASTSIDESHOOTFROMTHELEFTEYEOFTHEDEATHSHEADABEELINE FROMTHETREETHROUGHTHESHOTFIFTYFEETOUT

### The Solution (cont'd)



#### Runctuation is important:

A GOOD GLASS IN THE BISHOP'S HOSTEL IN THE DEVIL'S SEA, TWENY ONE DEGREES AND THIRTEEN MINUTES NORTHEAST AND BY NORTH,

MAIN BRANCH SEVENTH LIMB, EAST SIDE, SHOOT FROM THE LEFT EYE OF

THE DEATH'S HEAD A BEE LINE FROM THE TREE THROUGH THE SHOT,

FIFTY FEET OUT.

### Solving the Gold Bug Problem



#### Rerequisites to solve the problem:

- Meed to know the relative frequencies of single letters, and combinations of two and three letters in English
- Knowledge of all the words in the English dictionary is
  highly desired to make accurate inferences

#### Motif Finding and The Gold Bug Problem: Similarities



#### **Motif Finding:**

- In order to solve the problem, we analyze the frequencies of patterns in the nucleotide sequences
- Knowledge of established regulatory motifs makes the Motif Finding problem simpler.

#### **™** Gold Bug Problem:

- In order to solve the problem, we analyze the frequencies of patterns in the text written in English
- Knowledge of the words in the English dictionary helps to solve the Gold Bug problem.

### Outline

#### 03

- Motif finding problem

- **™** The Motif Finding Problem
- Rrute Force Motif Finding
- Search Trees
- Branch-and-Bound Motif Search

#### Motif Finding and The Gold Bug Problem: Differences



#### Motif Finding is harder than Gold Bug problem:

- We don't have the complete dictionary of motifs
- The "genetic" language does not have a standard "grammar"
- Only a small fraction of nucleotide sequences encode for motifs; the size of data is enormous



Given a random sample of DNA sequences:

Find the pattern that is implanted in each of the individual sequences, namely, the motif

(cont'd)

#### **Additional** information:

The hidden sequence is of length 8

The pattern is not exactly the same in each array because random point mutations may occur in the sequences

(cont'd)

™ The patterns revealed with no mutations:

acgtacgt

Consensus String

(cont'd)

#### The patterns with 2 point mutations:

 $cctgatagacgctatctggctatcc\underline{aGgtacTt} aggtcctctgtgcgaatctatgcgtttccaaccat\\ agtactggtgtacatttgat\underline{CcAtacgt} acaccggcaacctgaaacaaacgctcagaaccagaagtgc\\ aa\underline{acgtTAgt} gcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt\\ agcctccgatgtaagtcatagctgtaactattacctgccacccctattacatctt\underline{acgtCcAt} ataca\\ ctgttatacaacgcgtcatggcggggtatgcgttttggtcgtcgtacgctcgatcgtta\underline{CcgtacgGc} \\ ctgttatacaacgcgtcatggcggggtatgcgttttggtcgtcgtacgctcgatcgtta\underline{CcgtacgGc} \\ ctgttatacaacgcgtcatggcggggtatgcgttttggtcgtcgtacgctcgatcgtta\underline{CcgtacgGc} \\ ctgttatacaacgcgtcatggcggggtatgcgttttggtcgtcgtacgctcgatcgtta\underline{CcgtacgGc} \\ ctgttatacaacgcgtcatggcggggtatgcgttttggtcgtcgtacgctcgatcgtta\underline{CcgtacgGc} \\ ctgttatacaacgcgtcatggcgggggtatgcgttttggtcgtcgtacgctcgatcgttaCcgtacgCc$ 

(cont'd)

#### The patterns with 2 point mutations:

cctgatagacgctatctggctatccaGgtacTtaggtcctctgtgcgaatctatgcgtttccaaccat agtactggtgtacatttgatCcAtacgtacaccggcaacctgaaacaaacgctcagaaccagaagtgcaaacgtTAgtgcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaattttagcctccgatgtaagtcatagctgtaactattacctgccacccctattacatcttacgtCcAtatacacctgttatacaacgcgtcatggcggggtatgcgttttggtcgtcgtacgctcgatcgttaCcgtacgGc

Can we still find the motif, now that we have 2 mutations?

### Defining Motifs



- To define a motif, lets say we know where the motif starts in the sequence
- The motif start positions in their sequences can be represented as  $s = (s_1, s_2, s_3, ..., s_t)$



### Motifs: Profiles and Consensus



C c A t a c g t

Alignment a c g t T A g t

a c g t C c A t

C c g t a c g G

aGgtacTt

$$\mathbf{s} = (s_1, s_2, ..., s_t)$$

- A 3 0 1 0 3 1 1 0

  Profile C 2 4 0 0 1 4 0 0

  G 0 1 4 0 0 0 3 1

  T 0 0 0 5 1 0 1 4
- Construct matrix profile with frequencies of each nucleotide in columns

Consensus A C G T A C G T

Consensus nucleotide in each position has the highest score in column

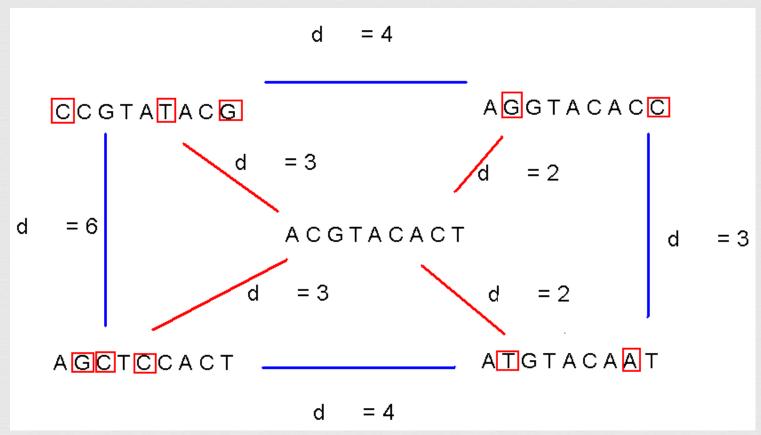
### Consensus



- Think of consensus as an "ancestor" motif, from which mutated motifs emerged
- The *distance* between a real motif and the consensus sequence is generally less than that for two real motifs

## Consensus (cont'd)





## Evaluating Motifs

03

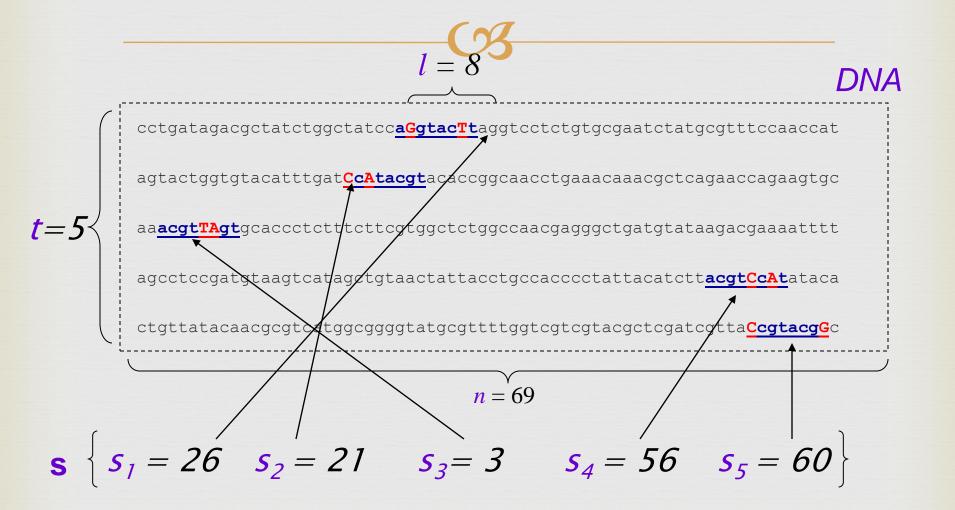
- We have a guess about the consensus sequence, but how "good" is this consensus?
- Need to introduce a scoring function to compare different guesses and choose the "best" one.

## Defining Some Terms



- $\bowtie t$  number of sample DNA sequences
- $\alpha n$  length of each DNA sequence
- $\bigcirc$  **DNA** sample of DNA sequences ( $t \times n$  array)
- $\bowtie s_i$  starting position of an  $\ell$ -mer in sequence i
- $\approx s = (s_1, s_2, ..., s_t)$  array of motif's starting positions

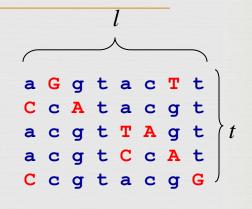
### Parameters



## Scoring Motifs

03

$$Score(s,DNA) = \sum_{i=1}^{l} \max_{k \in \{A,T,C,G\}} count(k,i)$$



```
A 3 0 1 0 3 1 1 0
C 2 4 0 0 1 4 0 0
G 0 1 4 0 0 0 3 1
T 0 0 0 5 1 0 1 4
```

Consensus acgtacgt

Score 3+4+4+5+3+4+3+4=30

# The Motif Finding Problem



- If starting positions  $\mathbf{s} = (s_1, s_2, \dots s_t)$  are given, finding consensus is easy even with mutations in the sequences because we can simply construct the profile to find the motif (consensus)
- But... the starting positions s are usually not given. How can we find the "best" profile matrix?

#### The Motif Finding Problem: Formulation



- Goal: Given a set of DNA sequences, find a set of L-mers, one from each sequence, that maximizes the consensus score
- Input: A *t* x *n* matrix of *DNA*, and *l*, the length of the pattern to find
- Output: An array of t starting positions  $\mathbf{s} = (s_1, s_2, \dots s_t)$  maximizing  $Score(\mathbf{s}, DNA)$

### The Matrix



- A position weight matrix (PWM)
  - also called position-specific weight matrix (PSWM)
  - also called position-frequency matrix (PFM)
  - also called position-specific scoring matrix (PSSM)
  - or just matrix
- There is a matrix element for all possible bases at every position.

	1	2	3	4	5	6	7	8	9	10	11
A	4	13	5	3	0	0	0	0	17	0	6
C	4	1	2	0	0	0	0	0	0	1	0
G	3	3	0	0	18	0	0	0	1	4	3
T	7	1	11	15	0	18	18	18	0	13	9

## Pattern Matching



#### Counts

A	4	13	5	3	0	0	0	0	17	0	6
С	4	1	2	0	0	0	0	0	0	1	0
G	3	3	0	0	18	0	0	0	1	4	3
T	7	1	11	15	0	18	18	18	0	13	9

Frequency

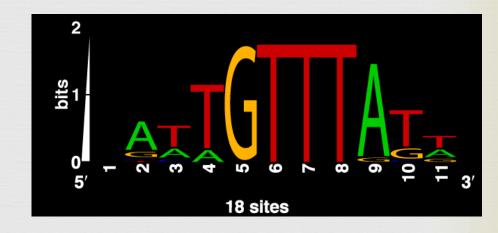
A	0.2	0.7	0.3	0.2	0.0	0.0	0.0	0.0	0.9	0.0	0.3
С	0.2	0.1	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.0
G	0.2	0.2	0.0	0.0	1.0	0.0	0.0	0.0	0.1	0.2	0.2
T	0.4	0.1	0.6	0.8	0.0	1.0	1.0	1.0	0.0	0.7	0.5

## Motif Logos



- A visual representation of the motif
- Each column of the matrix is represented as a stack of letters whose size is proportional to the corresponding residue frequency
- The total height of each column is proportional to its information content.

Α	4	13	5	3	0	0	0	0	17	0	6
С	4	1	2	0	0	0	0	0	0	1	0
G	3	3	0	0	18	0	0	0	1	4	3
			11								



T C G G t t g C G t A C a C G G t C A a G G A t т C C G G A T a Motifs G G A T C G G G G A т a T G G A T C G Т G G A C a T G G A C C a SCORE(Motifs) 3 + 4 + 0 + 0 + 1 + 1 + 1 + 5 + 2 + 3 + 6 + 4 = 30A: 0 0 0 6 COUNT(Motifs) 10 0 9 10 0 0 T: .9 .1 0 0 .6 PROFILE(Motifs) 0 .9 0 0 .5 0 .8 CONSENSUS (Motifs) C G G G C C T G A T т



### Outline

### 03

- Motif finding problem

- **Representation** Representation of the Brute Force Motif Finding
- Search Trees
- Branch-and-Bound Motif Search
- Ranch-and-Bound Median String Search

#### The Motif Finding Problem: Brute Force Solution



- Compute the scores for each possible combination of starting positions s
- The best score will determine the best profile and the consensus pattern in *DNA*
- The goal is to maximize Score(s,DNA) by varying the starting positions  $s_i$ , where:

$$s_i = [1, ..., n-l+1]$$
  
 $i = [1, ..., t]$ 

#### BruteForceMotifSearch



- 1. BruteForceMotifSearch(DNA, t, n, l)
- **2.** bestScore ← 0
- 3. **for** each  $\mathbf{s} = (s_1, s_2, \dots, s_t)$  from  $(1, 1, \dots, 1)$  to  $(n-l+1, \dots, n-l+1)$
- 4. if (Score(s, DNA) > bestScore)
- 5.  $bestScore \leftarrow score(s, DNA)$
- 6. bestMotif  $\leftarrow (s_1, s_2, \ldots, s_t)$
- 7. return bestMotif

### Running Time of BruteForceMotifSearch



- Varying (n l + 1) positions in each of t sequences, we're looking at  $(n l + 1)^t$  sets of starting positions
- For each set of starting positions, the scoring function makes  $\ell$  operations, so complexity is  $\ell(n \ell + 1)^t = O(\ell n^t)$
- That means that for t = 8, n = 1000, l = 10 we must perform approximately  $10^{20}$  computations it will take billions years

### Outline

#### 03

- Motif finding problem

- **™** Brute Force Motif Finding
- **™** The Median String Problem
- Search Trees
- Real Branch-and-Bound Motif Search
- Ranch-and-Bound Median String Search

# The Median String Problem



Given a set of *t* DNA sequences find a pattern that appears in all *t* sequences with the minimum number of mutations

This pattern will be the motif

## Hamming Distance

03

(v,w) is the number of nucleotide pairs that do not match when v and w are aligned. For example:

 $d_H(AAAAAA,ACAAAC) = 2$ 

## Total Distance: Definition



Governormal For each DNA sequence i, compute all  $d_H(v, x)$ , where x is an i-mer with starting position  $s_i$ 

$$(1 \le s_i \le n - \ell + 1)$$

- S Find minimum of  $d_H(v, x)$  among all  $\ell$ -mers in sequence i
- ™ TotalDistance(v,DNA) is the sum of the minimum Hamming distances for each DNA sequence i
- TotalDistance(v,DNA) =  $min_s d_H(v, s)$ , where s is the set of starting positions  $s_1, s_2, ... s_t$

## Total Distance: An Example

03

 $\bowtie$  Given v ="acgtacgt" and s

```
d_{H}(v,\,x)=0 acgtacgt cottgates acgtacgt acgtacgt acgtacgt acgtact acgtacgt acgt
```

*v* is the sequence in red, *x* is the sequence in blue

TotalDistance(v,DNA) = 0

## Total Distance: Example



 $\bowtie$  Given v ="acgtacgt" and s

*v* is the sequence in red, *x* is the sequence in blue

$$\bowtie$$
 *TotalDistance*(*v*,*DNA*) = 1+0+2+0+1 = 4

#### The Median String Problem: Formulation



- $\bigcirc$  Input: A  $t \times n$  matrix DNA, and  $\mathcal{L}$ , the length of the pattern to find
- Output: A string v of l nucleotides that minimizes TotalDistance(v,DNA) over all strings of that length

## Median String Search Algorithm

#### CS

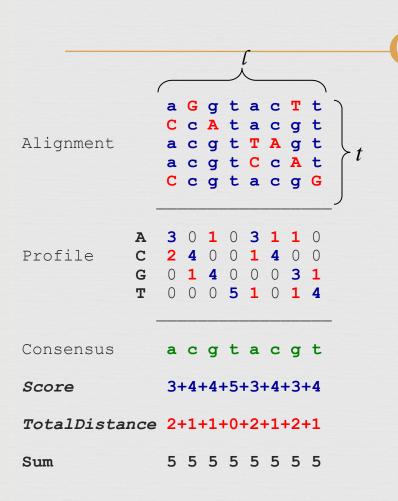
- 1. MedianStringSearch (DNA, t, n, l)
- 2.  $bestWord \leftarrow AAA...A$
- 3.  $bestDistance \leftarrow \infty$
- 4. for each *l*-mer s from AAA...A to TTT...T if TotalDistance(s,DNA) < bestDistance</p>
- 5.  $bestDistance \leftarrow TotalDistance(s,DNA)$
- 6. bestWord  $\leftarrow$  s
- 7. return bestWord

#### Motif Finding Problem == Median String Problem



- The *Motif Finding* is a maximization problem while *Median String* is a minimization problem
- However, the *Motif Finding* problem and *Median String* problem are computationally equivalent
- Need to show that minimizing *TotalDistance* is equivalent to maximizing *Score*

### We are looking for the same thing



```
At any column i
Score_i + TotalDistance_i = t
```

- Because there are  $\ell$  columns  $Score + TotalDistance = \ell * t$
- Rearranging:
  Score = ℓ \* t TotalDistance
- of the right side is equivalent to the maximization of the left side

### Motif Finding Problem vs. Median String Problem

- Why bother reformulating the Motif Finding problem into the Median String problem?
  - The Motif Finding Problem needs to examine all the combinations for **s**. That is  $(n l + 1)^t$  combinations!!!
  - The Median String Problem needs to examine all  $4^{\ell}$  combinations for v. This number is relatively smaller

### Motif Finding: Improving the Running Time

#### 03

#### Recall the BruteForceMotifSearch:

```
1. BruteForceMotifSearch (DNA, t, n, l)
2. bestScore \leftarrow 0
3. for each \mathbf{s} = (s_1, s_2, \ldots, s_t) from (1, 1, \ldots, 1) to (n - l + 1, \ldots, n - l + 1)
4. if (Score(\mathbf{s}, DNA) > bestScore)
5. bestScore \leftarrow Score(\mathbf{s}, DNA)
6. bestMotif \leftarrow (s_1, s_2, \ldots, s_t)
7. return bestMotif
```

### Median String: Improving the Running Time

#### 03

- 1. MedianStringSearch (DNA, t, n, l)
- 2.  $bestWord \leftarrow AAA...A$
- 3.  $bestDistance \leftarrow \infty$
- 4. **for** each *l*-mer **s from** AAA...A to TTT...T **if** *TotalDistance*(**s**,*DNA*) < *bestDistance*
- 5.  $bestDistance \leftarrow TotalDistance(s,DNA)$
- 6. bestWord  $\leftarrow$  s
- 7. return bestWord

# Structuring the Search

03

For the Median String Problem we need to consider all 4<sup>l</sup> possible *l*-mers:

aa... aa

aa... ac

aa... ag

aa... at

.

tt... tt

How to organize this search?

### Alternative Representation of the Search Space



$$\triangle$$
 Let **A** = 1, **C** = 2, **G** = 3, **T** = 4

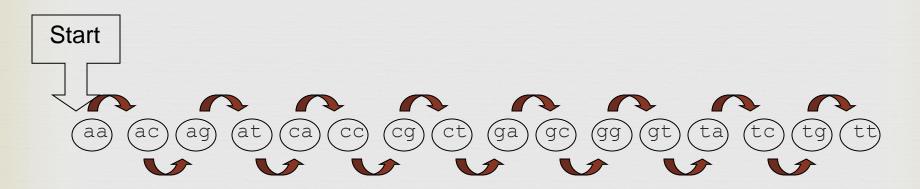
Then the sequences from AA...A to TT...T become:

№ Notice that the sequences above simply list all numbers as if we were counting on base 4 without using 0 as a digit

### Linked List



 $\bowtie$  Suppose l = 2



Need to visit all the predecessors of a sequence before visiting the sequence itself

## Linked List (cont'd)



- Call Let's try grouping the sequences by their prefixes



### Outline

### 03

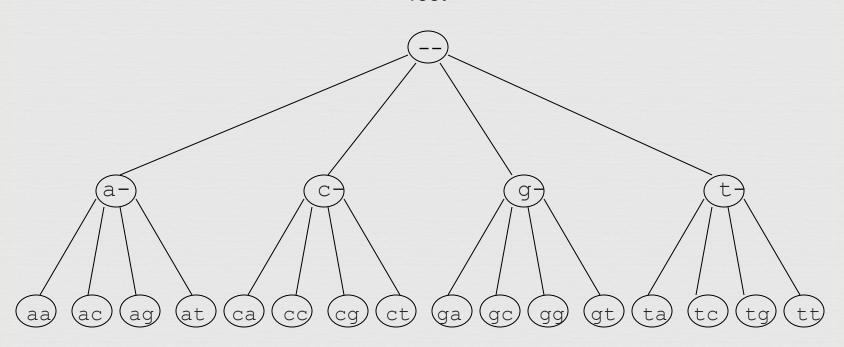
- **What is Motif**
- Motif finding problem

- Rrute Force Motif Finding
- **Search Trees**
- Branch-and-Bound Motif Search
- Ranch-and-Bound Median String Search

## Search Tree







## Analyzing Search Trees



- Characteristics of the search trees:
  - The sequences are contained in its leaves
  - The parent of a node is the prefix of its children
- How can we move through the tree?

### Moving through the Search Trees



- Four common moves in a search tree that we are about to explore:
  - Move to the next leaf
  - S Visit all the leaves
  - S Visit the next node
  - **3** Bypass the children of a node

### Visit the Next Leaf

Given a current leaf a, we need to compute the "next" leaf:

```
NextLeaf(\mathbf{a}, L, \mathbf{k})
                                   // a : the array of digits
2. for i \leftarrow 1 to L
                                   // L: length of the array
    if a_i < k
           a_i \leftarrow a_i + 1
           return a
           a_i \leftarrow 1
7. return a
```

```
// k: max digit value
```

## NextLeaf (cont'd)

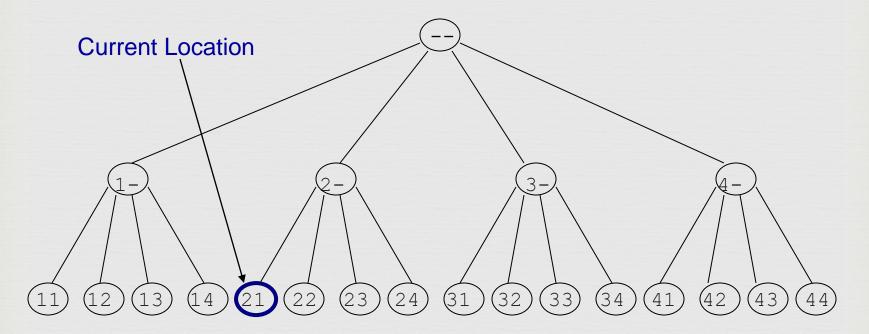


- $\bowtie$  The algorithm is common addition in radix k:
- "Carry the one" to the next digit position when the digit is at maximal value

## NextLeaf: Example

03

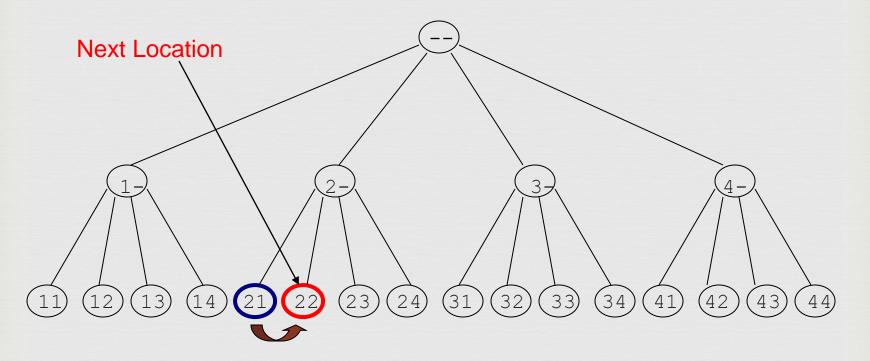
Moving to the next leaf:



## NextLeaf: Example (cont'd)



Moving to the next leaf:



### Visit All Leaves



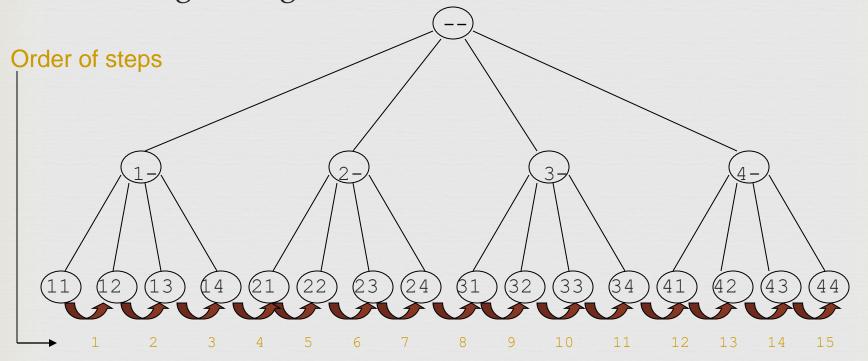
Rrinting all permutations in ascending order:

```
    AllLeaves(L,k) // L: length of the sequence
    a ← (1,...,1) // k: max digit value
    while forever // a: array of digits
    output a
    a ← NextLeaf(a,L,k)
    if a = (1,...,1)
    return
```

## Visit All Leaves: Example



Moving through all the leaves in order:



## Depth First Search



So we can search leaves

Mow about searching all vertices of the tree?

We can do this with a *depth first* search

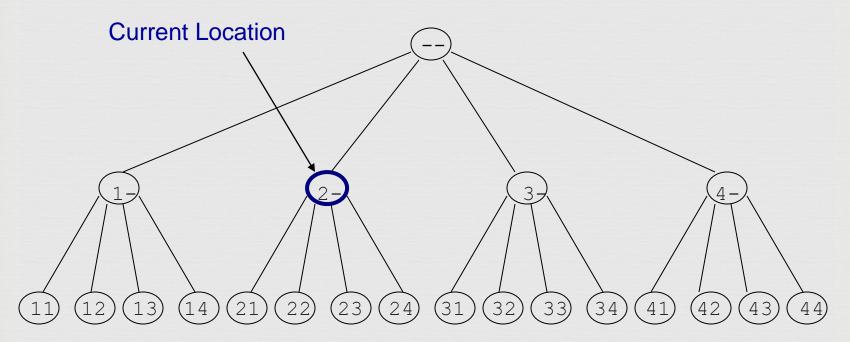
### Visit the Next Vertex

```
NextVertex(\mathbf{a}, \mathbf{i}, \mathbf{L}, \mathbf{k})
                               // a : the array of digits
    if i < /
                                 // i : prefix length
3.
                             // L: max length
     a_{i+1} \leftarrow 1
        return (a, i+1) // k: max digit value
    else
6. for j \leftarrow L to 1
         if a_i < k
           a_i \leftarrow a_i + 1
8.
           return(a,j)
9.
10. \text{ return}(a,0)
```

## Example

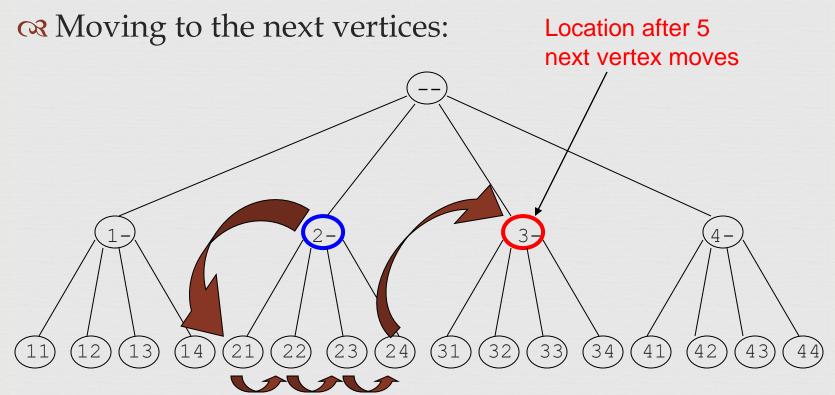
03

Moving to the next vertex:



## Example





## Bypass Move

#### 03

Given a prefix (internal vertex), find next vertex after skipping all its children

```
1. <u>Bypass(a, i, L, k)</u> // a: array of digits

2. for j \leftarrow 1 to i // i: prefix length

3. if a_j < k // L: maximum length

4. a_j \leftarrow a_j + 1 // k: max digit value

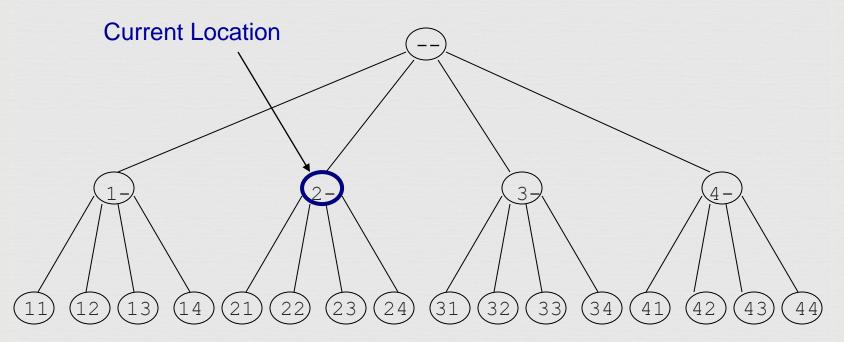
5. return(a, j)

6. return(a, 0)
```

## Bypass Move: Example



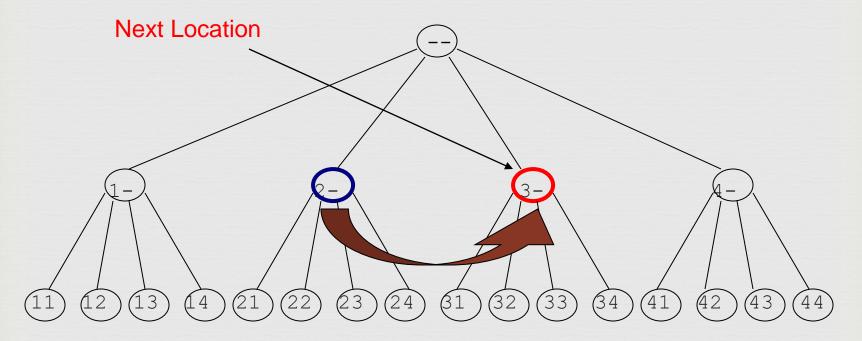
○ Bypassing the descendants of "2-":



## Example

03

○ Bypassing the descendants of "2-":



## Revisiting Brute Force Search



Now that we have method for navigating the tree, lets look again at BruteForceMotifSearch

## Brute Force Search Again

```
BruteForceMotifSearchAgain (DNA, t, n, l)
2.
       s \leftarrow (1,1,...,1)
3.
       bestScore ← Score(s, DNA)
       while forever
4.
5.
              s \leftarrow \text{NextLeaf}(s, t, n-l+1)
6.
               if (Score(s, DNA) > bestScore)
                      bestScore ← Score(s, DNA)
                      bestMotif \leftarrow (s_1, s_2, \ldots, s_t)
8.
       return bestMotif
9.
```

### Can We Do Better?

#### 03

- Sets of  $\mathbf{s}=(s_1, s_2, ..., s_t)$  may have a weak profile for the first i positions  $(s_1, s_2, ..., s_i)$
- Optimism: if all subsequent (t-i) positions  $(s_{i+1}, ...s_t)$  add

(t-i)\* l to Score(s, i, DNA)

If Score(s,i,DNA) + (t-i) \* l < BestScore, it makes no sense to search in vertices of the current subtree Use ByPass()

### Outline

- **What is Motif**
- Motif finding problem

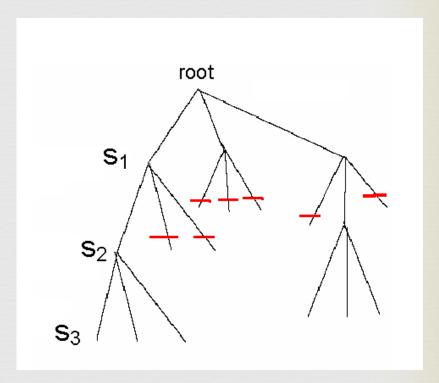
- Rrute Force Motif Finding

- **™** Branch-and-Bound Motif Search
- Ranch-and-Bound Median String Search

#### Branch and Bound Algorithm for Motif Search



- Since each level of the tree goes deeper into search, discarding a prefix discards all following branches
- This saves us from looking at  $(n l + 1)^{t-i}$  leaves
  - Use NextVertex() and ByPass() to navigate the tree



#### Pseudocode for Branch and Bound Motif Search

```
BranchAndBoundMotifSearch(DNA, t, n, l)
1. 2. 3. 4. 5.
            s \leftarrow (1,...,1)
            best$core ← 0
            i \leftarrow 1
            while i > 0
6.
7.
8.
9.
                   if i < t
                                      optimisticScore \leftarrow Score(s, i, DNA) +(t - i) * \ell if optimisticScore < bestScore (s, i) \leftarrow Bypass(s,i, n-\ell+1)
10.
                                       else
11.
                                           (s, i) \leftarrow \text{NextVertex}(s, i, n-\ell+1)
12.
                   else
                                    if Score(s,DNA) > bestScore

bestScore \leftarrow Score(s)

bestMotif \leftarrow (s_1, s_2, s_3, ..., s_t)

(s,i) \leftarrow NextVertex(s,i,t,n-l+1)
13.
14.
15.
16.
17.
            return bestMotif
```

### Median String Search Improvements



- Recall the computational differences between motif search and median string search
  - The Motif Finding Problem needs to examine all  $(n-l+1)^t$  combinations for **S**.
  - The Median String Problem needs to examine  $4^{\ell}$  combinations of v. This number is relatively small
- We want to use median string algorithm with the Branch and Bound trick!

### Outline

- **What is Motif**
- Motif finding problem

- Rrute Force Motif Finding

- Branch-and-Bound Motif Search
- **™** Branch-and-Bound Median String Search

# Branch and Bound Applied to Median String Search

Note that if the total distance for a prefix is greater than that for the best word so far:

TotalDistance (*prefix*, *DNA*) > *BestDistance* 

there is no use exploring the remaining part of the word

 $\bowtie$  We can eliminate that branch and BYPASS exploring that branch further

### Bounded Median String Search

```
BranchAndBoundMedianStringSearch(DNA, t, n, l)
1.
2.
      s \leftarrow (1,...,1)
3.
      bestDistance ← ∞
4
     i \leftarrow 1
     while i > 0
6.
        if i < l
7.
           prefix \leftarrow string corresponding to the first i nucleotides of s
8.
           optimisticDistance ← TotalDistance(prefix,DNA)
9.
           if optimisticDistance > bestDistance
10.
                         (s, i) \leftarrow Bypass(s, i, l, 4)
11.
            else
12.
                         (s, i) \leftarrow \text{NextVertex}(s, i, l, 4)
13.
       else
14.
            word ← nucleotide string corresponding to s
           if TotalDistance(s,DNA) < bestDistance</pre>
15.
16.
                         bestDistance ← TotalDistance(word, DNA)
                         bestWord ← word
17.
18.
           (s,i) \leftarrow \text{NextVertex}(s,i,l,4)
19.
      return bestWord
```

## Improving the Bounds



- $\bigcirc$ Given an  $\ell$ -mer w, divided into two parts at point i
  - $\boldsymbol{\omega}$   $\boldsymbol{u}$ : prefix  $w_1, ..., w_i$ ,
  - $v : suffix w_{i+1}, ..., w_{\ell}$
- $\bigcirc$  Find minimum distance for u in a sequence
- Note this doesn't tell us anything about whether *u* is part of any motif. We only get a minimum distance for prefix *u*

### Improving the Bounds (cont'd)



- Repeating the process for the suffix v gives us a minimum distance for v
- Since *u* and *v* are two substrings of *w*, and included in motif *w*, we can assume that the minimum distance of *u* plus minimum distance of *v* can only be less than the minimum distance for *w*

### Better Bounds



Searching for prefix VWe may find many instances of prefix V with a minimum distance q

$$\min_{\min} d(v) = q$$
  $\min_{\min} d(v) = q$   $\min_{\min} d(v) = q$ 

Likewise for U

$$\min \mathbf{d}(u) = z$$
  $\min \mathbf{d}(u) = z$ 

But for U and V combined, U is not at its minimum distance location, neither is V

$$\min d(q+1 z+2)$$

But at least we know w (prefix u suffix v) cannot have distance less than  $_{min}d(v) + _{min}d(u)$ 

### Better Bounds (cont'd)

03

 $\alpha$ If  $d(prefix) + d(suffix) \ge bestDistance$ :

✓ Motif w (prefix.suffix) cannot give a better (lower) score than d(prefix) + d(suffix)

In this case, we can ByPass()

### Better Bounded Median String Search

```
ImprovedBranchAndBoundMedianString(DNA,t,n,l)
1.
         s = (1, 1, ..., 1)
         best distance = \infty
4.
         i = 1
5.
         while i > 0
6.
               if i < \ell
                prefix = nucleotide string corresponding to (s_1, s_2, s_3, ..., s_i)
8.
                optimisticPrefixDistance = TotalDistance (prefix, DNA)
9.
                   if (optimisticPrefixDistance < bestsubstring[i])
10.
                      bestsubstring[i] = optimisticPrefixDistance
11.
                      if (l-i < i)
12.
                      optimisticSufxDistance = bestsubstring[[-i]]
13.
                      else
                        optimisticSufxDistance = 0;
14.
15.
                      if optimisticPrefixDistance + optimisticSufxDistance > bestDistance
                         (\mathbf{s}, i) = \text{Bypass}(\mathbf{s}, i, \ell, 4)
16.
17.
                      else
18.
                        (\mathbf{s}, i) = \text{NextVertex}(\mathbf{s}, i, l, 4)
19.
               else
20.
                 word = nucleotide string corresponding to (s_1, s_2, s_3, ..., s_t)
21.
                 if TotalDistance( word, DNA) < bestDistance
                  bestDistance = TotalDistance(word, DNA)
23.
                  bestWord = word
24.
                  (\mathbf{s},i) = \text{NextVertex}(\mathbf{s}, i, l, 4)
25.
        return bestWord
```

### More on the Motif Problem



- They always find the optimal solution, though they may be too slow to perform practical tasks
- Many algorithms sacrifice optimal solution for speed