

## Exhaustive search

Torgeir R. Hvidsten

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

1

## This lecture

- Restriction enzymes and the **partial digest problem**
- Finding **regulatory motifs** in DNA Sequences
- Exhaustive search methods

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

2

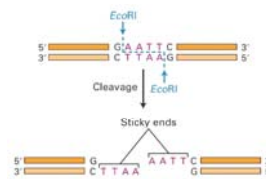
## Restriction enzymes and the partial digest problem

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

3

## Restriction enzymes

- *Hind*II - first restriction enzyme – was discovered accidentally in 1970 while studying how the bacterium *Haemophilus influenzae* takes up DNA from the virus
- Restriction enzymes are used as a defense mechanism by bacteria to break down the DNA of attacking viruses



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

4

## Uses of restriction enzymes

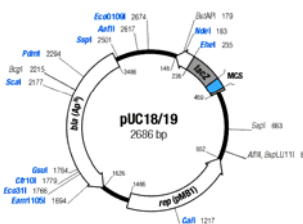
- Recombinant DNA technology (i.e. combining DNA sequences that would not normally occur together)
- Cloning
- cDNA/genomic library construction
- DNA mapping

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

5

## Restriction maps

- A map showing the positions of restriction sites in a DNA sequence
- If the DNA sequence is known then constructing a restriction map is trivial
- 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**



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

6

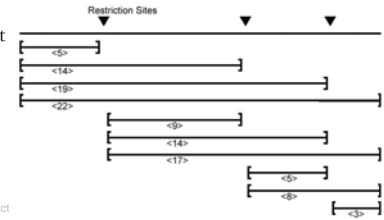
## Measuring length of restriction fragments



- 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

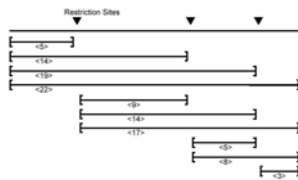
## Partial restriction digest

- 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



## 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
- **Restriction sites:**  $X = \{0, 5, 14, 19, 22\}$
- **Multiset:**  $\downarrow X = \{3, 5, 5, 8, 9, 14, 14, 17, 19, 22\}$



## Partial Digest Problem (PDP)

- $X$  the set of  $n$  integers representing the location of all cuts in the restriction map, including the start and end i.e.  $X = \{x_1 = 0, x_2, \dots, x_n\}$ .
- $n$  the total number of cuts
- $\downarrow X$  the multiset of integers representing lengths of each of the  $n(n-1)/2$  fragments produced from a partial digest i.e.  $\downarrow X = \{x_j - x_i \mid 1 \leq i < j \leq n\}$

**Problem:** Given the multiset  $L$ , find a set  $X$  such that  $\downarrow X = L$

## Partial Digest: Multiple Solutions (I)

It is not always possible to uniquely reconstruct a set  $X$  based only on  $\downarrow X$ .

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 (they are **homometric**).

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:

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

## Partial Digest: Multiple Solutions (II)

	0	1	2	5	7	9	12		0	1	5	7	8	10	12
0		1	2	5	7	9	12	0		1	5	7	8	10	12
1			1	4	6	8	11	1			4	6	7	9	11
2				3	5	7	10	2 <td></td> <td></td> <td></td> <td>2</td> <td>3</td> <td>5</td> <td>7</td>				2	3	5	7
5					2	4	7	5 <td></td> <td></td> <td></td> <td></td> <td>1</td> <td>3</td> <td>5</td>					1	3	5
7						2	5	7 <td></td> <td></td> <td></td> <td></td> <td></td> <td>2</td> <td>4</td>						2	4
9							3	9 <td></td> <td></td> <td></td> <td></td> <td></td> <td></td> <td>2</td>							2
12								12							

## BruteForcePDP

```

BruteForcePDP( $L, n$ )
1   $M \leftarrow$  Maximum element in  $L$ 
2  for every set of  $n - 2$  integers  $0 < x_2 < \dots < x_{n-1} < M$ 
   such that  $x_i \in L$  for  $1 < i < n$ 
3     $X \leftarrow \{0, x_2, \dots, x_{n-1}, M\}$ 
4    Form  $\Delta X$  from  $X$ 
5    if  $\Delta X = L$ 
6      return  $X$ 
7  output "No solution"
    
```

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

13

## Efficiency of BruteForcePDP

- The running time of BruteForcePDP is  $O(|L|^{n-2})$  or, since  $|L| = n(n-1)/2$ ,  $O(n^{2n-4})$
- Note that **without** restricting the elements of  $X$  to elements in  $L$ , the running time would be  $O(M^{n-2})$
- If  $L = \{2, 998, 1000\}$  ( $n = 3, M = 1000$ ), the running time would be very different
- Nonetheless, both algorithms have an exponential running time

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

14

## Branch and Bound Algorithm for PDP

1. Initiation: Add the start ( $0$ ) and end ( $width$ ) point of the sequence to  $X$  (and remove the end point from  $L$ )
2. Find the largest element  $y$  in  $L$
3. See if  $y$  fits on the right or left side of the restriction map by checking whether the other lengths (fragments) it creates are in  $L$
4. If it fits, remove these lengths from  $L$  and add  $y$  (or  $width - y$ ) to  $X$  (if not, backtrack)
5. Go back to step 2 until  $L$  is empty

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

15

## Definitions

Before describing PartialDigest, first define

$$\Delta(y, X)$$

as the multiset of all distances between point  $y$  and all other points in the set  $X$

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

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

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

16

PartialDigest( $L$ )

```

1   $width \leftarrow$  Maximum element in  $L$ 
2  Remove  $width$  from  $L$ 
3   $X \leftarrow \{0, width\}$ 
4  Place( $L, X$ )
    
```

## PartialDigest

Place( $L, X$ )

```

1  if  $L$  is empty
2    output  $X$ 
3    return
4   $y \leftarrow$  Maximum element in  $L$ 
5  if  $\Delta(y, X) \subseteq L$ 
6    Remove lengths  $\Delta(y, X)$  from  $L$  and add  $y$  to  $X$ 
7    Place( $L, X$ )
8  Remove  $y$  from  $X$  and add lengths  $\Delta(y, X)$  to  $L$ 
9  if  $\Delta(width - y, X) \subseteq L$ 
10   Remove lengths  $\Delta(width - y, X)$  from  $L$  and add  $width - y$  to  $X$  and
11   Place( $L, X$ )
12   Remove  $width - y$  from  $X$  and add lengths  $\Delta(width - y, X)$  to  $L$ 
13 return
    
```

17

## An Example

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

Initiation:

$width = 10$ , so delete  $10$  from  $L$  and add  $0$  and  $10$  to  $X$

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

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

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

18

## An Example

$$L = \{2, 2, 3, 3, 4, 5, 6, 7, 8\}$$
$$X = \{0, 10\}$$

$$y = 8$$
$$\Delta(y, X) = \{8, 2\}, \text{ which is a subset of } L$$

$$L = \{2, 3, 3, 4, 5, 6, 7\}$$
$$X = \{0, 8, 10\}$$

## An Example

$$L = \{2, 3, 3, 4, 5, 6, 7\}$$
$$X = \{0, 8, 10\}$$

$$y = 7$$
$$\Delta(y, X) = \{7, 1, 3\}, \text{ which is not a subset of } L$$
$$\Delta(\text{width} - y, X) = \{3, 5, 7\}, \text{ which is a subset of } L$$

$$L = \{2, 3, 4, 6\}$$
$$X = \{0, 3, 8, 10\}$$

## An Example

$$L = \{2, 3, 4, 6\}$$
$$X = \{0, 3, 8, 10\}$$

$$y = 6$$
$$\Delta(y, X) = \{6, 3, 2, 4\}, \text{ which is a subset of } L$$

$$L = \{\}$$
$$X = \{0, 3, 6, 8, 10\}$$

## An Example

$$L = \{\}$$
$$X = \{0, 3, 6, 8, 10\}$$

$L$  is empty!  
Output  $X$

And continue searching for more solutions ...

## An Example

$$L = \{2, 3, 4, 6\}$$
$$X = \{0, 3, 8, 10\}$$

$$y = 6$$
$$\Delta(y, X) = \{6, 3, 2, 4\}, \text{ which is a subset of } L$$
$$\Delta(\text{width} - y, X) = \{4, 1, 4, 6\}, \text{ which is not a subset of } L$$

**Backtrack!**

## An Example

$$L = \{2, 3, 3, 4, 5, 6, 7\}$$
$$X = \{0, 8, 10\}$$

$$y = 7$$
$$\Delta(y, X) = \{7, 1, 3\}, \text{ which is not a subset of } L$$
$$\Delta(\text{width} - y, X) = \{3, 5, 7\}, \text{ which is a subset of } L$$

**Backtrack!**

## An Example

$L = \{2, 2, 3, 3, 4, 5, 6, 7, 8\}$   
 $X = \{0, 10\}$

$y = 8$

$\neg(y, X) = \{8, 2\}$ , which is a subset of  $L$

$\neg(\text{width} - y, X) = \{2, 8\}$ , which is a subset of  $L$

$L = \{2, 3, 3, 4, 5, 6, 7\}$   
 $X = \{0, 2, 10\}$

## An Example

$L = \{2, 3, 3, 4, 5, 6, 7\}$   
 $X = \{0, 2, 10\}$

$y = 7$

$\neg(y, X) = \{7, 5, 3\}$ , which is a subset of  $L$

$L = \{2, 3, 4, 6\}$   
 $X = \{0, 2, 7, 10\}$

## An Example

$L = \{2, 3, 4, 6\}$   
 $X = \{0, 2, 7, 10\}$

$y = 6$

$\neg(y, X) = \{6, 4, 1, 4\}$ , which is not a subset of  $L$

$\neg(\text{width} - y, X) = \{4, 2, 3, 6\}$ , which is a subset of  $L$

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

## An Example

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

$L$  is empty!

Output  $X$

And continue searching for more solutions ...

## An Example

$L = \{2, 3, 3, 4, 5, 6, 7\}$   
 $X = \{0, 2, 10\}$

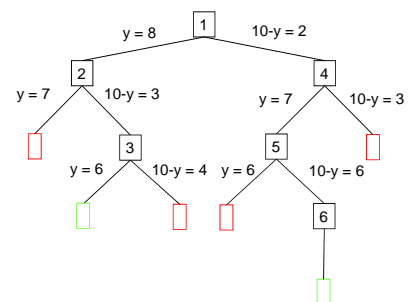
$y = 7$

$\neg(y, X) = \{7, 5, 3\}$ , which is a subset of  $L$

$\neg(\text{width} - y, X) = \{3, 1, 7\}$  which is not a subset of  $L$

**Backtrack and finish!**

## Recursion tree



## Analyzing PartialDigest (I)

- Let  $T(n)$  be the time PartialDigest takes to place  $n$  cuts
  - No branching case:  $T(n) = T(n-1) + O(n)$
  - Branching case:  $T(n) = 2T(n-1) + O(n)$
- The "No branching case" is quadratic  $O(n^2)$  (like SelectionSort)
- The "Branching case" is exponential  $O(2^n)$

## Analyzing PartialDigest (II)

$$T(n) = 2T(n-1) + O(n)$$

$$T(1) = 1$$

$$O(1) = 1$$

$$T(n) + O(n) = 2T(n-1) + O(n) + O(n) = 2(T(n-1) + O(n))$$

$$\text{Let } U(n) = T(n) + O(n)$$

$$U(n) = 2U(n-1)$$

$$U(1) = 2$$

This gives the sequence: 2, 4, 8, 16, 32, 64, ...  $\rightarrow U(n) = 2^n$   
 $T(n) = 2^n - O(n)$

## Finding regulatory motifs in DNA sequences

```
atgaccgggactactgatccgtatttggcctaggcgtacacattagataaacgtatgaagtacgttagactcggcgcgcg
accctatTTTTTggcagatttagtaccctggaasaaatttggatcaaacctttccgaaactggcctaagggtaca
tgaatccccgggatgacttAAAAAAGGGGGGgctctcccgatttTgaatgtaggatcattccagggtccga
gctgagaattggatgacctgtgaagtTttccacgaatcgcaaccaacgcggaccocaaagcagaccgataaggaga
tccctttgcgtaagtgcggaggcTggttagtgggaagcccaacggactaatggccacttagtccacttag
gtcaatcatgtcttTgaatggattTtaactgaggcattagccctTggcaccocaaatcagTtggcgcgcgca
cgtttTggccctgttagagcccccTactgatgaaactTcaatttagagagcTaatctatcgTgctgtTcat
aactTgattTggttTgaatgctTgggcacatacaaggagctTccttatcagTtaagtctgtatgacactatgta
ttggccattTgctaaagcccaactTgacaatTgaagatagactctTgcaatTcaactTatgccaagcgaaggag
ctggtgacacgacagattcttaactgattagctcgtctccgggtctaagtagcacgaagctctgggtactgataga
```

## Implanting motif AAAAAAAGGGGGG

```
atgaccgggactactgataAAAAAAGGGGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgcgcg
accctatTTTTTggcagatttagtaccctggaasaaatttggatcaaacctttccgaaataAAAAAAGGGGGGg
tgaatccccgggatgacttAAAAAAGGGGGGgctctcccgatttTgaatgtaggatcattccagggtccga
gctgagaattggatgAAAAAAGGGGGGgTccaacgaatcgcaaccaacgcggaccocaaagcagaccgataaggaga
tccctttgcgtaagtgcggaggcTggttagtgggaagcccaacggactaatAAAAAAGGGGGGgcttag
gtcaatcatgtcttTgaatggattAAAAAAGGGGGGgagccctTggcaccocaaatcagTtggcgcgcgca
cgtttTggccctgttagagcccccTAAAAAAGGGGGGgcaatttagagagcTaatctatcgTgctgtTcat
aactTgattAAAAAAGGGGGGgctTgggcacatacaaggagctTccttatcagTtaagtctgtatgacactatgta
ttggccattTgctaaagcccaactTgacaatTgaagatagactctTgcaatAAAAAAGGGGGGgaccgaaggag
ctggtgacacgacagattcttaactgattagctcgtctccgggtctaagtagcacgaagcttAAAAAAGGGGGGg
```

## Where is the implanted motif?

```
atgaccgggactactgataaaaaagggggggcgtacacattagataaacgtatgaagtacgttagactcggcgcgcg
accctatTTTTTggcagatttagtaccctggaasaaatttggatcaaacctttccgaaataaaaaaggggggg
tgaatccccgggatgacttaaaaaaggggggggctctcccgatttTgaatgtaggatcattccagggtccga
gctgagaattggatgaaaaaagggggggTccaacgaatcgcaaccaacgcggaccocaaagcagaccgataaggaga
tccctttgcgtaagtgcggaggcTggttagtgggaagcccaacggactaataaaaaagggggggcttag
gtcaatcatgtcttTgaatggattaaaaaagggggggccctTggcaccocaaatcagTtggcgcgcgca
cgtttTggccctgttagagcccccTaaaaaagggggggcaatttagagagcTaatctatcgTgctgtTcat
aactTgattaaaaaagggggggctgggcacatacaaggagctTccttatcagTtaagtctgtatgacactatgta
ttggccattTgctaaagcccaactTgacaatTgaagatagactctTgcaataaaaaaagggggggaccgaaggag
ctggtgacacgacagattcttaactgattagctcgtctccgggtctaagtagcacgaagcttaaaaaaaggggggg
```

## Implanting motif AAAAAAGGGGGG with four random mutations

```
atgaccgggactctgagAAAGAGGtGGGgggtacacattagataaacgtatgaagtcagttagactcggcggccgcg
accctatTTTTTgagcagtttagtgcctggaaasaaatttgatgacaaacttttccgataCAATAAAAGGGGGa
tgatccctgggatgacttAAAAAAGGGGGtGtctcctccgattttgaaatgtaggcatctccaggggtccga
gctgagaattggatg-AAAAAAGGGGGtTtccacgcaetcggaaccaacggcccaasaggcaagcogtaasaggaga
tccctttgcggttaetgcccgggggctggttacgtagggaagcccaacggaacttaataATAAAGGaaGGGcttatag
gtcaatcagttcttgtgaaatgattAAAGAGGGGGtGGaccgcttggccaccaaatcagttggggcggcga
cggtttggcccttggtagggcccccgtATAAAGAGGGGGcaattatgagagcctaactatcgctgctgttcat
aaacttgagttAAAAAAGGGGGcctcctgggacatacagaagggtcttcttcatgatttaetgctgtagcactatgta
ttggccattggctaaagcccaacttgacaatggagatagatccttgcataAAAAAGGGGGcaccgaaaggag
ctggtgagcaagcagacttctacgtgacttctcgtctccgggactaatagacagagcttActAAAAAGGGGGa
```

## Where is the motif?

```
atgaccgggactctgagaaaggttggggggtacacattagataaacgtatgaagtcagttagactcggcggccgcg
accctatTTTTTgagcagtttagtgcctggaaasaaatttgatgacaaactttccgataCAATAAAAGGGGGa
tgatccctgggatgacttAAAAAAGGGGGtGtctcctccgattttgaaatgtaggcatctccaggggtccga
gctgagaattggatgcaasaaagggttgcacgcaetcggaaccaacggcccaasaggcaagcogtaasaggaga
tccctttgcggttaetgcccgggggctggttacgtagggaagcccaacggaacttaataataaaggaggcttatag
gtcaatcagttcttgtgaaatgattaaacaatggggctggggccgcttggccaccaaatcagttggggcggcga
cggtttggcccttggtagggcccccgtataaacaaggaggcccaattatgagagcctaactatcgctgctgttcat
aaacttgagttAAAAAAGGGGGcctcctgggacatacagaagggtcttcttcatgatttaetgctgtagcactatgta
ttggccattggctaaagcccaacttgacaatggagatagatccttgcataAAAAAGGGGGcaccgaaaggag
ctggtgagcaagcagacttctacgtgacttctcgtctccgggactaatagacagagcttActAAAAAGGGGGa
```

## Why finding motif is difficult

```
atgaccgggactctgagAAAGAGGtGGGgggtacacattagataaacgtatgaagtcagttagactcggcggccgcg
accctatTTTTTgagcagtttagtgcctggaaasaaatttgatgacaaactttccgataCAATAAAAGGGGGa
tgatccctgggatgacttAAAAAAGGGGGtGtctcctccgattttgaaatgtaggcatctccaggggtccga
gctgagaattggatg-AAAAAAGGGGGtTtccacgcaetcggaaccaacggcccaasaggcaagcogtaasaggaga
tccctttgcggttaetgcccgggggctggttacgtagggaagcccaacggaacttaataATAAAGGaaGGGcttatag
gtcaatcagttcttgtgaaatgattAAAGAGGGGGtGGaccgcttggccaccaaatcagttggggcggcga
cggtttggcccttggtagggcccccgtATAAAGAGGGGGcaattatgagagcctaactatcgctgctgttcat
aaacttgagttAAAAAAGGGGGcctcctgggacatacagaagggtcttcttcatgatttaetgctgtagcactatgta
ttggccattggctaaagcccaacttgacaatggagatagatccttgcataAAAAAGGGGGcaccgaaaggag
ctggtgagcaagcagacttctacgtgacttctcgtctccgggactaatagacagagcttActAAAAAGGGGGa
```



## Gene regulation

- A microarray 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 a 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

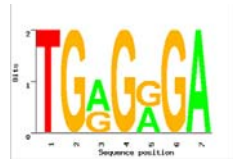
## Regulatory regions

- Every gene contains a **regulatory region** (RR) typically stretching 100-1000 bp upstream of the transcriptional start site
- Located somewhere 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 - and recruiting the DNA polymerase

## 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  
TGAGGGGA



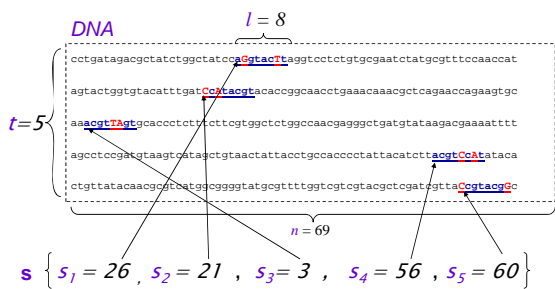
## The motif finding problem

Given a random sample of DNA sequences:

```
cctgatagacgctatctggctatccacgtacgtaggtcctctgtgcaatctatgctttccaacct
agtactgggtacatttgatagctacgtacaccggcaacctgaaacaaacgctcagaaccagaagtgc
aaacgtacgtcacccctcttctctggctctggccaacgagggcgtatgataagacgaaaattt
agcctccgatgtaagtcatagtctgtaactattacctgccaccctattacatctacgtacgtataca
ctgttatacaacgctcatggcggggatgctgttttggctgctacgctcgatcgttaacgtacgtc
```

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

## Parameters



## Definitions

- $t$  number of sample DNA sequences
- $n$  length of each DNA sequence
- DNA** sample of DNA sequences ( $t \times n$  array)
- $l$  length of the motif ( $l$ -mer)
- $s_i$  starting position of an  $l$ -mer in sequence  $i$
- $s = (s_1, s_2, \dots, s_t)$  array of motif starting positions

## 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
    
```

➤ Line up the patterns by their start indexes

$s = (s_1, s_2, \dots, s_t)$

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
    
```

➤ 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
- We need to introduce a scoring function to compare different motifs and choose the “best” one.



## Scoring motifs: consensus score

➤ Given  $s = (s_1, \dots, s_l)$  and  $DNA$ :

➤  $\text{Score}(s, DNA) =$

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

		$l$									
		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

➤ **Goal:** Given a set of DNA sequences, find a set of  $l$ -mers, one from each sequence, that maximizes the consensus score

➤ **Input:** A  $t \times n$  matrix  $DNA$ , and  $l$ , the length of the pattern to find

➤ **Output:** An array of  $t$  starting positions  $s = (s_1, s_2, \dots, s_t)$  maximizing  $\text{Score}(s, DNA)$

## BruteForceMotifSearch

$\text{BruteForceMotifSearch}(DNA, t, n, l)$

```

1 bestScore ← 0
2 for each  $s = (s_1, s_2, \dots, s_l)$  from  $(1, 1 \dots, 1)$  to  $(n-l+1, \dots, n-l+1)$ 
3   if  $(\text{Score}(s, DNA) > \text{bestScore})$ 
4     bestScore ←  $\text{Score}(s, DNA)$ 
5     bestMotif ←  $(s_1, s_2, \dots, s_l)$ 
6 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  $l$  operations, so complexity is  $l(n-l+1)^t = O(ln^t)$
- That means that for  $t = 8$ ,  $n = 1000$ , and  $l = 10$  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

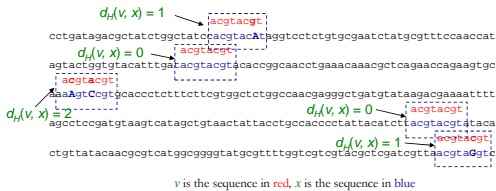
## Hamming Distance

- **Hamming distance:**
  - $d_H(v, w)$  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: Example

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



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

## The median string problem

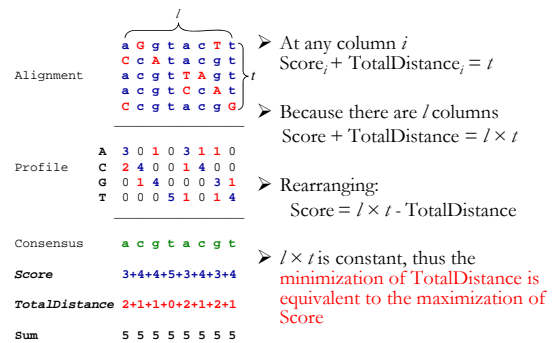
- Goal:** Given a set of DNA sequences, find a median string
- Input:** A  $t \times n$  matrix  $DNA$ , and  $l$ , the length of the pattern to find
- Output:** A string  $v$  of  $l$  nucleotides that **minimizes** TotalDistance( $v, DNA$ ) over all strings of that length

## Median string search algorithm

BruteForceMedianStringSearch ( $DNA, t, n, l$ )

- $bestWord \leftarrow AAA...A$
- $bestDistance \leftarrow \infty$
- for** each  $l$ -mer  $v$  **from**  $AAA...A$  to  $TTT...T$
- if** TotalDistance( $v, DNA$ ) <  $bestDistance$
- $bestDistance \leftarrow TotalDistance(v, DNA)$
- $bestWord \leftarrow v$
- return**  $bestWord$

## Motif finding problem = median string problem



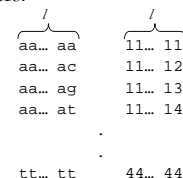
## Motif finding problem vs. median string problem

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  $(n - l + 1)^l$  combinations
- The median string problem needs only to examine all  $4^l$  combinations for  $v$ .

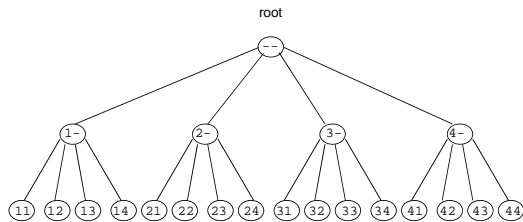
## Structuring the search

For the median string problem we need to consider all  $4^l$  possible  $l$ -mers:



How to organize this search?

## Search Tree



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

61

## 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?

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

62

## Visit the Next Leaf

Let  $\mathbf{a}$  be an array of digits:  $\mathbf{a} = \{a_1, a_2, \dots, a_l\}$  where  $a_i \in [1, k], 1 \leq i \leq l$ .

Given a current leaf  $\mathbf{a}$ , we need to compute the “next” leaf:

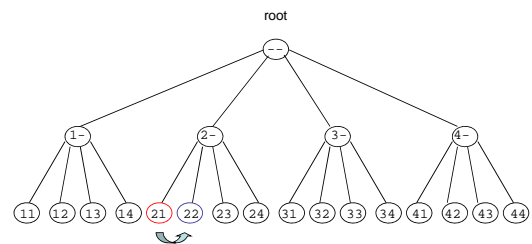
$\text{NextLeaf}(\mathbf{a}, l, k)$

- 1 **for**  $i \leftarrow l$  **to** 1
- 2   **if**  $a_i < k$
- 3      $a_i \leftarrow a_i + 1$
- 4   **return**  $\mathbf{a}$
- 5    $a_i \leftarrow 1$
- 6 **return**  $\mathbf{a}$

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

63

## NextLeaf: Example



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

64

## Visit all leaves

Printing all permutations in ascending order:

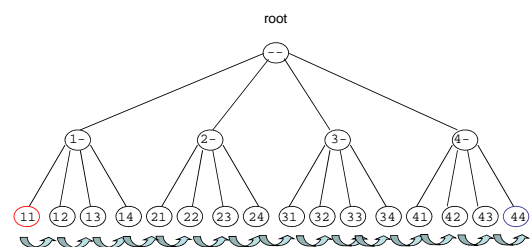
$\text{AllLeaves}(L, k)$

- 1  $\mathbf{a} \leftarrow (1, 1, \dots, 1)$
- 2 **while** forever
- 3   output  $\mathbf{a}$
- 4    $\mathbf{a} \leftarrow \text{NextLeaf}(\mathbf{a}, L, k)$
- 5   **if**  $\mathbf{a} = (1, 1, \dots, 1)$
- 6   **return**

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

65

## Visit all leaves: Example



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

66

## Visit the next vertex

We can search through all vertices of the tree with a *depth first* search

$i$  is the prefix length

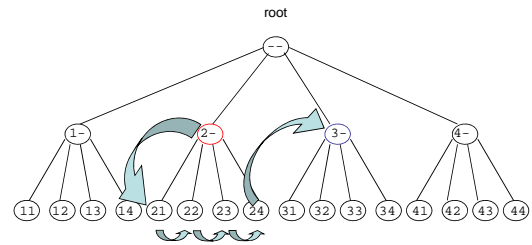
```

NextVertex( $a, i, L, k$ )
1  if  $i < L$ 
2     $a_{i+1} \leftarrow 1$ 
3    return ( $a, i+1$ )
4  else
5    for  $j \leftarrow L$  to 1
6      if  $a_j < k$ 
7         $a_j \leftarrow a_j + 1$ 
8        return ( $a, j$ )
9  return ( $a, 0$ )
    
```

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

67

## Visit the next vertex: Example (five moves)



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

68

## Bypass Move

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

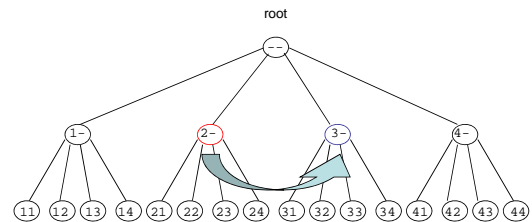
```

Bypass( $a, i, L, k$ )
1  for  $j \leftarrow i$  to 1
2    if  $a_j < k$ 
3       $a_j \leftarrow a_j + 1$ 
4      return ( $a, j$ )
5  return ( $a, 0$ )
    
```

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

69

## Bypass Move: Example



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

70

## Brute force search again

```

BruteForceMedianStringSearchAgain( $DNA, l, n, l$ )
1   $s \leftarrow (1, 1, \dots, 1)$ 
2   $bestDistance \leftarrow \infty$ 
3  while forever
4     $s \leftarrow \text{NextLeaf}(s, l, \mathcal{A})$ 
5     $word \leftarrow$  Nucleotide string corresponding to  $(s_1, s_2, \dots, s_l)$ 
6    if ( $\text{TotalDistance}(word, DNA) < bestDistance$ )
7       $bestDistance \leftarrow \text{TotalDistance}(word, DNA)$ 
8       $bestWord \leftarrow word$ 
9  if  $s = (1, 1, \dots, 1)$ 
10   return  $bestWord$ 
    
```

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

71

## Can We Do Better?

- Let  $\text{TotalDistance}(\text{prefix}, DNA)$  be the distance for a nucleotide string corresponding to  $(s_1, s_2, \dots, s_l)$
- Note that if the total distance for a prefix is greater than that for the best word so far:

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

there is no use exploring the remaining part of the word

- Use `ByPass()`!

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

72

## Bounded Median String Search

```
BranchAndBoundMedianStringSearch(DNA, l, n, l)
1  s ← (1, ..., 1)
2  bestDistance ← ∞
3  i ← 1
4  while i > 0
5    if i < l
6      prefix ← Nucleotide string corresponding to (s1, s2, ..., si)
7      optimisticDistance ← TotalDistance(prefix, DNA)
8      if optimisticDistance > bestDistance
9        (s, i) ← Bypass(s, i, l, l)
10     else
11       (s, i) ← NextVertex(s, i, l, l)
12   else
13     word ← Nucleotide string corresponding to (s1, s2, ..., sl)
14     if TotalDistance(s, DNA) < bestDistance
15       bestDistance ← TotalDistance(word, DNA)
16       bestWord ← word
17     (s, l) ← NextVertex(s, l, l, l)
18   return bestWord
```

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

73

## Running time

- As usual with branch-and-bound algorithms, there is **no improved running time in the worst case**
- However, it often results in a practical speedup

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

74