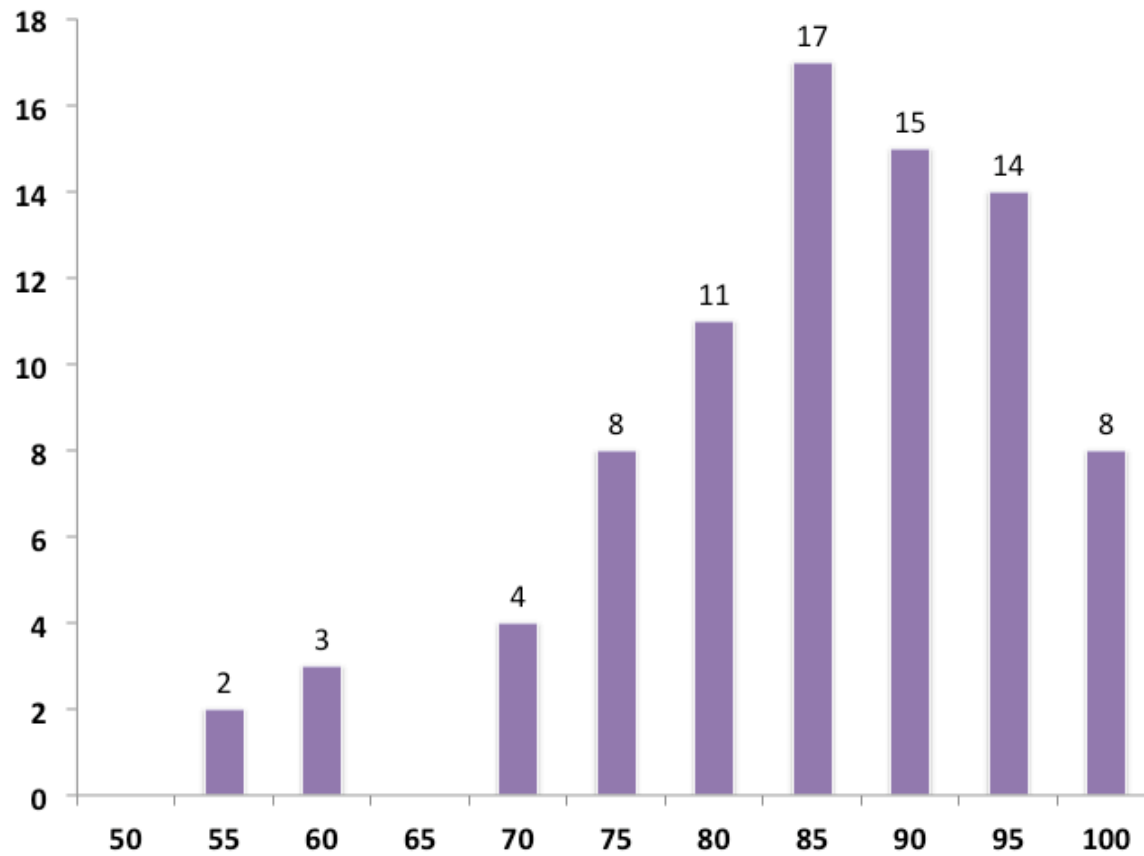# CSE 421 Midterm Scores



Mean 83
Sigma 11

# CSE 421 Algorithms

Sequence Alignment

# Sequence Alignment

Goal: position characters in strings so they "best" line up with one another

We can do this via Dynamic Programming

# What is an alignment?

Compare two strings and see how similar they are

Maximize the # of chars in a string that line up

ATGTTAT vs ATCGTAC

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

3

# What is an alignment?

Compare two strings and see how similar they are

Maximize the # of chars in a string that line up

## ATGTTAT vs ATCGTAC

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

matches                    mismatches

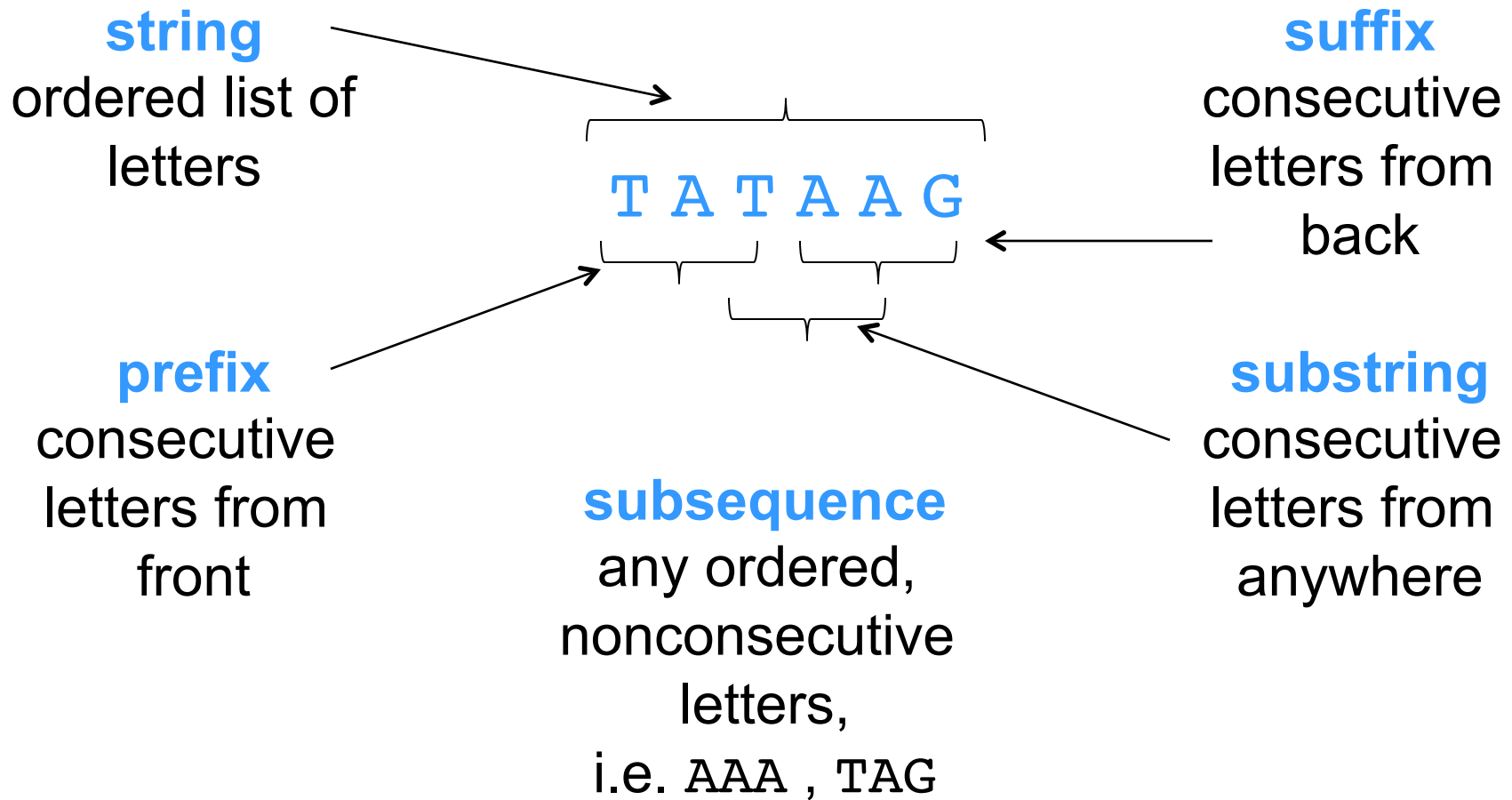# Why do we align?

Biology

Most widely used comp. tools in biology

New sequences always compared to databases

**Similar sequences often have similar origin and/or function**

Other

spell check, diff, svn/git/…, plagiarism, …

# Terminology

**string**
ordered list of letters

**suffix**
consecutive letters from back

$$T\ A\ T\ A\ A\ G$$

**prefix**
consecutive letters from front

**subsequence**
any ordered, nonconsecutive letters,
i.e. `AAA` , `TAG`

**substring**
consecutive letters from anywhere

# Formal definition of an alignment

```
a c g c t g        a c - - g c t g
                       |       |   |
c a t g t          - c a t g - t -
```

An alignment of strings S, T is represented as a
   pair of strings S', T' with gaps "–"  s.t.

1.  |S'| = |T'|, and                (|S| = "length of S")

2.  Removing gaps leaves S, T

(Note that this is a definition for a general alignment, not optimal.)

# Scoring an arbitrary alignment

Want to determine whether an alignment is "good" or "bad" so we define a cost function

score of (mis)aligning chars x & y $= \sigma(x, y) =$ $\begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$

Total value/score of an alignment

$$\Sigma \; \sigma(S'[i], T'[i])$$

Optimal alignment

Max alignment score of all poss. alignments

# Scoring an arbitrary alignment

```
a   c   -   -   g   c   t   g
    |           |       |
-   c   a   t   g   -   t   -
```

-1  +2  -1  -1  +2  -1  +2  -1

Score = +1

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

# Can we use
# Dynamic Programming?

1. Identify **subproblems**

   We can reuse the solution to smaller substrings (prefixes in this case)

2. Argue that we have **optimal substructure**

   Appending two optimal alignments should also be optimally aligned (some may change at the interface)

# Arguing for Optimal Substructure

Assume strings S & T are optimally aligned except for the last character

3 options for the last character:

1. match -- `S[i]` & `T[j]` aligned

2. mismatch -- `S[i]` & "–" aligned

3. mismatch -- `T[j]` & "–" aligned

\* Never align "–" & "–"; i.e. σ("–", "–") << 0

# "Recipe" for using DP for problems like this

1.  Argue for optimal substructure (☑)
2.  Find a recursive relation for subproblem costs
    > Use (1), find all subproblems that might contribute to an optimal cost
3.  Implement a bottom-up use of (2) to fill in a table of subproblem costs
4.  Write a recursive algorithm using the table from (3) to construct actual solutions to subproblems ("traceback")

# Setting up Optimal Alignment in $O(n^2)$ via DP

Input:      strings S, T

|S| = n, |T| = m

Output:    optimal alignment score

→ Generate the score first and then trace backwards to recover the actual alignment

# Setting up Optimal Alignment in O(n²) via DP

Compute optimal alignment of **all combinations of prefixes**, & store in a table for the future

T →
S ∨

| | - | A | C | G | T | ... | T |
|---|---|---|---|---|---|---|---|
| - | 0 | -1 | -2 | -3 | -4 | | -n |
| A | -1 | 2 | 1 | 0 | -1 | | |
| C | -2 | 1 | 4 | 3 | | | |
| G | -3 | 0 | 3 | ★ | | | |
| T | -4 | -1 | | | | | |
| ... | | | | | | ... | |
| T | -n | | | | | | |

Start UL, nothing aligned
End LR, w/ optimal score

Move diagonally → align chars
Move vert/horiz → introduce gap

$V(i,j)$ = optimal alignment score of
$S[1]…S[i]$ and $T[1]…T[j]$

i.e. all possible prefixes of S and T

14

# Computing the table: Base Case

| T → S v | - | A | C | G | T | ... | T |
|---|---|---|---|---|---|---|---|
| - | 0 | -1 | -2 | -3 | -4 | | -n |
| A | -1 | 2 | 1 | 0 | -1 | | |
| C | -2 | 1 | 4 | 3 | | | |
| G | -3 | 0 | 3 | | | | |
| T | -4 | -1 | | | | | |
| ... | | | | | | ... | |
| T | -n | | | | | | |

**Column:**

S aligns with nothing in T
all mismatches

$V(i,0) = \Sigma\sigma(S[k], \text{``-''})$

$= i*\sigma(S[k], \text{``-''})$

**Row:**

T aligns with nothing in S
all mismatches

$V(0,j) = \Sigma\sigma(\text{``-''}, T[k])$

$= j*\sigma(\text{``-''}, T[k])$

# Computing the table: General Case

T →
S ∨

| | - | A | C | G | T | ... | T |
|---|---|---|---|---|---|---|---|
| - | 0 | -1 | -2 | -3 | -4 | | -n |
| A | -1 | 2 | 1 | 0 | -1 | | |
| C | -2 | 1 | 4 | 3 | | | |
| G | -3 | 0 | 3 | ★ | | | |
| T | -4 | -1 | | | | | |
| ... | | | | | | ... | |
| T | -n | | | | | | |

At any given point in computing the table, we can choose whether it's best to

Align 2 characters

Take a gap

# Computing the table: General Case

$$\bigstar = V(i, j) = \max \begin{cases} V(i-1,\ j-1) + \sigma(S[i],\ T[j]) & \text{match} \\ V(i-1,\ j) \quad + \sigma(S[i],\ \text{``-''}) & \text{mismatch} \\ V(i,\ j-1) \quad + \sigma(\text{``-''},\ T[j]) & \text{mismatch} \end{cases}$$

Cost of ops so far

Cost of next op (match/mismatch)

|   | - | A | C | G | T | ... | T |
|---|---|---|---|---|---|-----|---|
| - | 0 | -1 | -2 | -3 | -4 | | -n |
| A | -1 | 2 | 1 | 0 | -1 | | |
| C | -2 | 1 | 4 | 3 | | | |
| G | -3 | 0 | 3 | ★ | | | |
| T | -4 | -1 | | | | | |
| ... | | | | | ... | | |
| T | -n | | | | | | |

Need these 3 positions filled in to determine ★

# Example: base case

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

T →
S v

|  |  | C | A | T | G | T |
|---|---|---|---|---|---|---|
|  |  | i=0 | 1 | 2 | 3 | 4 | 5 |
|  | j=0 | 0 | -1 | -2 | -3 | -4 | -5 |
| A | 1 | -1 |  |  |  |  |  |
| C | 2 | -2 |  |  |  |  |  |
| G | 3 | -3 |  |  |  |  |  |
| C | 4 | -4 |  |  |  |  |  |

```
V(i,0) = i*σ(S[k], "-")
V(0,j) = j*σ("-",,T[k])
```

`8

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

# Example: general step

|  |  | C | A | T | G | T |
|---|---|---|---|---|---|---|
|  | i=0 | 1 | 2 | 3 | 4 | 5 |
| j=0 | 0 | -1 | -2 | -3 | -4 | -5 |
| A  1 | -1 | -1 |  |  |  |  |
| C  2 | -2 |  |  |  |  |  |
| G  3 | -3 |  |  |  |  |  |
| C  4 | -4 |  |  |  |  |  |

T →
S ∨

19

# Example: general step

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

|  |  | **C** | **A** | **T** | **G** | **T** |
|---|---|---|---|---|---|---|
| T → <br> S v |  | i=0 1 | 2 | 3 | 4 | 5 |
|  | j=0 | 0 -1 | -2 | -3 | -4 | -5 |
| **A** | 1 | -1 -1 | | | | |
| **C** | 2 | -2 | | | | |
| **G** | 3 | -3 | | | | |
| **C** | 4 | -4 | | | | |

$$V(i, j) = \max \begin{cases} V(i-1,\ j-1)\ +\ \sigma(S[i],\ T[j]) & \nwarrow \\ V(i-1,\ j)\quad +\ \sigma(S[i],\ ``\text{-}") & \uparrow \\ V(i,\ j-1)\quad\ +\ \sigma(``\text{-}",\ T[j]) & \leftarrow \end{cases}$$

# Example: general step

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

T →
S v

|  |  | C | A | T | G | T |
|---|---|---|---|---|---|---|
|  | i=0 | 1 | 2 | 3 | 4 | 5 |
| j=0 | 0 | -1 | -2 | -3 | -4 | -5 |
| A 1 | -1 | -1 |  |  |  |  |
| C 2 | -2 |  |  |  |  |  |
| G 3 | -3 |  |  |  |  |  |
| C 4 | -4 |  |  |  |  |  |

$$V(i, j) = \max \begin{cases} V(0,1) + \sigma(S[1], T[2]) \\ V(0,2) + \sigma(S[1], \text{``-''}) \\ V(1,1) + \sigma(\text{``-''}, T[2]) \end{cases}$$

21

# Example: general step

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

$$T \rightarrow$$
$$S \vee$$

|  | | C | A | T | G | T |
|---|---|---|---|---|---|---|
| | i=0 | 1 | 2 | 3 | 4 | 5 |
| j=0 | 0 | -1 | -2 | -3 | -4 | -5 |
| A 1 | -1 | -1 | **1** | | | |
| C 2 | -2 | | | | | |
| G 3 | -3 | | | | | |
| C 4 | -4 | | | | | |

-1 + 2 = **1, match**

$$V(i, j) = \max \begin{cases} \\ -2 -1 = -3 \\ \\ -1 -1 = -2 \end{cases}$$

22

$$\sigma(x, y) = \begin{cases} \text{match} & 2 \\ \text{mismatch} & -1 \end{cases}$$

# Example: completed table

T →
S v

|  |  | **C** | **A** | **T** | **G** | **T** |
|---|---|---|---|---|---|---|
|  |  | i=0 | 1 | 2 | 3 | 4 | 5 |
| | j=0 | 0 | -1 | -2 | -3 | -4 | -5 |
| **A** | 1 | -1 | -1 | 1 | 0 | -1 | -2 |
| **C** | 2 | -2 | 1 | 0 | 0 | -1 | -2 |
| **G** | 3 | -3 | 0 | 0 | -1 | 2 | 1 |
| **C** | 4 | -4 | -1 | -1 | -1 | 1 | 1 |

Time = O(mn) = O(|S|*|T|)

# How do we find the alignment itself? Traceback

Trace LR to UL following
highest score path

Can go  ↖ ↑ ←

Multiple optimal
alignments are possible

We can break ties
arbitrarily

|     |     | C | A | T | G | T |
|-----|-----|-----|-----|-----|-----|-----|
|     | i=0 | 1 | 2 | 3 | 4 | 5 |
| j=0 | 0 | -1 | -2 | -3 | -4 | -5 |
| A 1 | -1 | -1 | 1 | 0 | -1 | -2 |
| C 2 | -2 | 1 | 0 | 0 | -1 | -2 |
| G 3 | -3 | 0 | 0 | -1 | 2 | 1 |
| C 4 | -4 | -1 | -1 | -1 | 1 | 1 |

Corresponding
Alignment:
   CATGT
   –ACGC

24

# Example

Mismatch = -1
Match = 2

| j | | 0 | 1 | 2 | 3 | 4 | 5 | |
|---|---|---|---|---|---|---|---|---|
| i | | | c | a | t | g | t | ←T |
| 0 | | 0 | -1 | -2 | -3 | -4 | -5 | |
| 1 | a | -1 | -1 | 1 | 0 | -1 | -2 | |
| 2 | c | -2 | 1 | 0 | 0 | -1 | -2 | |
| 3 | g | -3 | 0 | 0 | -1 | 2 | 1 | |
| 4 | c | -4 | -1 | -1 | -1 | 1 | 1 | |
| 5 | t | -5 | -2 | -2 | 1 | 0 | 3 | |
| 6 | g | -6 | -3 | -3 | 0 | 3 | 2 | |

↑
S

# Complexity Notes

Time = O(mn), (value and alignment)

Space = O(mn)

Easy to get value in Time = O(mn) and Space = O(min(m,n))

Possible to get value *and alignment* in Time = O(mn) and Space =O(min(m,n)) (KT section 6.7)

# Significance of Alignments

Is "42" a good score? *Compared to what?*

Easier to compare when using standardized scoring
   functions, esp. for DNA

Usual approach: compared to a specific "null
model", such as "random sequences"

Interesting stats problem; much is known

# Variations

## Local Alignment

Preceding gives *global* alignment, i.e. full length of both strings;

Might well miss strong similarity of part of strings amidst dissimilar flanks

## Gap Penalties

10 adjacent spaces cost 10 x one space?

## Many others

## Similarly fast DP algs often possible

# Summary: Alignment

Functionally similar proteins/DNA often have recognizably similar sequences even after eons of divergent evolution

Ability to find/compare/experiment with "same" sequence in other organisms is a huge win

Surprisingly simple scoring works well in practice: score positions separately & add, usually w/ fancier gap model like affine

Simple dynamic programming algorithms can find *optimal* alignments under these assumptions in poly time (product of sequence lengths)

This, and heuristic approximations to it like BLAST, are workhorse tools in molecular biology, and elsewhere.

# Summary: Dynamic Programming

## Keys to D.P. are to

a) identify the subproblems (usually repeated/overlapping)

b) solve them in a careful order so all small ones solved before they are needed by the bigger ones, and

c) build table with solutions to the smaller ones so bigger ones just need to do table lookups (*no* recursion, despite recursive formulation implicit in (a))

d) Implicitly, optimal solution to whole problem devolves to optimal solutions to subproblems

## A *really* important algorithm design paradigm