# Algorithms in Comp. Bio

Partial Digest Problem
(cont.)
Profiles & Motif Finding

# Recap:

- Essay due Thursday
- New HW, due next Thursday:
  will cover 1$^{st}$ part of
  Chapter 4

  ### Recall:
  - use whatever, but <u>cite</u>
    <u>source</u> + put in <u>your</u>
    <u>own</u> words

  - groups are fine : rule of
    thumb is to write up
    late w/out notes, +
    only use them if you
    need to

- A couple of you had <u>questions</u> +
  suggestions about future topics -

  please email me a reminder!

Now - Ch 4, on Exhaustive Search...

# Last time: Notation

Dfn: A multiset :

ex: $\{2, 2, 2, 3, 3, 4, 5\}$
$\{2_3, 3_2, 4, 5\}$

Dfn: If $X$ is a set of $n$ points on a line segment,
$$\Delta X = \{x_i - x_j : 1 \leq i < j \leq n\}$$

Aside: How big is $\Delta X$?
$$\binom{n}{2} = \frac{n(n-1)}{2} = \frac{n!}{2!(n-2)!}$$

Ex: Let $X = \{0, 2, 4, 7, 10\}$.
$$\Delta X = \{2, 2, 3, 3, 4, 5, 6, 7, 8, 10\}$$
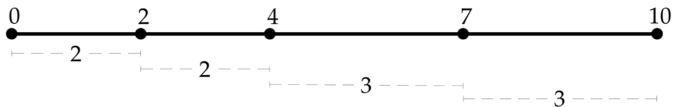$$= \{2_2, 3_2, 4, 5, 6, 7, 8, 10\}$$

**Partial Digest Problem**:
*Given all pairwise distances between points on a line, reconstruct the positions of those points.*
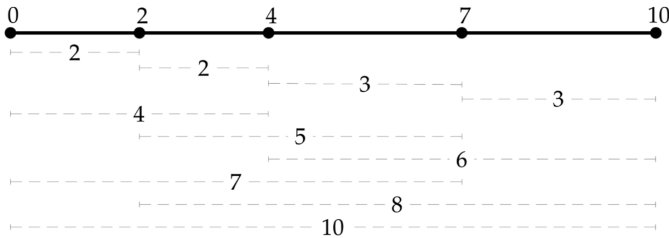
**Input:** The multiset of pairwise distances $L$, containing $\binom{n}{2}$ integers.

**Output:** A set $X$, of $n$ integers, such that $\Delta X = L$

Why?



(a) Complete digest.

(b) Partial digest.

**Figure 4.1** Different methods of digesting a DNA molecule. A complete digest produces only fragments between consecutive restriction sites, while a partial digest yields fragments between any two restriction sites. Each of the dots represents a restriction site.

# Two (slow) solutions:

BRUTEFORCEPDP($L, n$)
1   $M \leftarrow$ maximum element in $L$
2   **for** every set of $n - 2$ integers $0 < x_2 < \cdots < x_{n-1} < M$
3         $X \leftarrow \{0, x_2, \ldots, x_{n-1}, M\}$
4         Form $\Delta X$ from $X$
5         **if** $\Delta X = L$
6                 **return** $X$
7   **output** "No Solution"

$\sim O(M^{n-2})$

+ observe: do we really need
   to try every value $< M$?
Since $0$ is in $X$, if some
        $x \notin L$, then $x \notin X$.

So:

ANOTHERBRUTEFORCEPDP($L, n$)
1   $M \leftarrow$ maximum element in $L$
2   **for** every set of $n - 2$ integers $0 < x_2 < \cdots < x_{n-1} < M$ from $L$
3         $X \leftarrow \{0, x_2, \ldots, x_{n-1}, M\}$
4         Form $\Delta X$ from $X$
5         **if** $\Delta X = L$
6                 **return** $X$
7   **output** "No Solution"

$\sim O(L^{n-2})$

A better way [Skiena 1990]:

Include 0 & largest value in L.
<span style="color:red">remove M from</span> <span style="color:red">L</span>
Consider the next largest "M
  value, called $\delta$.

Where could $\delta$ be?

X


$$M-\delta$$
0
$\delta$

<span style="color:red">b/c this was empty (not in L)</span>

<span style="color:red">~~max{$\ell \in L$}~~</span>
<span style="color:red">~~M~~</span>

Then what?
<span style="color:red">remove $\delta$, check</span>
<span style="color:red">$(0, M-\delta) \in L$</span>

Ex: $\Delta X = L = \{ \cancel{2}, \cancel{2}, \cancel{3}, \cancel{5}, \cancel{4}, \cancel{5}, \cancel{6}, \cancel{7}, \cancel{8}, \cancel{10} \}$

$$|L| = 10 = \binom{n}{2} = \frac{n(n-1)}{2}$$

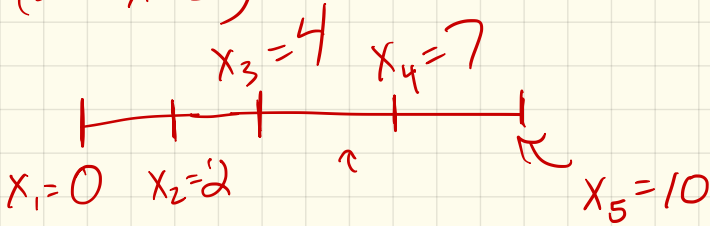$$= \frac{n^2}{2} - \frac{n}{2}$$

solve...

$$n = 5$$

Build $X$:
(s.t. $\Delta X = L$)

$x_3 = 4$ $\quad$ $x_4 = 7$



$x_1 = 0$ $\quad$ $x_2 = 2$ $\qquad\qquad$ $x_5 = 10$

$\delta_1 = \cancel{2}$, so $x_2 = 2$

now $\delta_2 = 7$

now $\delta_3 = 6$

$\vdots$

# Formal algorithm:

PARTIALDIGEST($L$)
1  $width \leftarrow$ Maximum element in $L$
2  DELETE($width, L$) ← *remove $\overline{M}$ from $L$*
3  $X \leftarrow \{0, width\}$
4  PLACE($L, X$)

→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      Add $y$ to $X$ and remove lengths $\Delta(y, X)$ from $L$
  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      Add $width - y$ to $X$ and remove lengths $\Delta(width - y, X)$ from $L$
 11   →PLACE($L, X$)
 12      Remove $width - y$ from $X$ and add lengths $\Delta(width - y, X)$ to $L$
 13  **return**

*$|X|$ is one bigger*

# Note: 
- Recursive!
- Undoes mistakes along the way
- Lists <u>all</u> sets $X$ s.t. $\Delta X = L$

# Runtime:

Worst case:

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

2 recursive calls

check all $\Delta$'s in $L$, remove

(Towers of Hanoi, but worse — uses $O(1)$)

$$O(n \cdot 2^n)$$

If only one viable alternative, then considerably faster $\leftarrow O(n^2)$
(But both can be viable!)

This works much faster in practice, but polynomial time algorithms didn't come until 2002.

# Shifting to another problem:
## DNA profiles & motifs : Idea

Frequent (or rare) substrings
may correspond to regulatory
motifs in DNA. *ℓ-mer*

Picture :

```
CGGGGCTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAACCAAAGCGGACAAA
GGGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCTC
CTGCTGTACAACTGAGATCATGCTGCTTCAAC
TACATGATCTTTTGTGGATGAGGGAATGATGC
```

(a) Seven random sequences.

*8-mer*

```
CGGGGCTATGCAACTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAAATGCAACTCCAAAGCGGACAAA
GGGATGCAACTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGATGCAACTCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCATGCAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCAACTTTCAAC
TACATGATCTTTTGATGCAACTTGGATGAGGGAATGATGC
```

(b) The same DNA sequences with the implanted pattern ATGCAACT.

Brute force?

However, hard to spot (when
not underlined)!

```
CGGGGCTATGCAACTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAAATGCAACTCCAAAGCGGACAAA
GGATGCAACTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGATGCAACTCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCATGCAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCAACTTTCAAC
TACATGATCTTTTGATGCAACTTGGATGAGGGAATGATGC
```

(c) Same as (b), but hiding the implant locations. Suddenly this problem looks difficult to solve.

Even worse: DNA mutates!

```
CGGGGCTATcCAgCTGGGTCGTCACATTCCCCTTTCGATA
TTTGAGGGTGCCCAATAAggGCAACTCCAAAGCGGACAAA
GGATGgAtCTGATGCCGTTTGACGACCTAAATCAACGGCC
AAGGAaGCAACcCCAGGAGCGCCTTTGCTGGTTCTACCTG
AATTTTCTAAAAAGATTATAATGTCGGTCCtTGgAACTTC
CTGCTGTACAACTGAGATCATGCTGCATGCcAtTTTCAAC
TACATGATCTTTTGATGgcACTTGGATGAGGGAATGATGC
```

(d) Same as (b), but with the implanted pattern ATG-CAACT randomly mutated in two positions; no two implanted instances are the same. If we hide the locations as in (c), the difficult problem becomes nearly impossible.

*Formalize:* • t DNA sequences, ℓ-mers, n nucleotides each
• Select a position in each: $(s_1, s_2, ..., s_t)$
+ $1 \leq s_i \leq n - \ell + 1$

```
                        CGGGGCTATcCAgCTGGGTCGTCACATTCCCCTT...
         TTTGAGGGTGCCCAATAAggGCAACTCCAAAGCGGACAAA
                   GGATGgAtCTGATGCCGTTTGACGACCTA...
               AAGGAaGCAACcCCAGGAGCGCCTTTGCTGG...
AATTTTCTAAAAAGATTATAATGTCGGTCCtTGgAACTTC
    CTGCTGTACAACTGAGATCATGCTGCATGCcAtTTTCAAC
         TACATGATCTTTTGATGgcACTTGGATGAGGGAATGATGC
```

(a) Superposition of the seven highlighted 8-mers from figure 4.2 (d).

ℓ

**Alignment Matrix**

|     | A | T | C | C | A | G | C | T |
|-----|---|---|---|---|---|---|---|---|
|     | G | G | G | C | A | A | C | T |
|     | A | T | G | G | A | T | C | T |
|     | A | A | G | C | A | A | C | C |
|     | T | T | G | G | A | A | C | T |
|     | A | T | G | C | C | A | T | T |
|     | A | T | G | G | C | A | C | T |

t

**Profile matrix**

|   | A | T | C | C | A | G | C | T |
|---|---|---|---|---|---|---|---|---|
| A | 5 | 1 | 0 | 0 | 5 | 5 | 0 | 0 |
| T | 1 | 5 | 0 | 0 | 0 | 1 | 1 | 6 |
| G | 1 | 1 | 6 | 3 | 0 | 1 | 0 | 0 |
| C | 0 | 0 | 1 | 4 | 2 | 0 | 6 | 1 |

4

**Consensus string**

|  | A | T | G | C | A | A | C | T |
|--|---|---|---|---|---|---|---|---|

(b) The alignment matrix, profile matrix and consensus string formed from the 8-mers starting at positions **s** = (8, 19, 3, 5, 31, 27, 15) in figure 4.2 (d).

# Notation:

$P(\vec{s}) :=$ profile matrix wrt starting position vector $\vec{s}$

$M_{P(\vec{s})}(j) :=$ largest count in column $j$ of $P(s)$

|  |  | A | T | C | C | A | G | C | T |
|---|---|---|---|---|---|---|---|---|---|
| | | G | G | G | C | A | A | C | T |
| | | A | T | G | G | A | T | C | T |
| **Alignment** | | A | A | G | C | A | A | C | C |
| | | T | T | G | G | A | A | C | T |
| | | A | T | G | C | C | A | T | T |
| | | A | T | G | G | C | A | C | T |
| **Profile** | **A** | 5 | 1 | 0 | 0 | 5 | 5 | 0 | 0 |
| | **T** | 1 | 5 | 0 | 0 | 0 | 1 | 1 | 6 |
| | **G** | 1 | 1 | 6 | 3 | 0 | 1 | 0 | 0 |
| | **C** | 0 | 0 | 1 | 4 | 2 | 0 | 6 | 1 |
| **Consensus** | | A | T | G | C | A | A | C | T |

$P(s) \longrightarrow$

$$M_{P(s)}(1) = 5$$
$$M_{P(s)}(2) = 5$$
$$M_{P(s)}(8) = 6$$

Consensus score:
$$\text{Score } (\bar{s}, \text{DNA}) = \sum_{j=1}^{\ell} M_{P(s)}(j)$$

| | | A | T | C | C | A | G | C | T |
|---|---|---|---|---|---|---|---|---|---|
| | | G | G | G | C | A | A | C | T |
| | | A | T | G | G | A | T | C | T |
| **Alignment** | | A | A | G | C | A | A | C | C |
| | | T | T | G | G | A | A | C | T |
| | | A | T | G | C | C | A | T | T |
| | | A | T | G | G | C | A | C | T |
| | **A** | 5 | 1 | 0 | 0 | 5 | 5 | 0 | 0 |
| **Profile** | **T** | 1 | 5 | 0 | 0 | 0 | 1 | 1 | 6 |
| | **G** | 1 | 1 | 6 | 3 | 0 | 1 | 0 | 0 |
| | **C** | 0 | 0 | 1 | 4 | 2 | 0 | 6 | 1 |
| **Consensus** | | A | T | G | C | A | A | C | T |

Here,

$$\text{Score } (\bar{s}, \text{DNA}) = 5 + 5 + 6 + 4 + 5 + 5 + 6 + 6$$

Why? Strength of a profile:

$\ell \cdot t$ means best possible alignment — same letter in each spot

$\frac{\ell t}{4}$: worst alignment — equal mix of nucleotides per spot

---

**Motif Finding Problem:**

*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 of $DNA$, and $l$, the length of the pattern to find.

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

---

Note: In reality, often use
entropy: O.

Let $P_{i,j}$ be $(i,j)^{th}$ entry in profile.

$$Entropy = \sum_{j=1}^{\ell} \sum_{i=1}^{4} \left[ \frac{P_{i,j}}{t} \log \frac{P_{i,j}}{t} \right]$$

where $t = \#$ sequences

This is more statistically robust measure

(but algorithm is essentially unchange)

# Reframing:

Sift through large # of alternatives to find best one
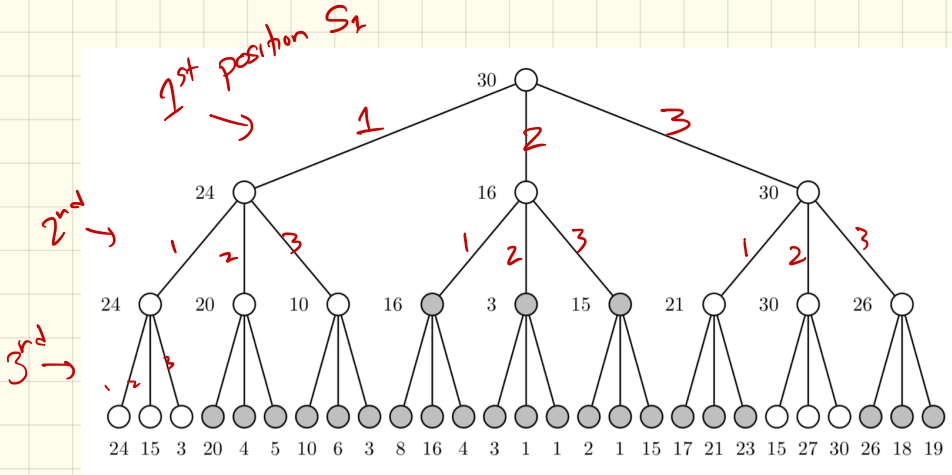
$(n-l+1)^t$ starting positions!

**POSSIBLE $\vec{S}$'s:**

| | | | | |
|---|---|---|---|---|
| ( | 1, | 1, ..., | 1, | 1 ) |
| ( | 1, | 1, ..., | 1, | 2 ) |
| ( | 1, | 1, ..., | 1, | 3 ) |
| | | $\vdots$ | | |
| ( | 1, | 1, ..., | 1, | $n-l+1$ ) |
| ( | 1, | 1, ..., | 2, | 1 ) |
| ( | 1, | 1, ..., | 2, | 2 ) |
| ( | 1, | 1, ..., | 2, | 3 ) |
| | | $\vdots$ | | |
| ( | 1, | 1, ..., | 2, | $n-l+1$ ) |
| | | $\vdots$ | | |
| ( | $n-l+1,$ | $n-l+1,$ ..., | $n-l+1,$ | 1 ) |
| ( | $n-l+1,$ | $n-l+1,$ ..., | $n-l+1,$ | 2 ) |
| ( | $n-l+1,$ | $n-l+1,$ ..., | $n-l+1,$ | 3 ) |
| | | $\vdots$ | | |
| ( | $n-l+1,$ | $n-l+1,$ ..., | $n-l+1,$ | $n-l+1$ ) |

(counting in base $n-l+1$?)

Branch & bound intuition:
What if we can go partway
but rule out entire subtree?



1st position $S_2$

30

1          2          3

2nd →

24         16         30

1    2    3      1    2    3      1    2    3

24   20   10    16   3   15    21   30   26

3rd →

1  2  3

24 15  3  20  4  5  10  6  3  8  16  4  3  1  1  2  1  15 17 21 23 15 27 30 26 18 19

(More details, next time,
plus connection to medians)

(through mid 4.6)
rest of Ch 4 on Thursday