

---

# CS481: Bioinformatics Algorithms

---

Can Alkan

EA224

`calkan@cs.bilkent.edu.tr`

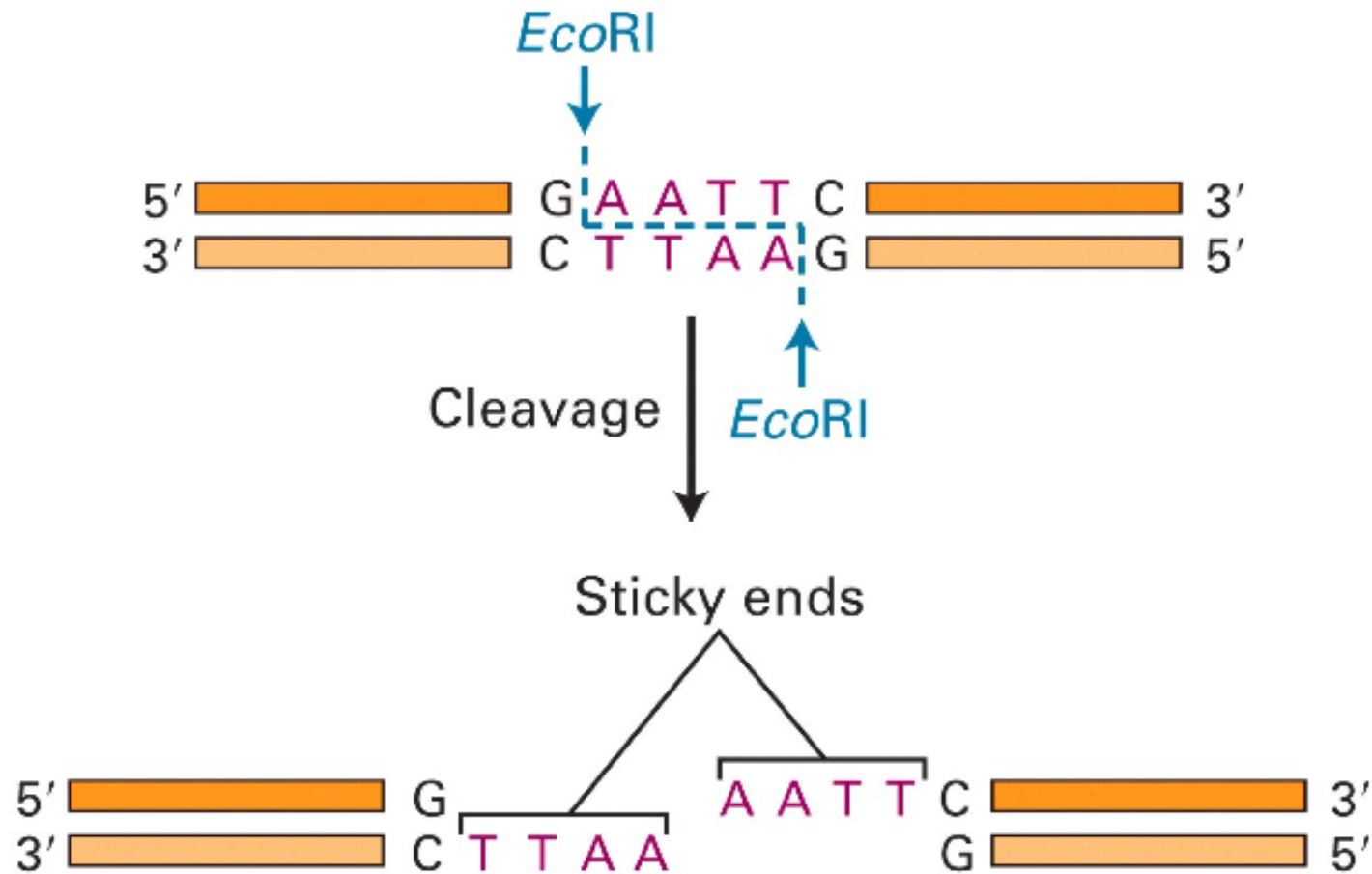
<http://www.cs.bilkent.edu.tr/~calkan/teaching/cs481/>

---

# DNA MAPPING

---

# Molecular Scissors



# Recognition Sites of Restriction Enzymes

Enzyme	Source Microorganism	Recognition Site <sup>a</sup>	Ends Produced
BamI II	<i>Bacillus amyloliquefaciens</i>	<div>↓</div> <div>-G-G-A-T-C-C-</div> <div>-C-C-T-A-G-G-</div> <div>↑</div>	Sticky
EcoRI	<i>Escherichia coli</i>	<div>↓</div> <div>G A A T T C</div> <div>C T T A A G</div> <div>↑</div>	Sticky
HindIII	<i>Haemophilus influenzae</i>	<div>↓</div> <div>-A-A-G-C-T-T-</div> <div>-T-T-C-G-A-A-</div> <div>↑</div>	Sticky
KpnI	<i>Klebsiella pneumonia</i>	<div>↓</div> <div>-G-G-T-A-C-C-</div> <div>-C-C-A-T-G-G-</div> <div>↑</div>	Sticky

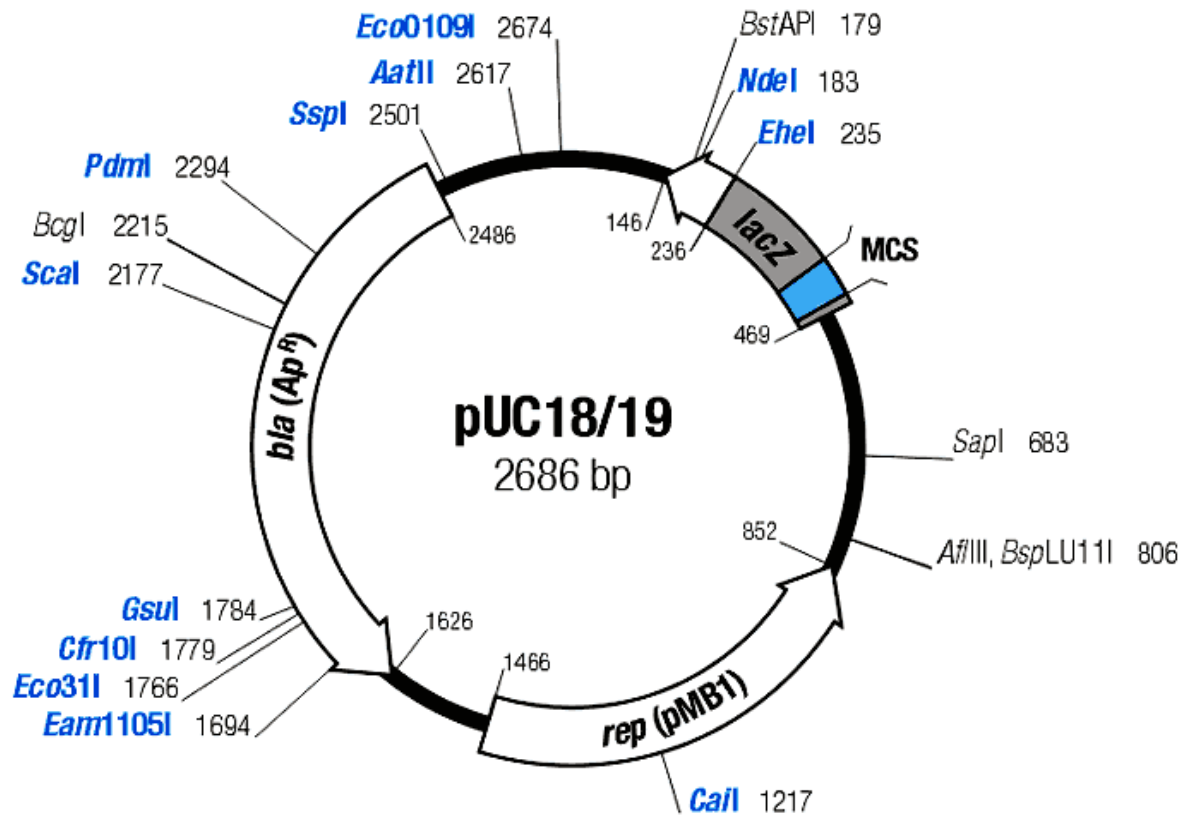
---

# Uses of Restriction Enzymes

- Recombinant DNA technology
  - Cloning
  - cDNA/genomic library construction
  - DNA mapping
-

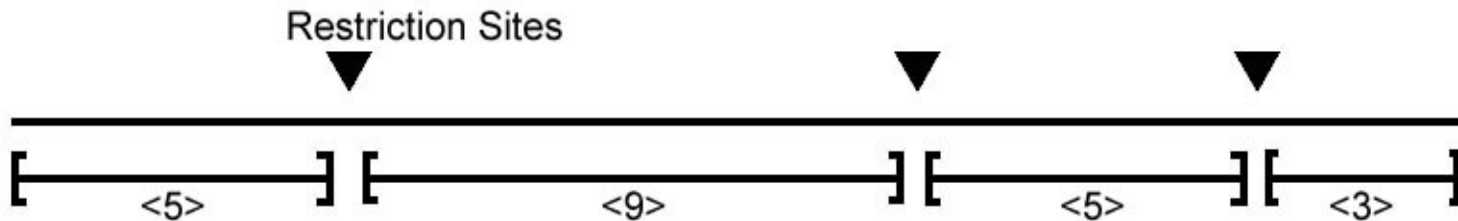
# Restriction Maps

- A map showing positions of restriction sites in a DNA sequence
- If DNA sequence is known then construction of restriction map is a trivial exercise
- In early days of molecular biology DNA sequences were often unknown
- Biologists had to solve the problem of constructing restriction maps **without knowing DNA sequences**



# Full Restriction Digest

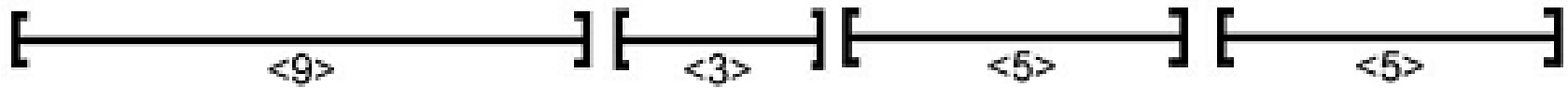
- Cutting DNA at each restriction site creates multiple **restriction fragments**:



- Is it possible to reconstruct the order of the fragments from the sizes of the fragments  $\{3, 5, 5, 9\}$  ?

# Full Restriction Digest: Multiple Solutions

- Alternative ordering of restriction fragments:



VS





---

# Measuring Length of Restriction Fragments

- Restriction enzymes break DNA into restriction fragments.
  - **Gel electrophoresis** is a process for separating DNA by size and measuring sizes of restriction fragments
  - Can separate DNA fragments that differ in length in only 1 nucleotide for fragments up to 500 nucleotides long
-

---

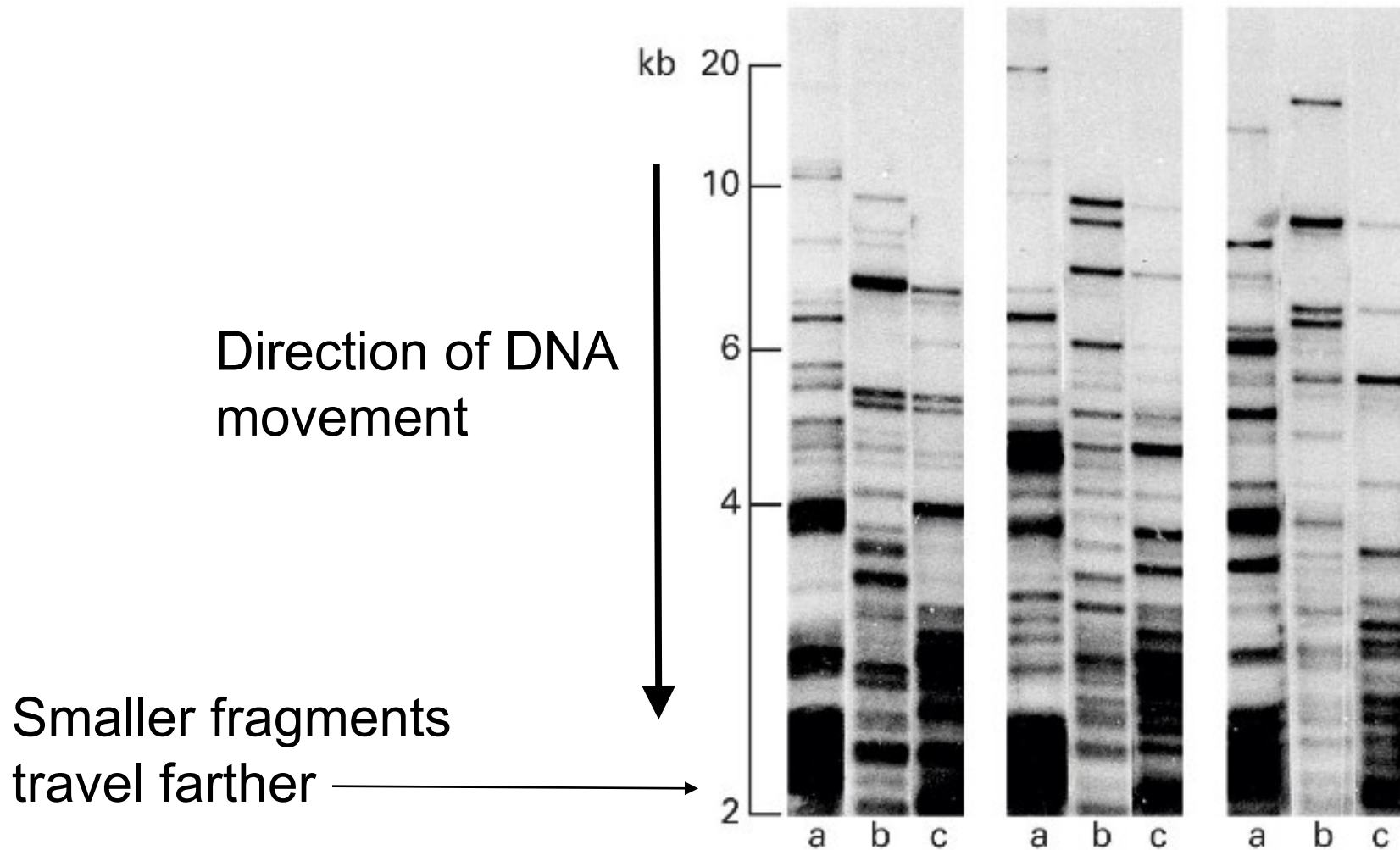
# Gel Electrophoresis

- DNA fragments are injected into a gel positioned in an electric field
  - DNA are negatively charged near neutral pH
    - The ribose phosphate backbone of each nucleotide is acidic; DNA has an overall negative charge
  - DNA molecules move towards the positive electrode
-

# Gel Electrophoresis (cont'd)

- DNA fragments of different lengths are separated according to size
  - Smaller molecules move through the gel matrix more readily than larger molecules
- The gel matrix restricts random diffusion so molecules of different lengths separate into different bands

# Gel Electrophoresis: Example



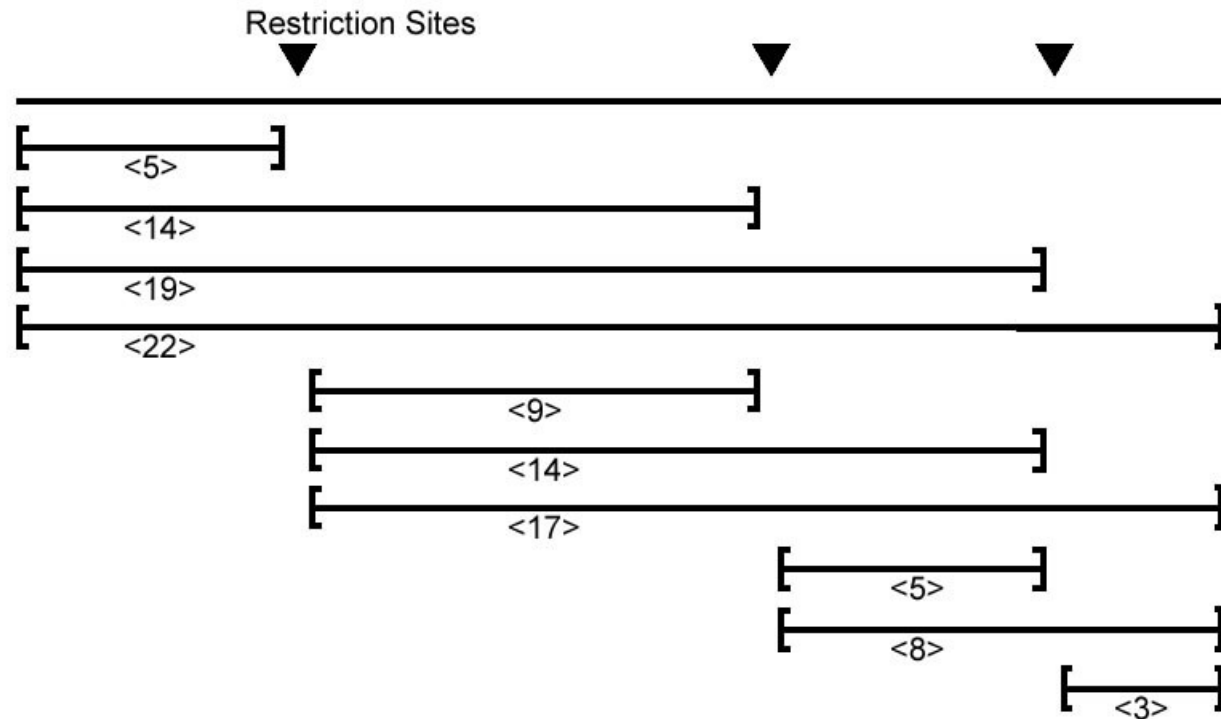
---

# Partial Restriction Digest

- The sample of DNA is exposed to the restriction enzyme for only a limited amount of time to prevent it from being cut at all restriction sites
  - This experiment generates the set of all possible restriction fragments between every two (not necessarily consecutive) cuts
  - This set of fragment sizes is used to determine the positions of the restriction sites in the DNA sequence
-

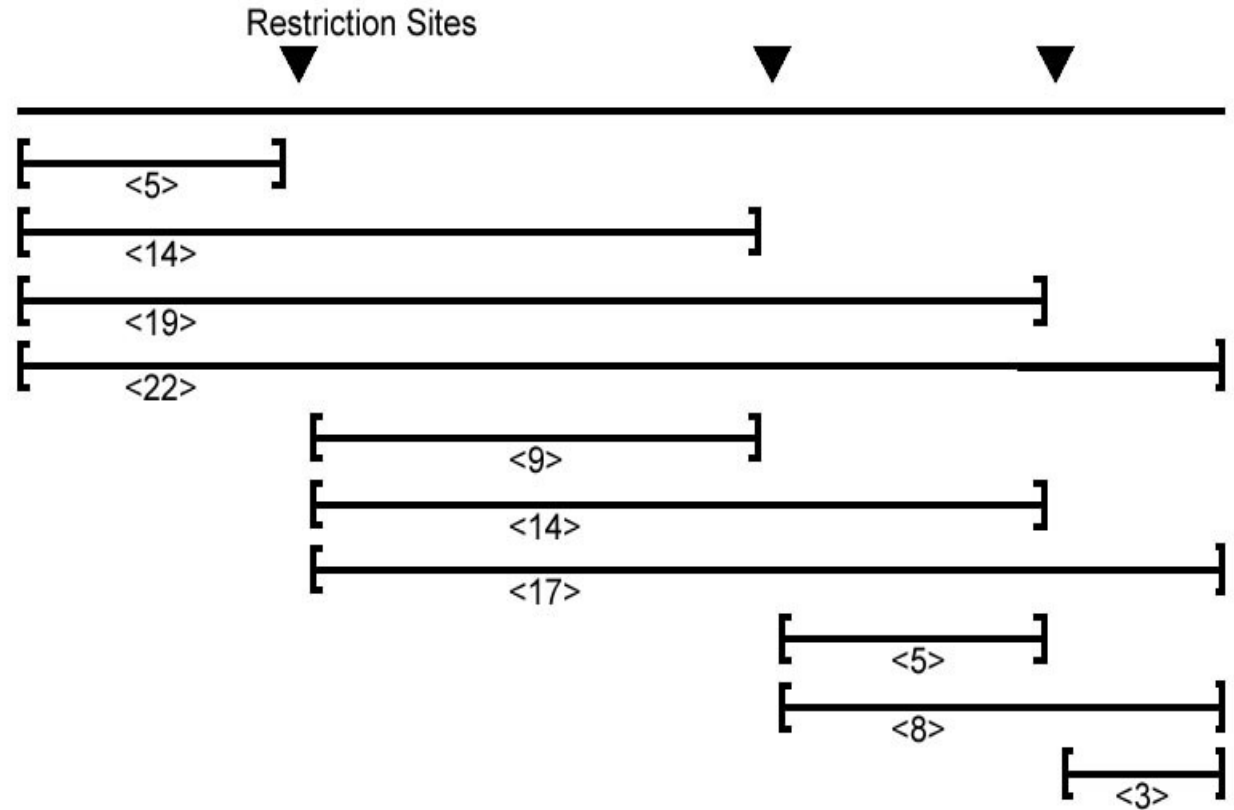
# Partial Digest Example

- Partial Digest results in the following 10 restriction fragments:



# Multiset of Restriction Fragments

- We assume that multiplicity of a fragment can be detected, i.e., the number of restriction fragments of the same length can be determined (e.g., by observing twice as much fluorescence intensity for a double fragment than for a single fragment)



**Multiset:** {3, 5, 5, 8, 9, 14, 14, 17, 19, 22}

# Partial Digest Fundamentals

- $X$** : the set of  **$n$**  integers representing the location of all cuts in the restriction map, including the start and end
- $n$** : the total number of cuts
- $\Delta X$** : the multiset of integers representing lengths of each of the  $C(n, 2)$  fragments produced from a partial digest



# One More Partial Digest Example

X	0	2	4	7	10
0		2	4	7	10
2			2	5	8
4				3	6
7					3
10					

Representation of  $\Delta^X = \{2, 2, 3, 3, 4, 5, 6, 7, 8, 10\}$  as a two dimensional table, with elements of

$$X = \{0, 2, 4, 7, 10\}$$

along both the top and left side. The elements at  $(i, j)$  in the table is  $x_j - x_i$  for  $1 \leq i < j \leq n$ .

# Partial Digest Problem: Formulation

Goal: Given all pairwise distances between points on a line, reconstruct the positions of those points

- Input: The multiset of pairwise distances  $L$ , containing  $n(n-1)/2$  integers
- Output: A set  $X$ , of  $n$  integers, such that  $\Delta X = L$

# Partial Digest: Multiple Solutions

- It is not always possible to uniquely reconstruct a set  $X$  based only on  $\Delta X$ .
- For example, the set

$$X = \{0, 2, 5\}$$

and

$$(X + 10) = \{10, 12, 15\}$$

both produce  $\Delta X = \{2, 3, 5\}$  as their partial digest set.

- The sets  $\{0, 1, 2, 5, 7, 9, 12\}$  and  $\{0, 1, 5, 7, 8, 10, 12\}$  present a less trivial example of non-uniqueness. They both digest into:

$$\{1, 1, 2, 2, 2, 3, 3, 4, 4, 5, 5, 5, 6, 7, 7, 7, 8, 9, 10, 11, 12\}$$

# Homometric Sets

	0	1	2	5	7	9	12
0		1	2	5	7	9	12
1			1	4	6	8	11
2				3	5	7	10
5					2	4	7
7						2	5
9							3
12							

	0	1	5	7	8	10	12
0		1	5	7	8	10	12
1			4	6	7	9	11
5				2	3	5	7
7					1	3	5
8						2	4
10							2
12							

---

# Brute Force Algorithms

- Also known as exhaustive search algorithms; examine every possible variant to find a solution
  - Efficient in rare cases; usually impractical
-

# Partial Digest: Brute Force

1. Find the restriction fragment of maximum length  $M$ .  $M$  is the length of the DNA sequence.

2. For every possible set

$$\mathbf{X} = \{0, x_2, \dots, x_{n-1}, M\}$$

compute the corresponding  $\Delta \mathbf{X}$

5. If  $\Delta \mathbf{X}$  is equal to the experimental partial digest  $L$ , then  $\mathbf{X}$  is the correct restriction map

# BruteForcePDP

1. BruteForcePDP( $L, n$ ):
2.      $M \leftarrow$  maximum element in  $L$
3.     for every set of  $n - 2$  integers  $0 < x_2 < \dots x_{n-1} < M$
4.          $X \leftarrow \{0, x_2, \dots, x_{n-1}, M\}$
5.         Form  $\Delta X$  from  $X$
6.         if  $\Delta X = L$
7.             return  $X$
8.     output “no solution”

# Efficiency of BruteForcePDP

- BruteForcePDP takes  $O(\mathbf{M}^{n-2})$  time since it must examine all possible sets of positions.
- One way to improve the algorithm is to limit the values of  $x_i$  to only those values which occur in  $\mathbf{L}$ .



# AnotherBruteForcePDP

1. AnotherBruteForcePDP( $L, n$ )
2.  $M \leftarrow$  maximum element in  $L$
3. for every set of  $n - 2$  integers  $0 < x_2 < \dots x_{n-1} < M$
4.  $X \leftarrow \{0, x_2, \dots, x_{n-1}, M\}$
5. Form  $\Delta X$  from  $X$
6. if  $\Delta X = L$
7. return  $X$
8. output “no solution”

# AnotherBruteForcePDP

1. AnotherBruteForcePDP( $L, n$ )
2.  $M \leftarrow$  maximum element in  $L$
3. for every set of  $n - 2$  integers  $0 < x_2 < \dots x_{n-1} < M$  from  $L$
4.  $X \leftarrow \{0, x_2, \dots, x_{n-1}, M\}$
5. Form  $\Delta X$  from  $X$
6. if  $\Delta X = L$
7. return  $X$
8. output “no solution”

# Efficiency of AnotherBruteForcePDP

- It's more efficient, but still slow
- If  $L = \{2, 998, 1000\}$  ( $n = 3$ ,  $M = 1000$ ), BruteForcePDP will be extremely slow, but AnotherBruteForcePDP will be quite fast
- Fewer sets are examined, but runtime is still exponential:  $O(n^{2n-4})$

# Branch and bound algorithm for PDP

- By Steven Skiena (Stony Brook Univ.)
- We first define  $\Delta(\mathbf{y}, \mathbf{X})$   
as the multiset of all distances between point  $\mathbf{y}$   
and all other points in the set  $\mathbf{X}$

$$\Delta(y, X) = \{|y - x_1|, |y - x_2|, \dots, |y - x_n|\}$$

$$\text{for } X = \{x_1, x_2, \dots, x_n\}$$

# PartialDigest Algorithm

PartialDigest( $L$ ):

$width \leftarrow$  Maximum element in  $L$

DELETE( $width$ ,  $L$ )

$X \leftarrow \{0, width\}$

PLACE( $L$ ,  $X$ )

# PartialDigest Algorithm (cont'd)

1. PLACE( $L, X$ )
2. if  $L$  is empty
3.     output  $X$
4.     return
5.  $y \leftarrow$  maximum element in  $L$
6. Delete( $y, L$ )
7. if  $\Delta(y, X) \notin L$
8.     Add  $y$  to  $X$  and remove lengths  $\Delta(y, X)$  from  $L$
9.     PLACE( $L, X$ )
10.    Remove  $y$  from  $X$  and add lengths  $\Delta(y, X)$  to  $L$
11. if  $\Delta(\text{width}-y, X) \notin L$
12.     Add  $\text{width}-y$  to  $X$  and remove lengths  $\Delta(\text{width}-y, X)$  from  $L$
13.     PLACE( $L, X$ )
14.     Remove  $\text{width}-y$  from  $X$  and add lengths  $\Delta(\text{width}-y, X)$  to  $L$
15. return

---

# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0 \}$$

---

---

# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0 \}$$

Remove 10 from ***L*** and insert it into ***X***. We know this must be the length of the DNA sequence because it is the largest fragment.

---



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 10 \}$$



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 10 \}$$

Take 8 from  $L$  and make  $y = 2$  or 8. But since the two cases are symmetric, we can assume  $y = 2$ .



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 10 \}$$

We find that the distances from  $y=2$  to other elements in  $X$  are  $\Delta(y, \mathbf{X}) = \{8, 2\}$ , so we remove  $\{8, 2\}$  from  $L$  and add 2 to  $X$ .



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 10 \}$$



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 10 \}$$

Take 7 from  $L$  and make  $y = 7$  or  $y = 10 - 7 = 3$ . We will explore  $y = 7$  first, so  $\Delta(y, X) = \{7, 5, 3\}$ .



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 10 \}$$

For  $y = 7$  first,  $\Delta(y, X) = \{7, 5, 3\}$ . Therefore we remove  $\{7, 5, 3\}$  from  $L$  and add 7 to  $X$ .



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 7, 10 \}$$



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 7, 10 \}$$

Next: take 6 from  $L$  and make  $y = 6$  or  $y = 10 - 6 = 4$ .

$\Delta(y, X) = \{6, 4, 1, 4\}$ , which is NOT a subset of  $L$  so we will NOT explore this branch





# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 7, 10 \}$$

This time make  $y = 4$ .  $\Delta(y, X) = \{4, 2, 3, 6\}$ , which is a subset of  $L$  so we will explore this branch. We remove  $\{4, 2, 3, 6\}$  from  $L$  and add 4 to  $X$ .



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 4, 7, 10 \}$$



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 4, 7, 10 \}$$

$L$  is now empty, so we have a solution, which is  $X$ .



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 7, 10 \}$$

To find other solutions, we backtrack.



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 10 \}$$

More backtrack.



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 2, 10 \}$$

This time we will explore  $y = 3$ .  $\Delta(y, X) = \{3, 1, 7\}$ , which is not a subset of  $L$ , so we won't explore this branch.



# An Example

$$L = \{ 2, 2, 3, 3, 4, 5, 6, 7, 8, 10 \}$$

$$X = \{ 0, 10 \}$$

We backtracked back to the root. Therefore we have found all the solutions.



# Analyzing PartialDigest Algorithm

- Still exponential in worst case, but is very fast on average
- Informally, let  $T(n)$  be time PartialDigest takes to place  $n$  cuts
  - No branching case:  $T(n) < T(n-1) + O(n)$ 
    - Quadratic
  - Branching case:  $T(n) < 2T(n-1) + O(n)$ 
    - Exponential

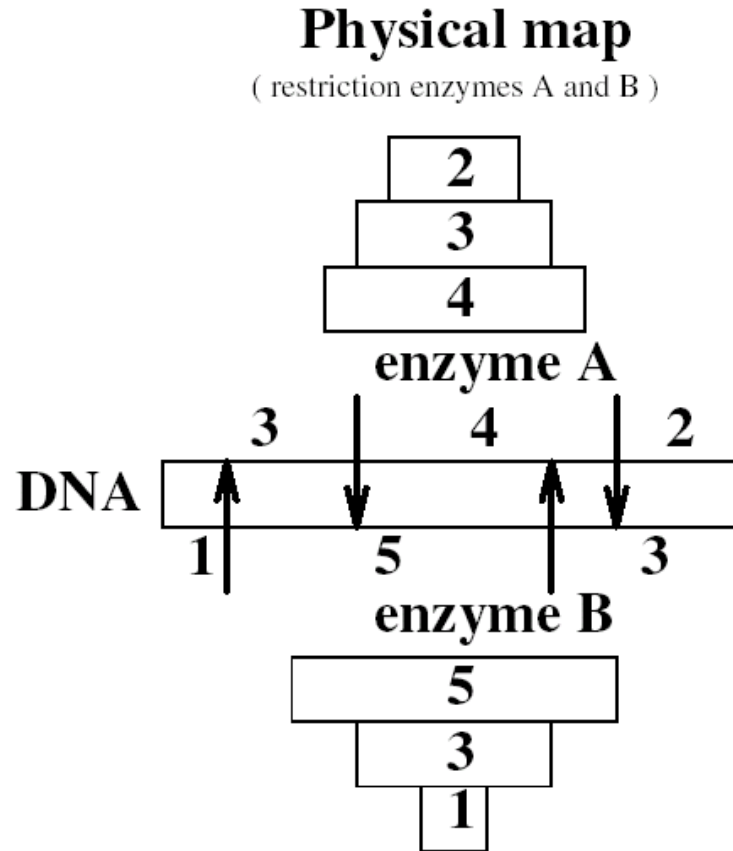


---

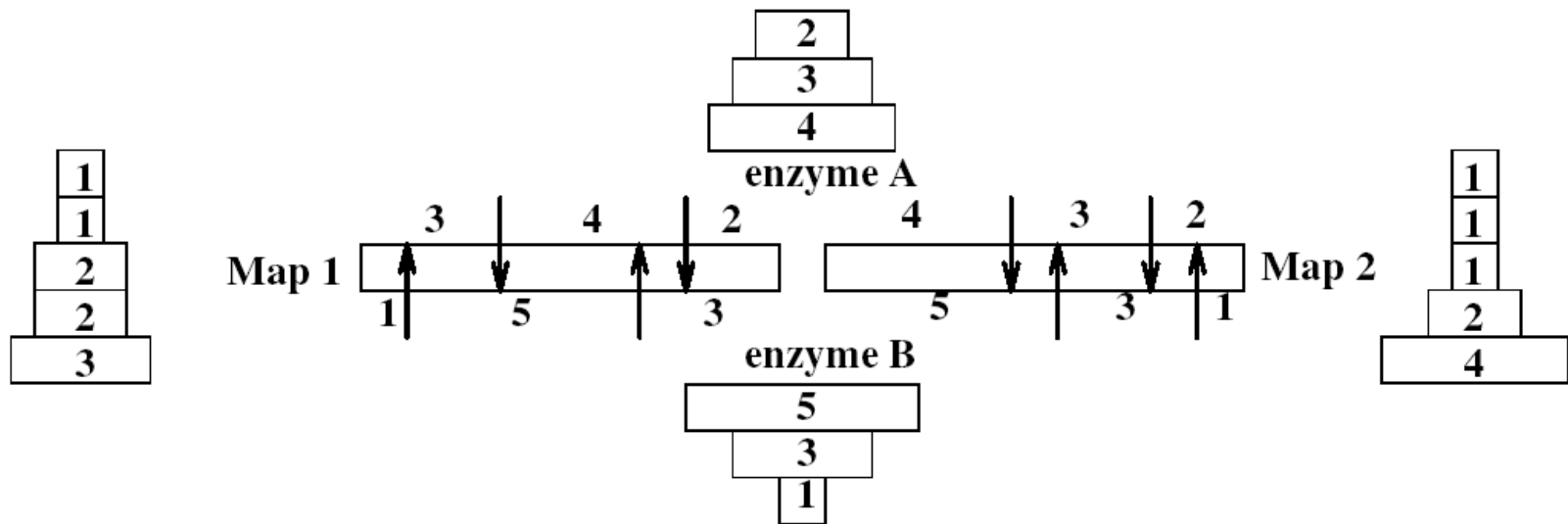
# Double Digest Mapping

- Double Digest is yet another experimentally method to construct restriction maps
    - Use two restriction enzymes; three **full** digests:
      - One with only first enzyme
      - One with only second enzyme
      - One with both enzymes
  - Computationally, Double Digest problem is more complex than Partial Digest problem
-

# Double Digest: Example



# Double Digest: Example



Without the information about ***X*** (i.e. ***A+B***), it is impossible to solve the double digest problem as this diagram illustrates

# Double Digest Problem

Input:  $dA$  – fragment lengths from the digest with enzyme  $A$ .

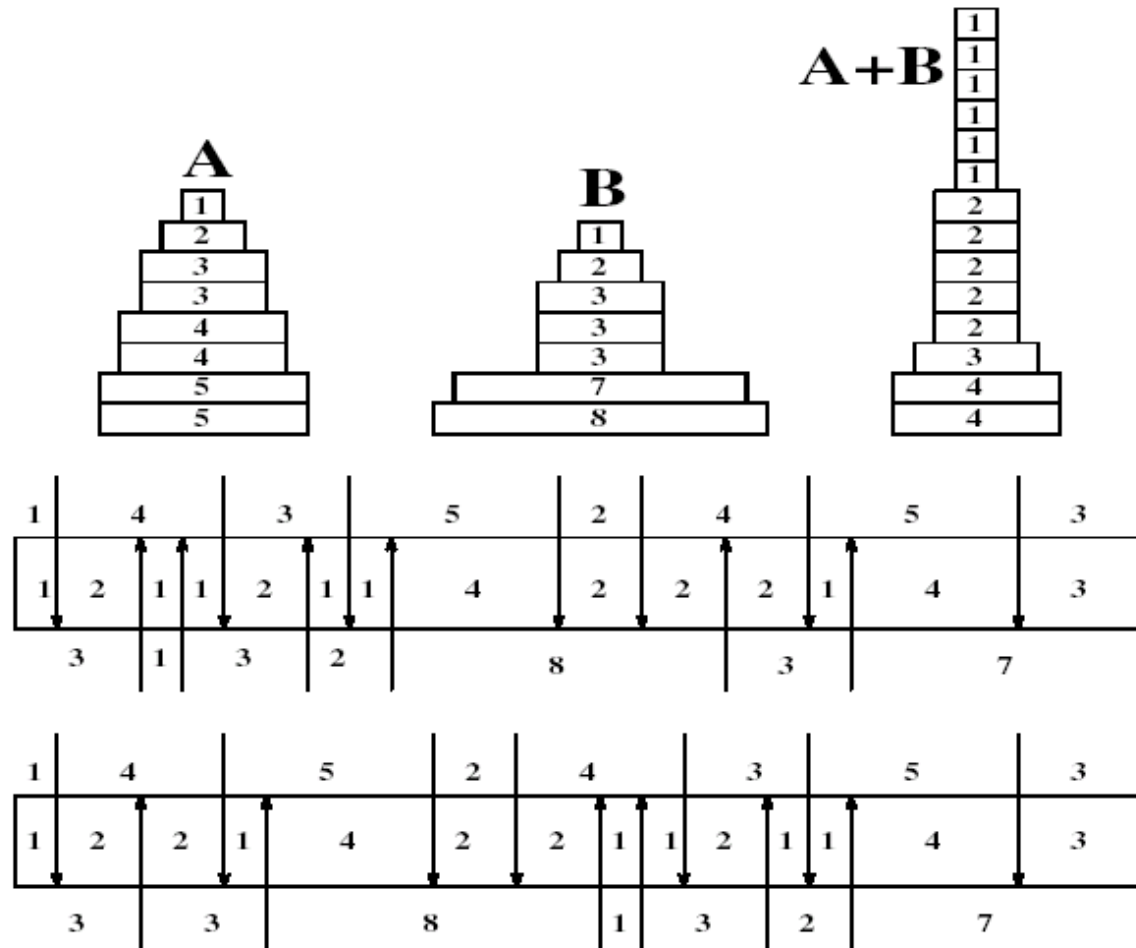
$dB$  – fragment lengths from the digest with enzyme  $B$ .

$dX$  – fragment lengths from the digest with *both*  $A$  and  $B$ .

Output:  $A$  – location of the cuts in the restriction map for the enzyme  $A$ .

$B$  – location of the cuts in the restriction map for the enzyme  $B$ .

# Double Digest: Multiple Solutions



---

# MOTIFS

---

# Random Sample

atgaccgggatactgataccgtatTTTggcctaggcggtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatTTTTtgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaatactgggcataaggtaca  
tgagtatccctgggatgactTTTgggaacactatagtgtctctcccgatTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgaccttgtaagtgtTTTccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga  
tccTTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatggcccacttagtccacttatag  
gtcaatcatgttcttgtgaatggatTTTTaactgagggcatagaccgcttggcgcacccaaattcagtgtgggcgagcgcaa  
cggTTTTggcccttgtagaggccccgtactgatggaaactTTTcaattatgagagagctaatttatcgcggtgcgtgttcat  
aacttgagttggtTTTcgaaaatgctctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatTTTcaacgtatgccgaaccgaaaggggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttctgggtactgatagca

# Implanting Motif AAAAAAAAAAGGGGGGGG

atgaccgggatactgatAAAAAAAAAGGGGGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatttttttgagcagatttagtgacctggaaaaaaatttgagtacaaaacttttccgaataAAAAAAAAAGGGGGGGGa  
tgagtatccctgggatgacttAAAAAAAAAGGGGGGGtgctctcccgatttttgaatatgtaggatcattcgccagggtccga  
gctgagaattggatgAAAAAAAAAGGGGGGGtccacgcaatcgcgaaacacgcggacccaaaggcaagaccgataaaggaga  
tcccttttgcggaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAAAAAAAAAGGGGGGGcttatag  
gtcaatcatgttcttgtgaatggatttAAAAAAAAAGGGGGGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa  
cggttttgcccttggttagaggccccctAAAAAAAAAGGGGGGGcaattatgagagagctaattctatcgctgctgttcat  
aacttgagttAAAAAAAAAGGGGGGGctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatAAAAAAAAAGGGGGGGaccgaaagggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcgaagcttAAAAAAAAAGGGGGGGa



# Where is the Implanted Motif?

atgaccgggatactgataaaaaaagggggggggcgctacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataaaaaaaaaggggggga  
tgagtatccctgggatgacttaaaaaaaggggggggtgctctcccgatTTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgaaaaaaaggggggggtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga  
tccTTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaataaaaaaagggggggcttatag  
gtcaatcatgttcttgtgaatggatttaaaaaaaggggggggaccgcttggcgcacccaaattcagtgtgggcgagcgcaa  
cggTTTTggcccttgttagaggcccccgtaaaaaaaggggggggcaattatgagagagctaatttatcgctgctgttcat  
aacttgagttaaaaaaaggggggggtggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataaaaaaagggggggaccgaaagggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttaaaaaaaggggggga

# Implanting Motif AAAAAAAGGGGGGGG with Four Mutations

atgaccgggatactgatAgAAgAAAGGttGGGggcggtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatttttttgagcagatttagtgacctggaaaaaaatttgagtacaaaacttttccgaataCAAtAAACGGcGGGa  
tgagtatccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgatttttgaatatgtaggatcattcgccagggtccga  
gctgagaattggatgCAAAAAAGGGAttGtccacgcaatcggaaccaacgcggacccaaaggcaagaccgataaaggaga  
tcccttttgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAtAAtAAAGGaaGGGccttatag  
gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgacccaaattcagtgtgggcgagcgcaa  
cggttttgcccttggttagaggcccccgAtAAACAAAGGaGGGcCaattatgagagagctaattctatcgctgctgttcat  
aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcAtAAAAAGGaGcGGaccgaaagggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa

# Where is the Motif???

atgaccgggatactgatagaagaaagggttgggggcgtagacacattagataaacgtatgaagtacgttagactcggcgccgccg  
acccctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaatacaataaaaacggcgga  
tgagtatccctgggatgacttaaaataatggagtggtgctctcccgattTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgcaaaaaagggttgtccacgcaatcggaaccaacgcggacccaaaggcaagaccgataaaggaga  
tccTTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaataataaaaggaagggttatag  
gtcaatcatgttcttgtgaatggatttaacaataagggtgggaccgcttggcgcacccaaattcagtgtgggagcgcaa  
cggtTTTggcccttggtagaggccccgtataaacaaggaggggccaattatgagagagctaatttatcgctgctgttcat  
aacttgagttaaaaaatagggagccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatactaaaaaggagcggaccgaaagggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttactaaaaaggagcgga

# Finding (15,4) Motif

atgaccgggatactgatAgAAgAAAGGttGGGggcggtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctattttttgagcagatttagtgacctggaaaaaaatttgagtacaaaacttttccgaataCAAtAAAAcGGcGGGa  
tgagtatccctgggatgacttAAAAtAAtGGAgtGGtgctctcccgatttttgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgCAAAAAAGGGattGtccacgcaatcggaaccaacgcggacccaaaggcaagaccgataaaggaga  
tcccttttgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAtAAtAAAGGaaGGGcttatag  
gtcaatcatgttcttgtgaatggatttAAcAAtAAGGGctGGgaccgcttggcgcacccaaattcagtgtgggagagcgcaa  
cggttttggcccttggtagaggccccgtAtAAACAAAGGaGGGcCaattatgagagagctaattctatcgctgctgttcat  
aacttgagttAAAAAAtAGGGaGccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataActAAAAGGaGcGGgaccgaaaggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAGGaGcGGa

AgAAgAAAGGttGGG  
CAAtAAAAcGGcGGG

---

# Challenge Problem

- 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.
-

---

# Combinatorial Gene Regulation

- An experiment showed that when gene X is knocked out, 20 other genes are not expressed
  - **How can one gene have such drastic effects?**
-

---

# Regulatory Proteins

- Gene X encodes regulatory protein, a.k.a. a ***transcription factor (TF)***
  - The 20 unexpressed genes rely on gene X's TF to induce transcription
  - A single TF may regulate multiple genes
-

# Regulatory Regions

- Every gene contains a regulatory region (RR) typically stretching 100-1000 bp upstream of the transcriptional start site
- Located within the RR are the **Transcription Factor Binding Sites** (TFBS), also known as **motifs**, specific for a given transcription factor
- TFs influence gene expression by binding to a specific location in the respective gene's regulatory region - TFBS



---

# Transcription Factor Binding Sites

- A TFBS can be located anywhere within the Regulatory Region.
  - TFBS may vary slightly across different regulatory regions since non-essential bases could mutate
-

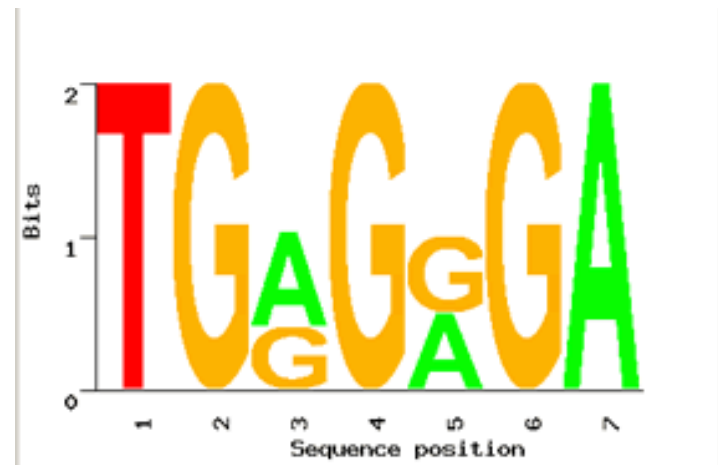
# Motifs and Transcriptional Start Sites



# Motif Logo

- Motifs can mutate on non important bases
- The five motifs in five different genes have mutations in position 3 and 5
- Representations called *motif logos* illustrate the conserved and variable regions of a motif

TGGGGGA  
TGAGAGA  
TGGGGGA  
TGAGAGA  
TGAGGGA



---

# Identifying Motifs

- Genes are turned on or off by regulatory proteins
  - These proteins bind to upstream regulatory regions of genes to either attract or block an RNA polymerase
  - Regulatory protein (TF) binds to a short DNA sequence called a motif (TFBS)
  - So finding the same motif in multiple genes' regulatory regions suggests a regulatory relationship amongst those genes
-

# 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
- How to discern it from “random” motifs?

# The Motif Finding Problem

- Given a random sample of DNA sequences:

```
cctgatagacgctatctggctatccacgtacgtaggtcctctgtgCGaatctatgcgtttccaacCAT  
agtactgggtgtacattttgatacgtacgtacaccggcaacctgaaacaaacgctcagaaccagaagtgc  
aaacgtacgtgcaccctcttttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt  
agcctccgatgtaagtcatagctgtaactattacctgccaccctattacatcttacgtacgtataca  
ctgttatacaacgcgctcatggcgggggtatgcgttttggtcgctcgtacgctcgatcgttaacgtacgtc
```

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

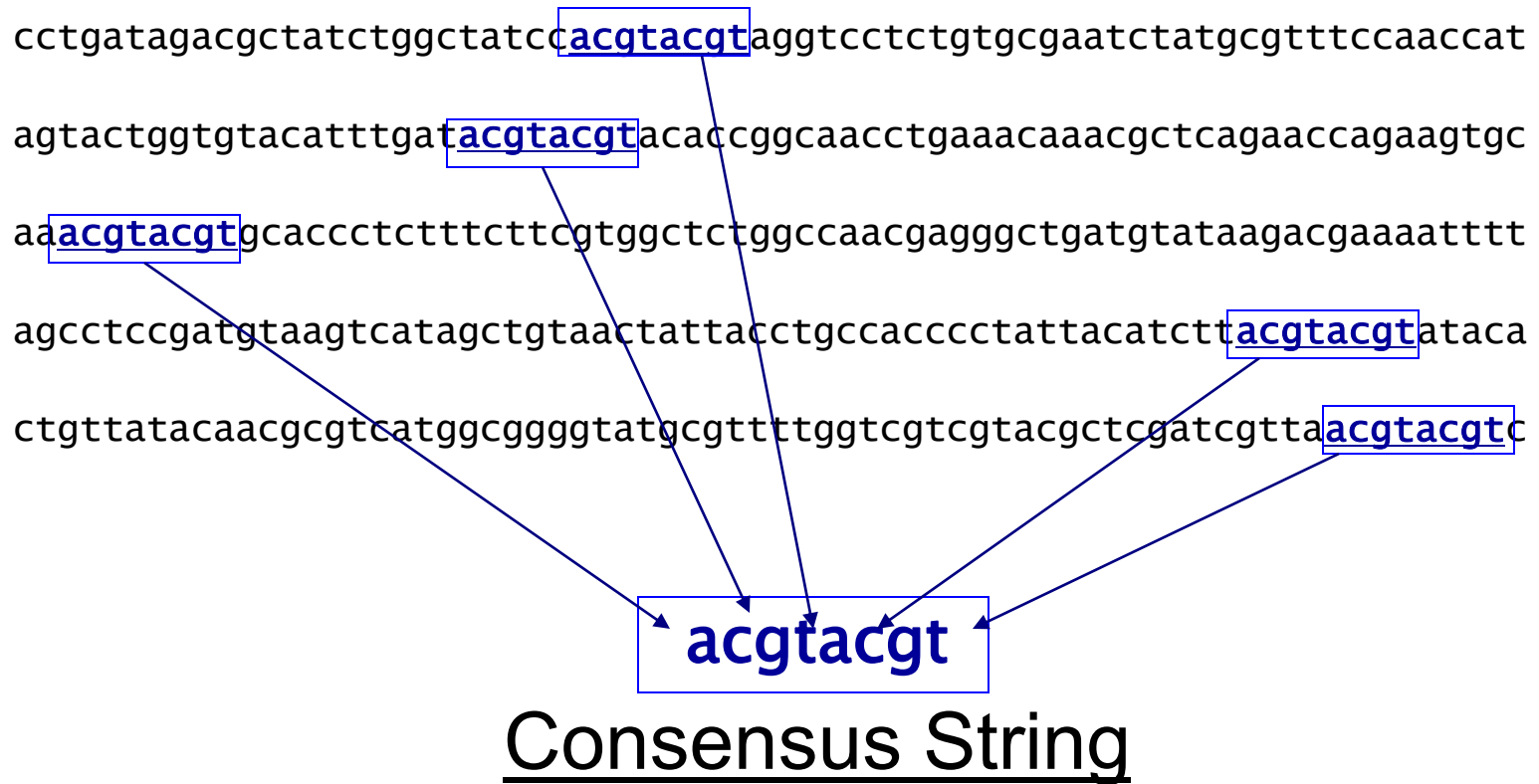
---

# The Motif Finding Problem (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
-

# The Motif Finding Problem (cont'd)

## ■ The patterns revealed with no mutations:





# The Motif Finding Problem (cont'd)

- The patterns with 2 point mutations:

cctgatatagacgctatctggctatccaaGgtacTtaggtcctctgtgCGaatctatgCGtttccaacccat  
agtactgggtgtacatttgatCcAtacgtacaccggcaacctgaaacaaacgctcagaaccagaagtgc  
aaacgtTAgtgcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt  
agcctccgatgtaagtcatagctgtaactattacctgccacccctattacatcttacgtCcAtataca  
ctgttatacaacgcgctcatggcggggtatgcgttttggtcgctcgctacgctcgatcgттаCcgtacgGc

# The Motif Finding Problem (cont'd)

- The patterns with 2 point mutations:

cctgatatagacgctatctggctatccaGgtacTtaggtcctctgtgCGaatctatgCGtttccaaccat  
agtactgggtgtacatttgatCcAtacgtacaccggcaacctgaaacaaacgctcagaaccagaagtgc  
aaacgtTAgtgcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt  
agcctccgatgtaagtcatagctgtaactattacctgccaccctattacatcttacgtCcAtataca  
ctgttatacaacgcgctcatggcggggtatgcgttttggtcgctcgctacgctcgatcgттаCcgtacgGc

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

# Defining Motifs

- To define a motif, let's say we know where the motif starts in the sequence
- The motif start positions in their sequences can be represented as  $\mathbf{s} = (s_1, s_2, s_3, \dots, s_t)$



# Motifs: Profiles and Consensus

Alignment

a	G	g	t	a	c	T	t
C	c	A	t	a	c	g	t
a	c	g	t	T	A	g	t
a	c	g	t	C	c	A	t
C	c	g	t	a	c	g	G

Profile

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

A C G T A C G T

- Line up the patterns by their start indexes

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

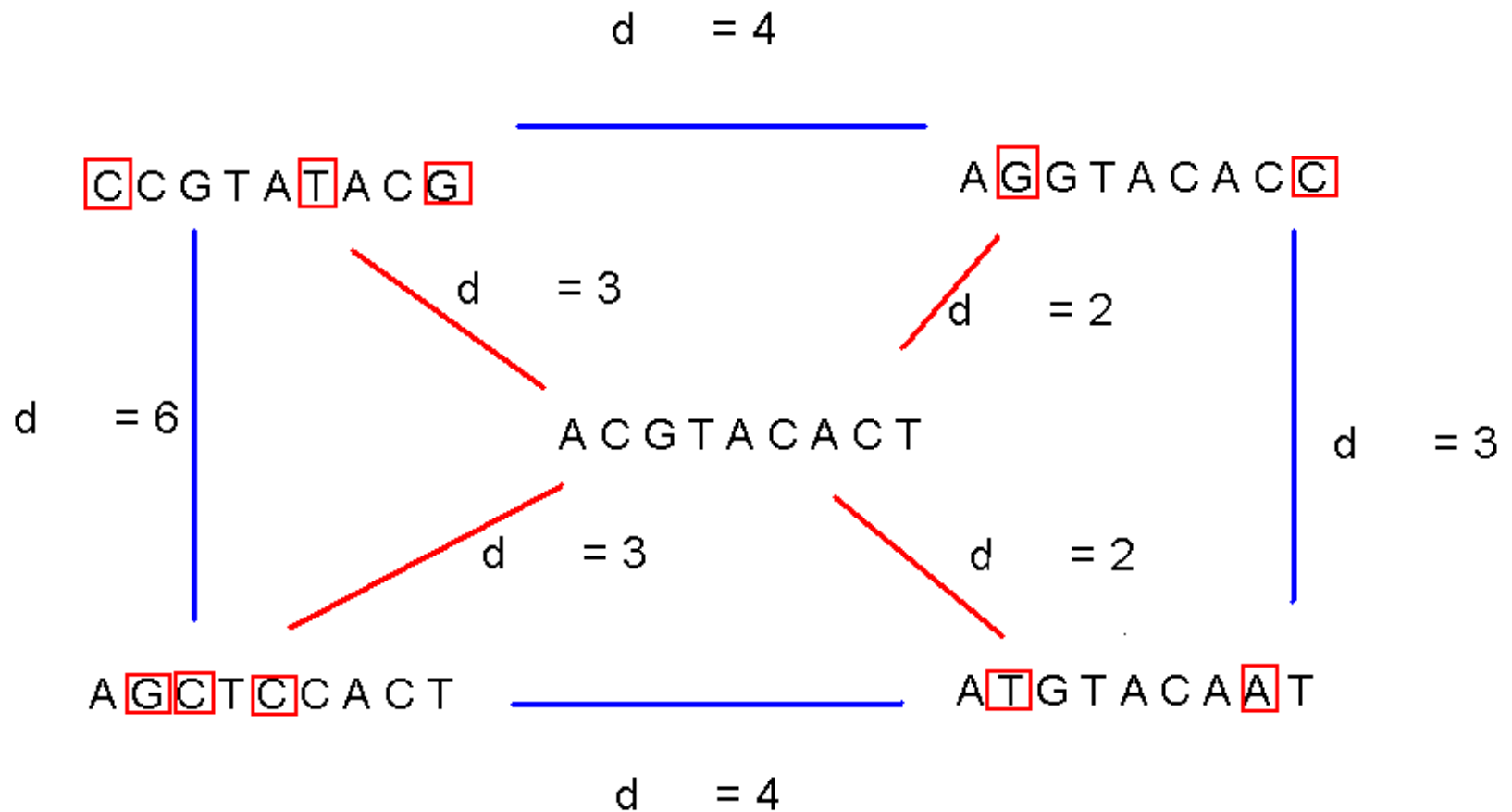
- Construct matrix profile with frequencies of each nucleotide in columns
- Consensus nucleotide in each position has the highest score in column

---

# 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

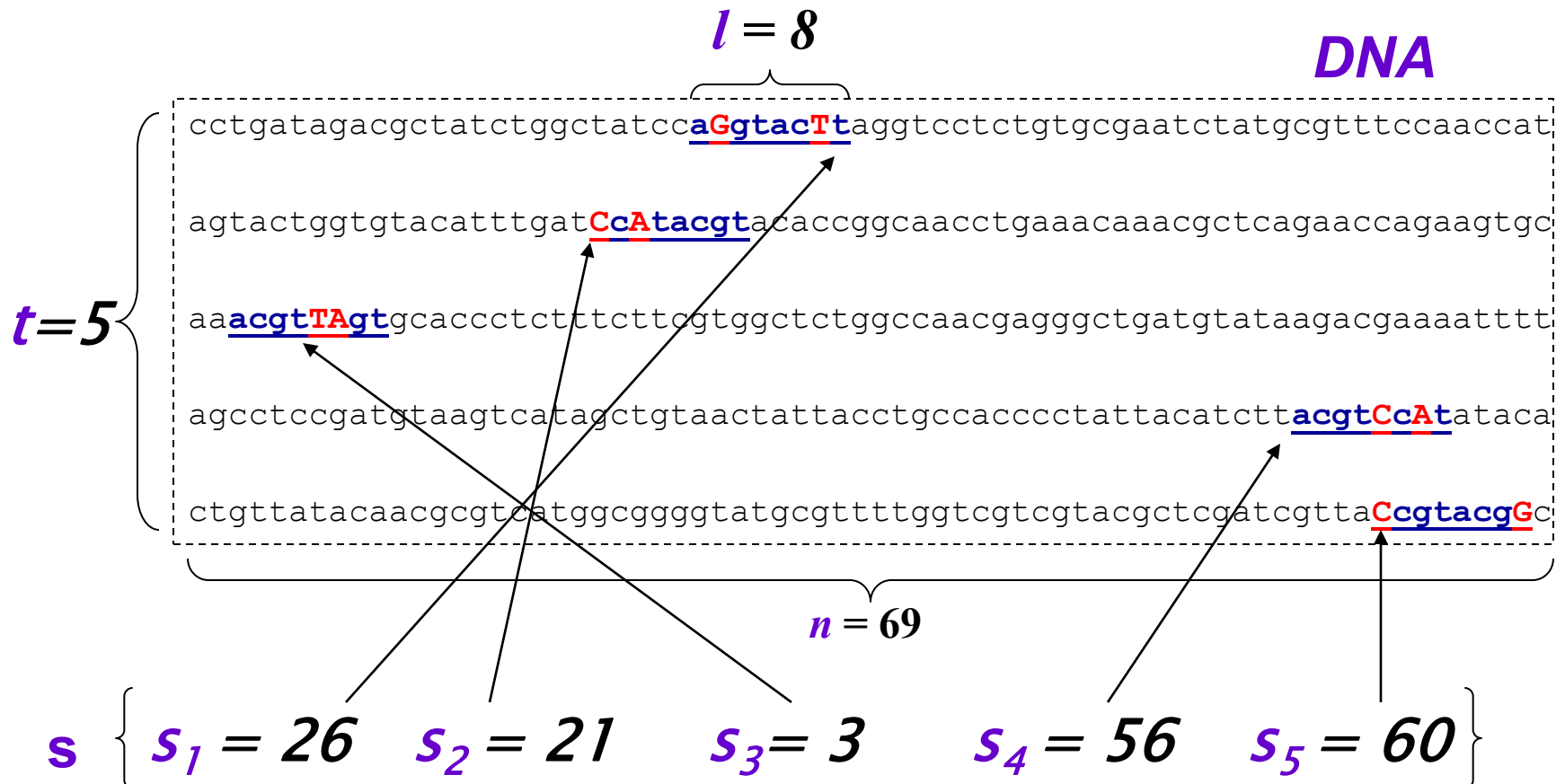
- We have a guess about the consensus sequence, but how “good” is this consensus?
- Need to introduce a scoring function to compare different guesses and choose the “best” one.

# Defining Some Terms

- $t$  - number of sample DNA sequences
- $n$  - length of each DNA sequence
- **DNA** - sample of DNA sequences ( $t \times n$  array)
  
- $\ell$  - length of the motif ( $\ell$ -mer)
- $s_i$  - starting position of an  $\ell$ -mer in sequence  $i$
- $\mathbf{s} = (s_1, s_2, \dots, s_t)$  - array of motif's starting positions



# Parameters



# Scoring Motifs

- Given  $\mathbf{s} = (s_1, \dots, s_t)$  and  $DNA$ :

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

$l$							
a	G	g	t	a	c	T	t
C	c	A	t	a	c	g	t
a	c	g	t	T	A	g	t
a	c	g	t	C	c	A	t
C	c	g	t	a	c	g	G
$t$							

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    a c g t a c g t

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

## 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  $\text{Score}(\mathbf{s}, \mathbf{DNA})$  by varying the starting positions  $\mathbf{s}_i$ , where:

$$\begin{aligned} \mathbf{s}_i &= [1, \dots, n-\ell+1] \\ i &= [1, \dots, t] \end{aligned}$$

# BruteForceMotifSearch

1. BruteForceMotifSearch(*DNA*, *t*, *n*, *l*)
2. *bestScore*  $\leftarrow$  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*  $\leftarrow$   $(s_1, s_2, \dots, s_t)$
7. return *bestMotif*

# Running Time of BruteForceMotifSearch

- Varying  $(n - \ell + 1)$  positions in each of  $t$  sequences, we're looking at  $(n - \ell + 1)^t$  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  $t = 8$ ,  $n = 1000$ ,  $\ell = 10$  we must perform approximately  $10^{20}$  computations – it will take billions of years

---

# 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

- Hamming distance:
  - $d_H(\mathbf{v}, \mathbf{w})$  is the number of nucleotide pairs that do not match when  $\mathbf{v}$  and  $\mathbf{w}$  are aligned. For example:

$$d_H(\text{AAAAAA}, \text{ACAAAC}) = 2$$

# Total Distance: An Example

- Given  $v = \text{"acgtacgt"}$  and  $s$

$$d_H(v, x) = 0$$

acgtacgt

cctgatagacgctatctggctatccacgtacgttaggtcctctgtgcaatctatgcgttttccaacat

$$d_H(v, x) = 0$$

acgtacgt

agtactgggtgtacatttgatacgtacgtacaccggcaacctgaaacaaacgctcagaaccagaaagtgc

acgtacgt

$$d_H(v, x) = 0$$

aaacgtacgtgcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt

$$d_H(v, x) = 0$$

acgtacgt

agcctccgatgtaagtcatactgtaactattacctgccaccctattacatcttacgtacgtataca

$$d_H(v, x) = 0$$

acgtacgt

ctgttatacaacgcgtcatggcggggtatgcgttttggctcgtctacgtcgcgacgttaacgtacgtc

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

- $TotalDistance(v, DNA) = 0$

# Total Distance: Example

- Given  $v = \text{"acgtacgt"}$  and  $s$

$$d_H(v, x) = 1$$

cctgatagacgctatctggctatccacgtacgtacggtcctctgtgcgaatctatgcggtttccaaccat

$$d_H(v, x) = 0$$

agtagtggtgtacatttgatacgtacgtacacccggcaacctgaaacaaacgctcagaaccagaagtgc

$$d_H(v, x) = 2$$

aaaAgtCcgtagcaccctctttcttcgtggctctggccaacgagggctgatgtataagacgaaaatttt

$$d_H(v, x) = 0$$

agcctccgatgtaagtcatactgtaactattacctgccacccctattacatcttacgtacgtataca

$$d_H(v, x) = 1$$

ctgttatacaacgcgctcatggcgggggatgcggtttggtcgctcgtacgctcgatcggttacgtacgtac

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

- $TotalDistance(v, DNA) = 1 + 0 + 2 + 0 + 1 = 4$

# Total Distance: Definition

- For each DNA sequence  $i$ , compute all  $d_H(\mathbf{v}, \mathbf{x})$ , where  $\mathbf{x}$  is an  $\ell$ -mer with starting position  $s_i$  ( $1 \leq s_i \leq n - \ell + 1$ )
- Find minimum of  $d_H(\mathbf{v}, \mathbf{x})$  among all  $\ell$ -mers in sequence  $i$
- $TotalDistance(\mathbf{v}, \mathbf{DNA})$  is the sum of the minimum Hamming distances for each DNA sequence  $i$
- $TotalDistance(\mathbf{v}, \mathbf{DNA}) = \min_{\mathbf{s}} d_H(\mathbf{v}, \mathbf{s})$ , where  $\mathbf{s}$  is the set of starting positions  $s_1, s_2, \dots, s_t$

# The Median String Problem: Formulation

- Goal: Given a set of DNA sequences, find a median string
- Input: A  $t \times n$  matrix  $DNA$ , and  $\ell$ , the length of the pattern to find
- Output: A string  $\mathbf{v}$  of  $\ell$  nucleotides that **minimizes**  $TotalDistance(\mathbf{v}, \mathbf{DNA})$  over all strings of that length

# Median String Search Algorithm

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***

---

## 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

Alignment

a	G	g	t	a	c	T	t
C	c	A	t	a	c	g	t
a	c	g	t	T	A	g	t
a	c	g	t	C	c	A	t
C	c	g	t	a	c	g	G

Profile

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      a c g t a c g t

Score            3+4+4+5+3+4+3+4

TotalDistance 2+1+1+0+2+1+2+1

Sum              5 5 5 5 5 5 5 5

- At any column  $i$   
 $Score_i + TotalDistance_i = t$
- Because there are  $l$  columns  
 $Score + TotalDistance = l * t$
- Rearranging:  
 $Score = l * t - TotalDistance$
- $l * t$  is constant the minimization 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  $(n - \ell + 1)^t$  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

## Recall the BruteForceMotifSearch:

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

# Structuring the Search

- How can we perform the line

for each  $\mathbf{s}=(s_1, s_2, \dots, s_t)$  from  $(1, 1 \dots 1)$  to  $(n-\ell+1, \dots, n-\ell+1)$  ?

- We need a method for efficiently structuring and navigating the many possible motifs
- This is not very different than exploring all  $t$ -digit numbers

# Median String: Improving the Running Time

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

- For the Median String Problem we need to consider all  $4^l$  possible  $l$ -mers:

$\overbrace{\text{aa... aa}}^l$   
aa... ac  
aa... ag  
aa... at  
.  
.  
tt... tt

How to organize this search?

## Alternative Representation of the Search Space

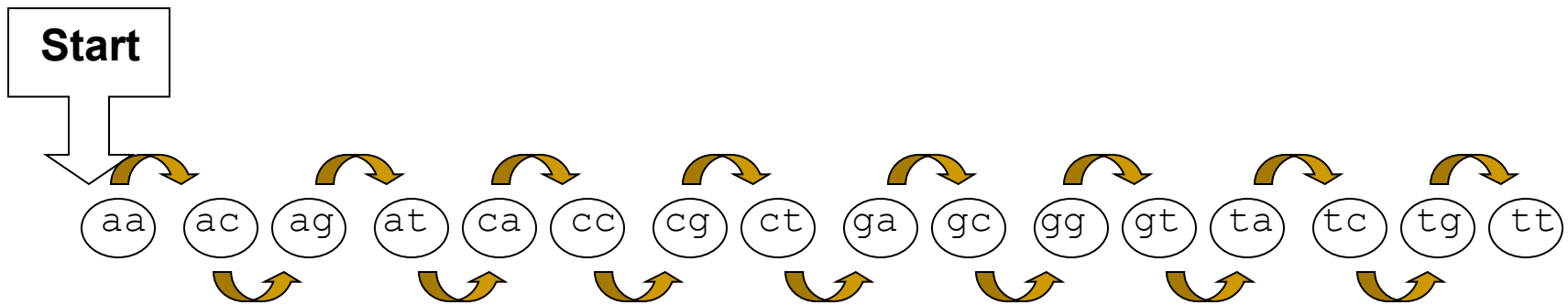
- Let **A** = 1, **C** = 2, **G** = 3, **T** = 4
- Then the sequences from AA...A to TT...T become:

$\ell$   
 $\overbrace{11 \dots 11}$   
11...12  
11...13  
11...14  
.  
.  
44...44

- 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

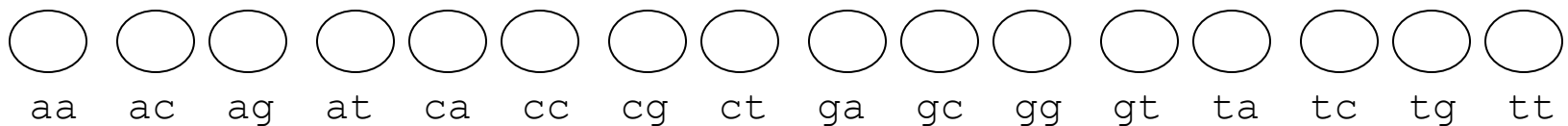
- Suppose  $l = 2$



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

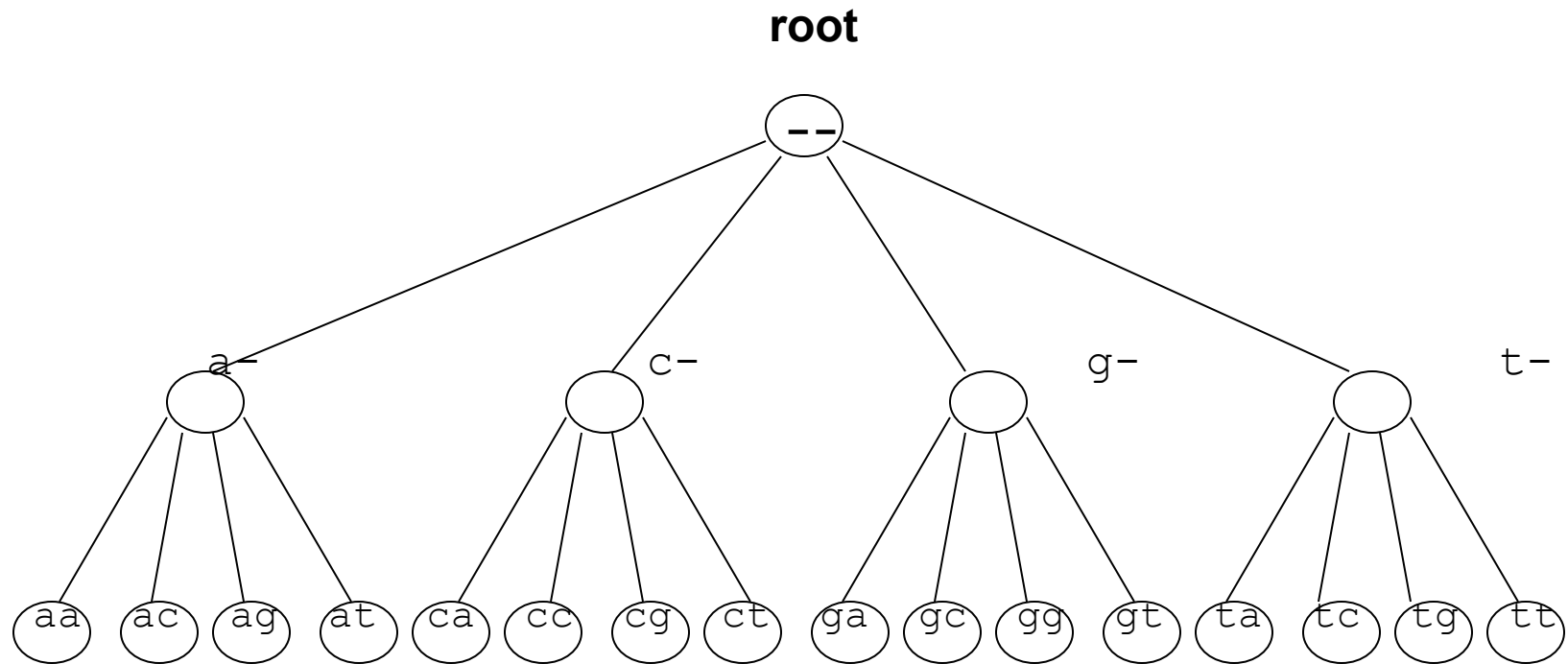
# Linked List (cont'd)

- Linked list is not the most efficient data structure for motif finding
- Let's try grouping the sequences by their prefixes





# 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
    - Visit all the leaves
    - Visit the next node
    - Bypass the children of a node
-

# Visit the Next Leaf

**Given a current leaf  $a$ , we need to compute the “next” leaf:**

1. NextLeaf(  $a, L, k$  )      //  $a$  : the array of digits
2. **for**  $i \leftarrow L$  to 1      //  $L$ : length of the array
3.     **if**  $a_i < k$       //  $k$ : max digit value
4.          $a_i \leftarrow a_i + 1$
5.         **return**  $a$
6.          $a_i \leftarrow 1$
7. **return**  $a$

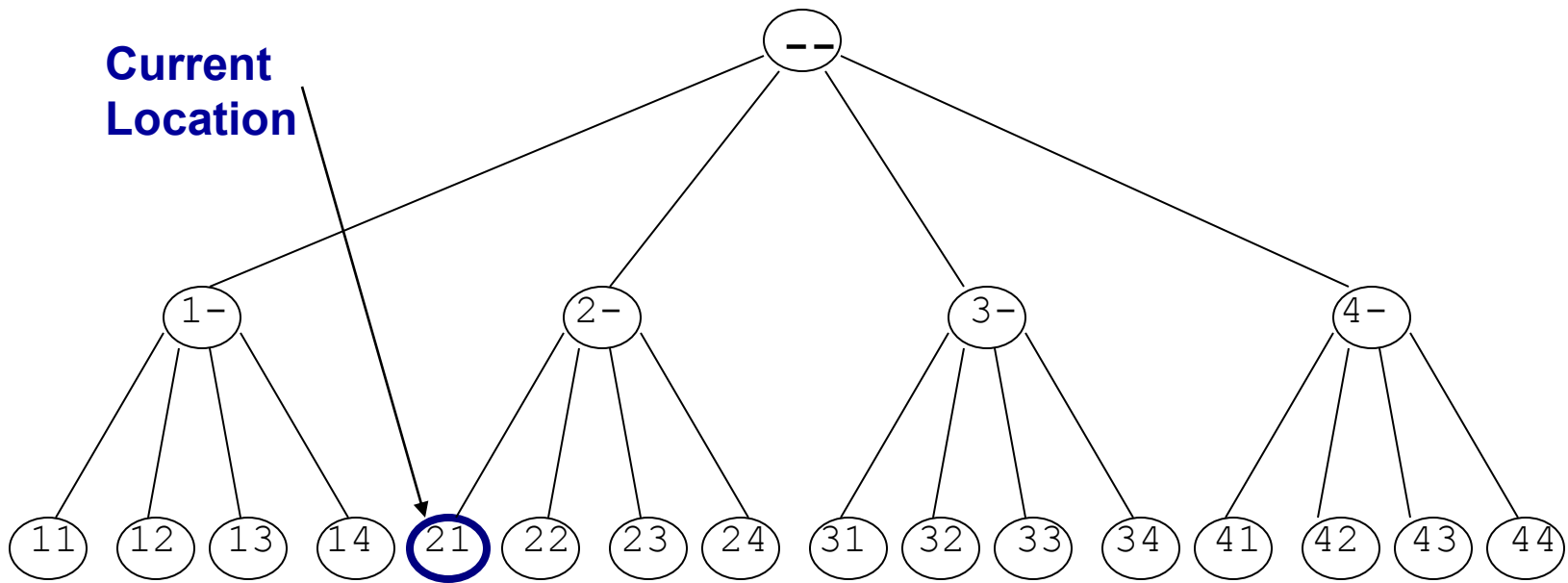
---

## NextLeaf (cont'd)

- The algorithm is common addition in radix  $k$ :
  - Increment the least significant digit
  - “Carry the one” to the next digit position when the digit is at maximal value
-

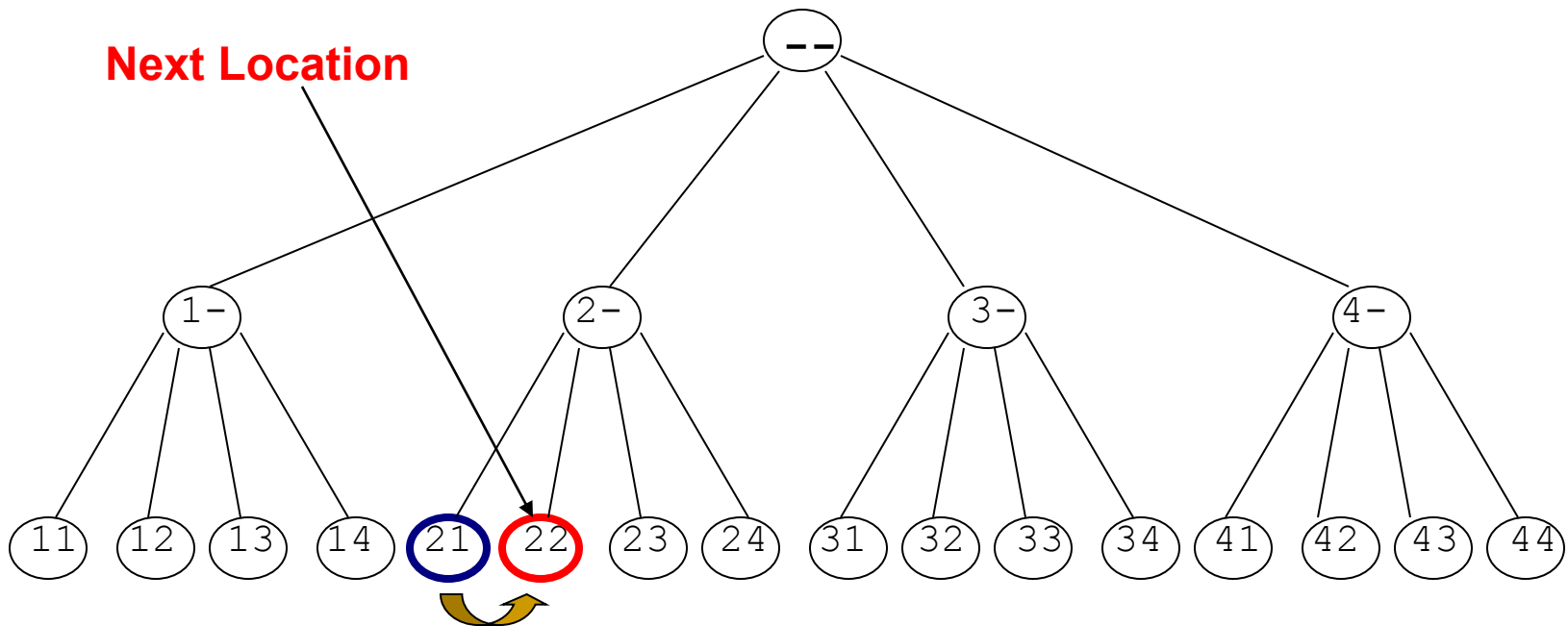
# NextLeaf: Example

- Moving to the next leaf:



# NextLeaf: Example (cont'd)

- Moving to the next leaf:



# Visit All Leaves

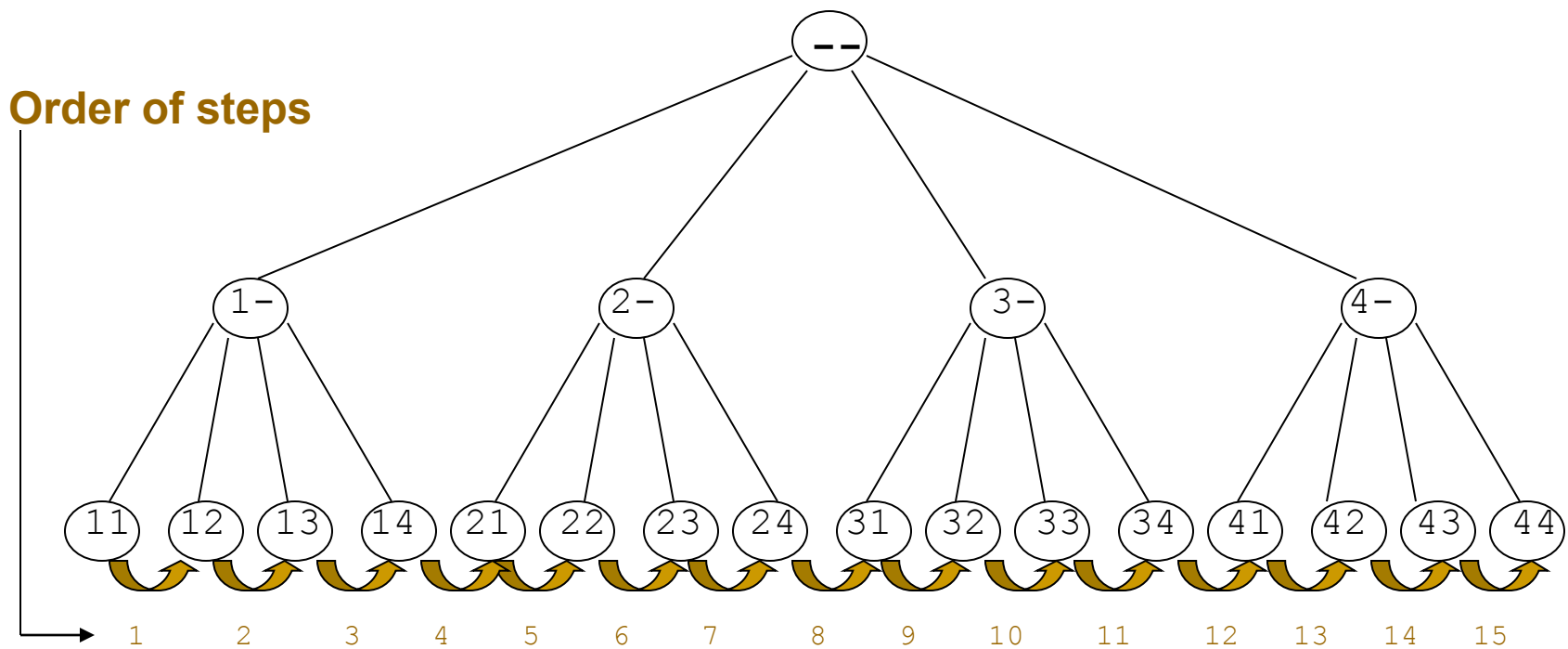
- Printing all permutations in ascending order:

```
1.  AllLeaves(L,k)  // L: length of the sequence
2.   $\mathbf{a} \leftarrow (1, \dots, 1)$   //  $\mathbf{k}$ : max digit value
3.  while forever  //  $\mathbf{a}$ : array of digits
4.      output  $\mathbf{a}$ 
5.       $\mathbf{a} \leftarrow \text{NextLeaf}(\mathbf{a}, L, k)$ 
6.      if  $\mathbf{a} = (1, \dots, 1)$ 
7.          return
```



# Visit All Leaves: Example

- Moving through all the leaves in order:



# Depth First Search

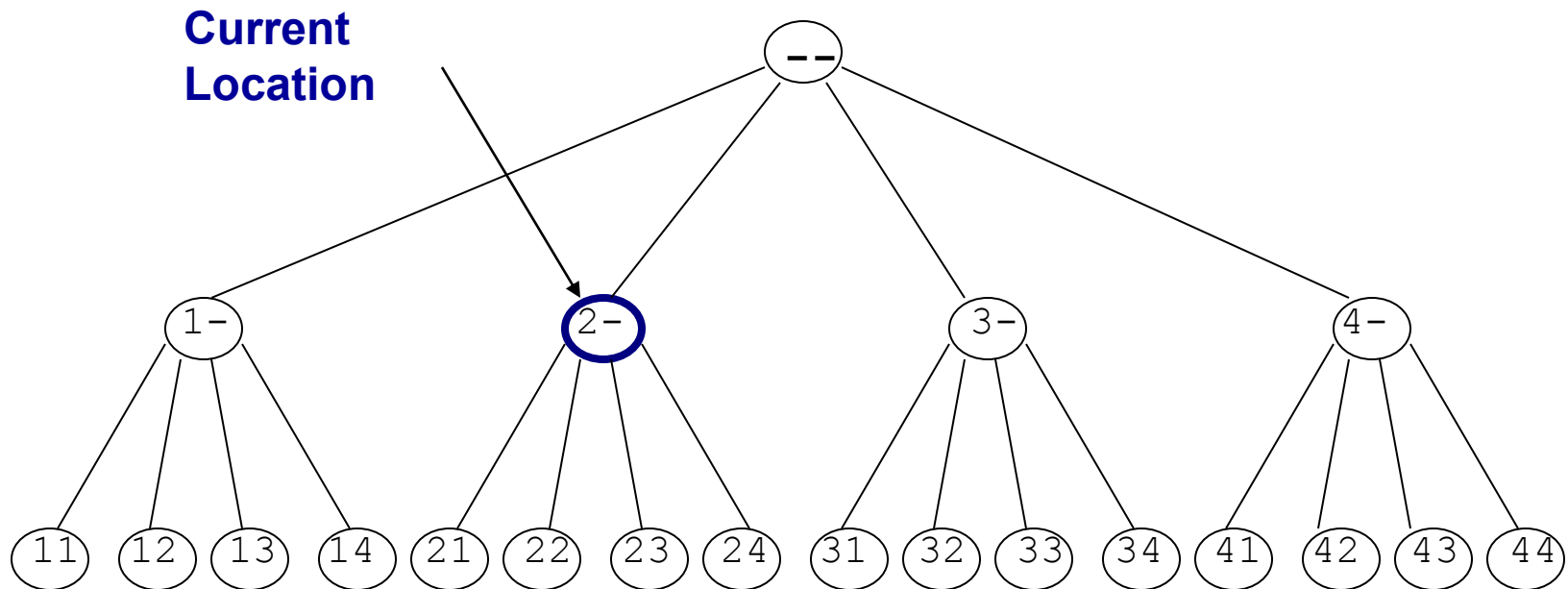
- So we can search leaves
- How about searching all vertices of the tree?
- We can do this with a *depth first* search

# Visit the Next Vertex

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

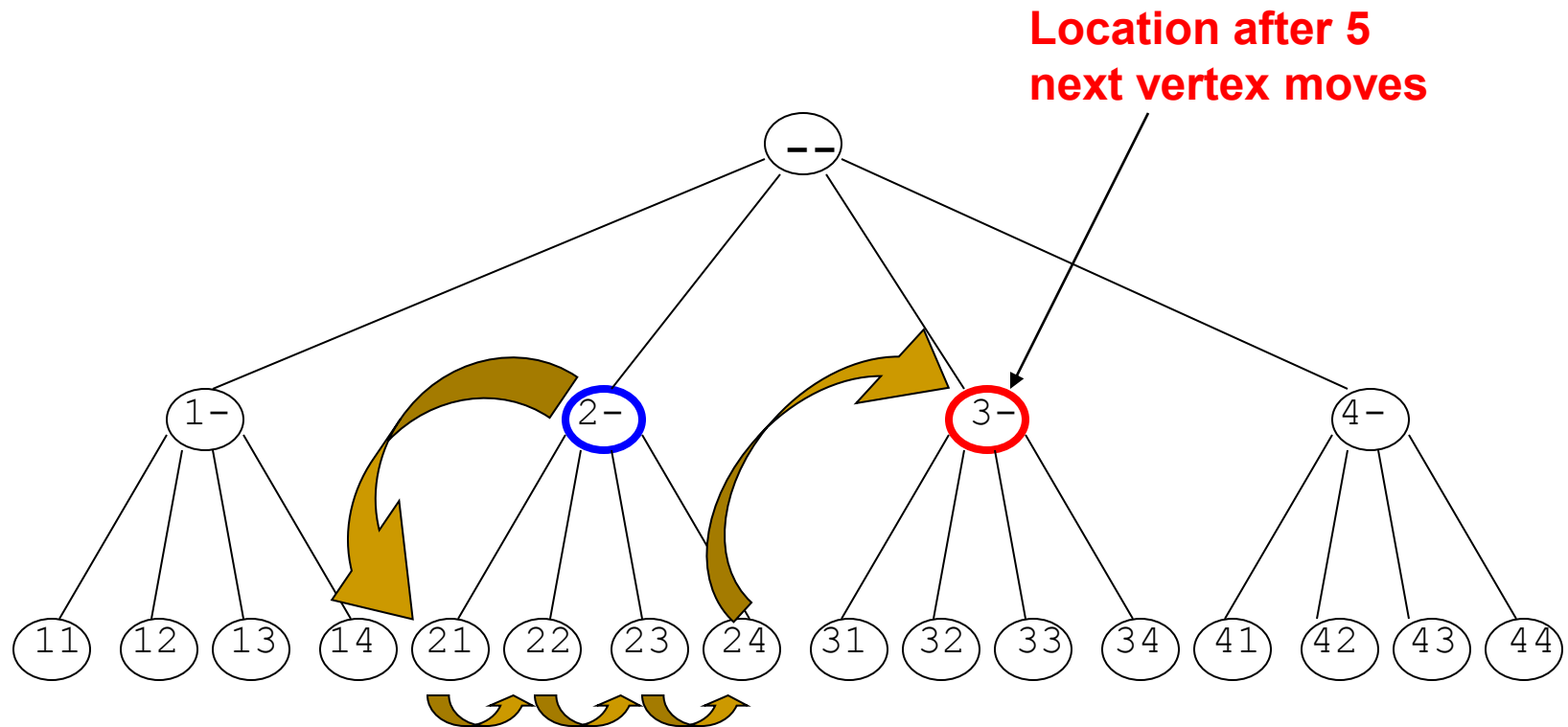
# Example

- Moving to the next vertex:



# Example

- Moving to the next vertices:



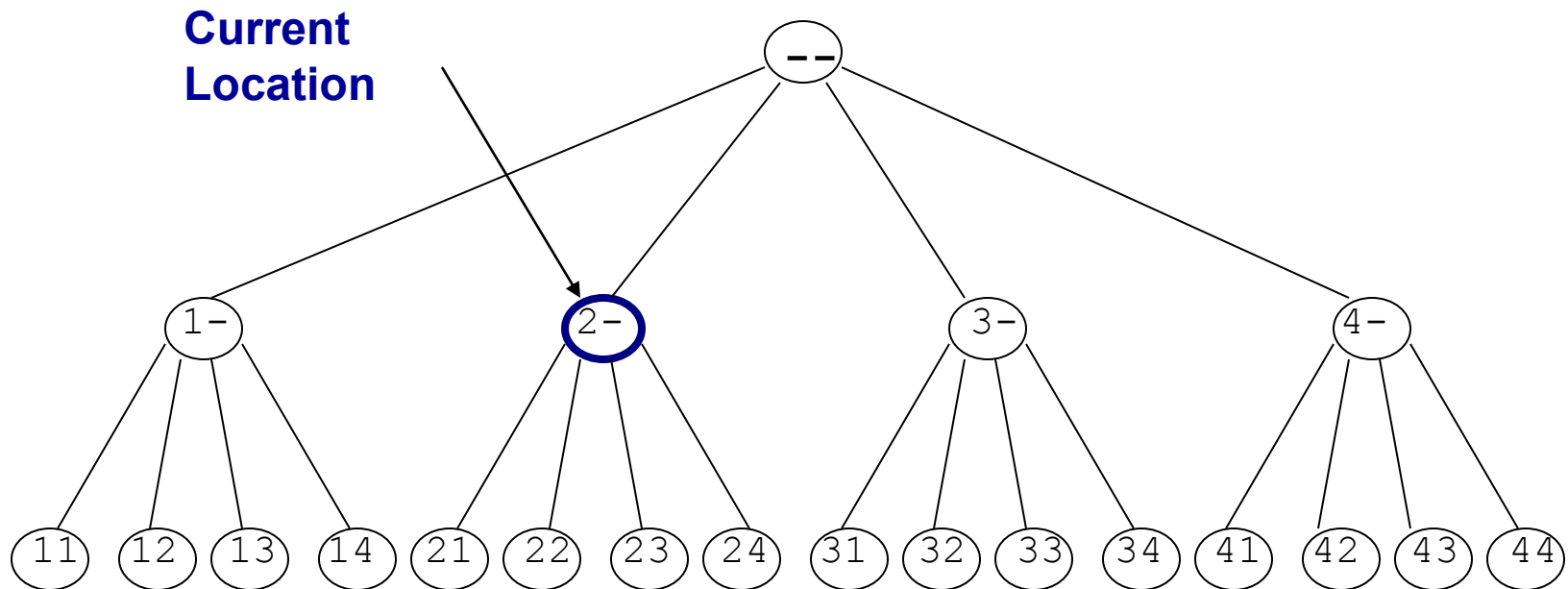
# Bypass Move

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

```
1.  Bypass(a, i, L, k)    // a: array of digits
2.  for  $j \leftarrow i$  to  $l$     //  $i$ : 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)
```

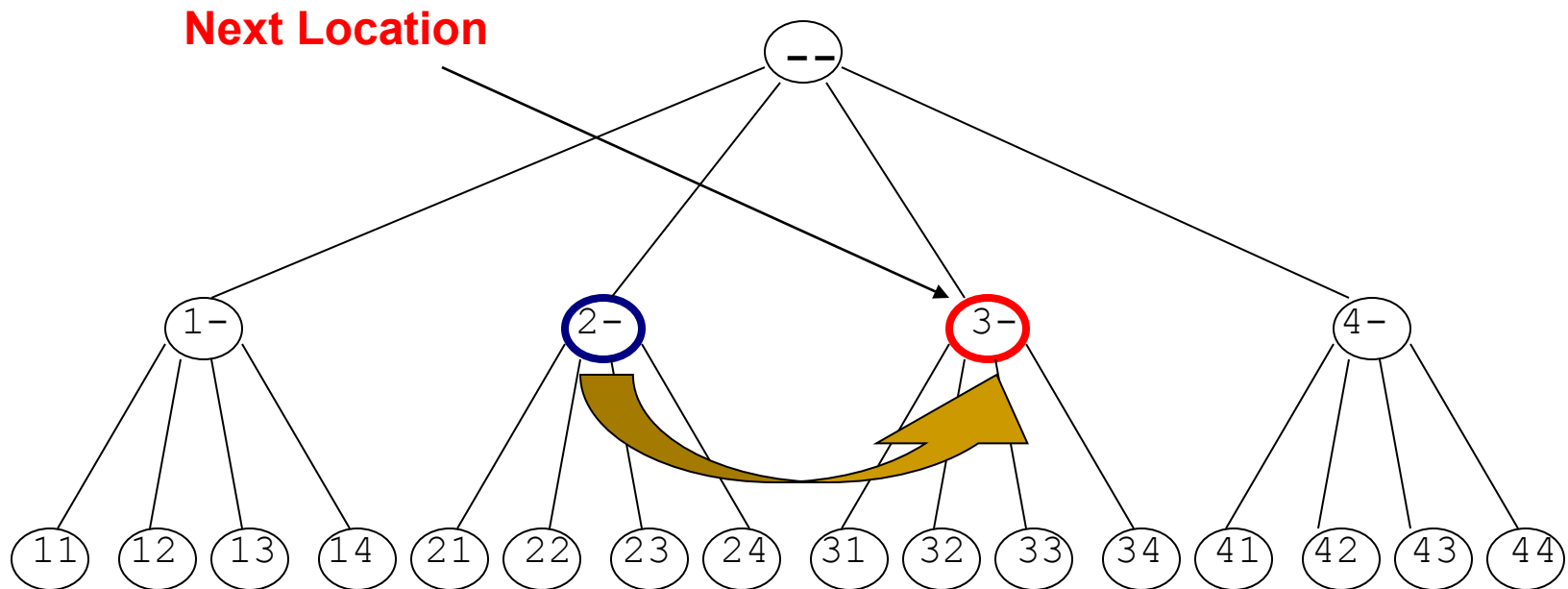
# Bypass Move: Example

- Bypassing the descendants of “2-”:



# Example

- Bypassing the descendants of “2-”:





# Brute Force Search Again

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

# Can We Do Better?

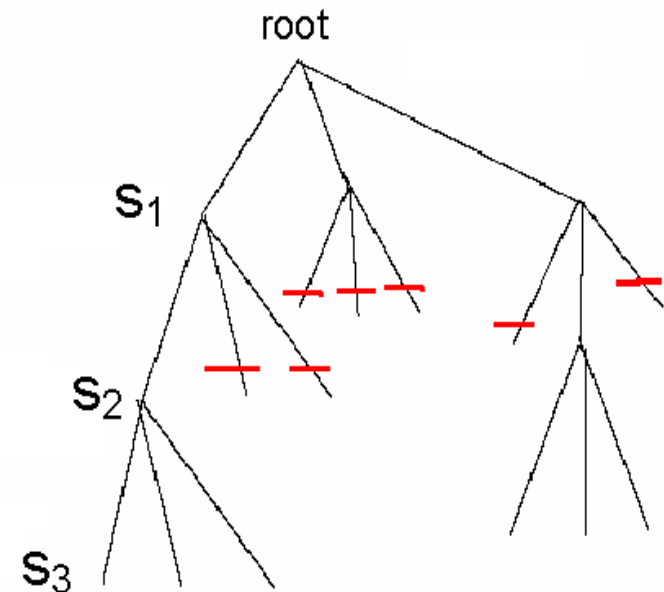
- Sets of  $\mathbf{s}=(s_1, s_2, \dots, s_t)$  may have a weak profile for the first  $i$  positions  $(s_1, s_2, \dots, s_i)$
- Every row of alignment may add at most  $\ell$  to Score
- Optimism: if all subsequent  $(t-i)$  positions  $(s_{i+1}, \dots, s_t)$  add

$$(t - i) * \ell \text{ to } \text{Score}(\mathbf{s}, i, \mathbf{DNA})$$

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

# 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  $(n - \ell + 1)^{t-i}$  leaves
  - Use **NextVertex()** and **ByPass()** to navigate the tree



# Pseudocode for Branch and Bound Motif Search

```
1.  BranchAndBoundMotifSearch(DNA, t, n, l)
2.  s  $\leftarrow$  (1,...,1)
3.  bestScore  $\leftarrow$  0
4.  i  $\leftarrow$  1
5.  while i > 0
6.      if i < t
7.          optimisticScore  $\leftarrow$  Score(s, i, DNA) + (t - i) * l
8.          if optimisticScore < bestScore
9.              (s, i)  $\leftarrow$  Bypass(s, i, n - l + 1)
10.         else
11.             (s, i)  $\leftarrow$  NextVertex(s, i, n - l + 1)
12.         else
13.             if Score(s, DNA) > bestScore
14.                 bestScore  $\leftarrow$  Score(s)
15.                 bestMotif  $\leftarrow$  (s1, s2, s3, ..., st)
16.                 (s, i)  $\leftarrow$  NextVertex(s, i, t, n - l + 1)
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  $(n-\ell+1)^t$  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!

# 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:

$$\text{TotalDistance}(\textit{prefix}, \textit{DNA}) > \textit{BestDistance}$$

there is no use exploring the remaining part of the word

- We can eliminate that branch and BYPASS exploring that branch further

# Bounded Median String Search

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

# Improving the Bounds

- Given an  $\ell$ -mer  $\mathbf{w}$ , divided into two parts at point  $i$ 
  - $\mathbf{u}$  : prefix  $w_1, \dots, w_i$ ,
  - $\mathbf{v}$  : suffix  $w_{i+1}, \dots, w_\ell$
- Find minimum distance for  $\mathbf{u}$  in a sequence
- No instances of  $\mathbf{u}$  in the sequence have distance less than the minimum distance
- Note this doesn't tell us anything about whether  $\mathbf{u}$  is part of any motif. We only get a minimum distance for prefix  $\mathbf{u}$



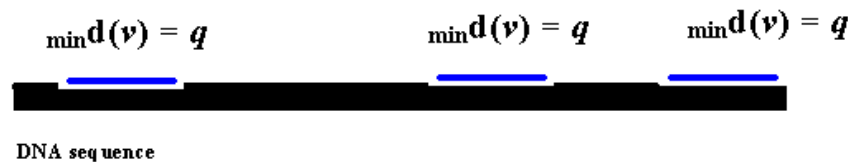
# Improving the Bounds (cont'd)

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

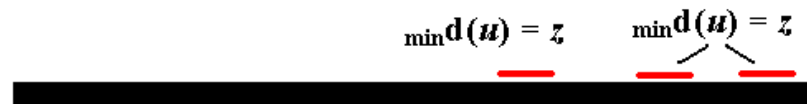
# Better Bounds

Searching for prefix  $V$

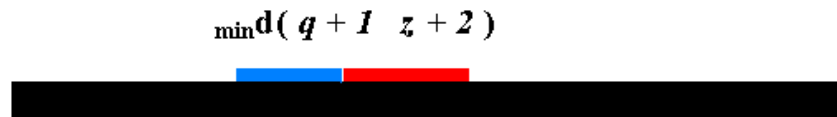
We may find many instances of prefix  $V$  with a minimum distance  $q$



Likewise for  $U$



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



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)

- If  $d(\textit{prefix}) + d(\textit{suffix}) \geq \textit{bestDistance}$ :
  - Motif  $w$  ( $\textit{prefix.suffix}$ ) cannot give a better (lower) score than  $d(\textit{prefix}) + d(\textit{suffix})$
  - In this case, we can **ByPass()**

# Better Bounded Median String Search

```
1. ImprovedBranchAndBoundMedianString(DNA,t,n,l)
2.   s = (1, 1, ..., 1)
3.   bestdistance =  $\infty$ 
4.   i = 1
5.   while i > 0
6.     if i < l
7.       prefix = nucleotide string corresponding to (s1, s2, s3, ..., si)
8.       optimisticPrefixDistance = TotalDistance (prefix, DNA)
9.       if (optimisticPrefixDistance < bestsubstring[ i ])
10.        bestsubstring[ i ] = optimisticPrefixDistance
11.        if (l - i < i)
12.          optimisticSufxDistance = bestsubstring[l - i ]
13.        else
14.          optimisticSufxDistance = 0;
15.        if optimisticPrefixDistance + optimisticSufxDistance  $\geq$  bestDistance
16.          (s, i) = Bypass(s, i, l, 4)
17.        else
18.          (s, i) = NextVertex(s, i, l, 4)
19.      else
20.        word = nucleotide string corresponding to (s1, s2, s3, ..., sl)
21.        if TotalDistance( word, DNA) < bestDistance
22.          bestDistance = TotalDistance(word, DNA)
23.          bestWord = word
24.          (s, i) = NextVertex(s, i, l, 4)
25.   return bestWord
```