# Algorithms in Bioinformatics

## More dynamic programming

# Recap

- HW - up after class
  due on Tuesday the $2^{nd}$
- No class next Thursday

# Dynamic Programming:
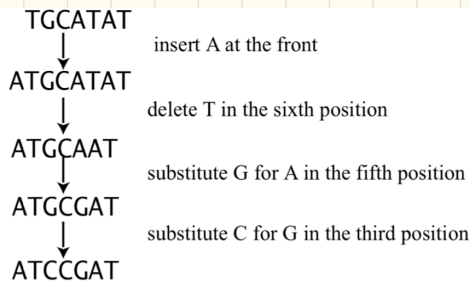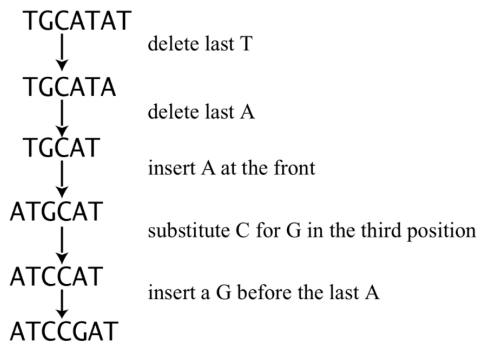## aka "Smart Recursion"

High level idea:
- Formulate recursion
- Notice recursive alg. is slow
- Store recursive call answers for lookup

{ Wait, could just do this iteratively

# Last time: Edit distance

## Example:

TGCATAT
↓ delete last T
TGCATA
↓ delete last A
TGCAT
↓ insert A at the front
ATGCAT
↓ substitute C for G in the third position
ATCCAT
↓ insert a G before the last A
ATCCGAT

TGCATAT
↓ insert A at the front
ATGCATAT
↓ delete T in the sixth position
ATGCAAT
↓ substitute G for A in the fifth position
ATGCGAT
↓ substitute C for G in the third position
ATCCGAT

## Alignment matrix:

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

(at most m+n columns)

## Another way:
### Write # of repititions:

| | | 0 | 1 | 2 | 2 | 3 | 4 | 5 | 6 | 7 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **v** | = | A | T | - | G | T | T | A | T | - | |
| | | | | | | | | | | | |
| **w** | = | A | T | C | G | T | - | A | - | C | |
| | | 0 | 1 | 2 | 3 | 4 | 5 | 5 | 6 | 6 | 7 |

This leads to formulation from
last time.

Book's view:



|   |   | 0 | 1 | 2 | 2 | 3 | 4 | 5 | 6 | 7 | 7 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **v** | = | A | T | - | G | T | T | A | T | - |   |

$$
\begin{array}{cccccccc}
| & | & & | & | & & | & \\
\end{array}
$$

|   |   | A | T | C | G | T | - | A | - | C |
|---|---|---|---|---|---|---|---|---|---|---|
| **w** | = |   |   |   |   |   |   |   |   |   |
|   |   | 0 | 1 | 2 | 3 | 4 | 5 | 5 | 6 | 6 | 7 |

# Now : Longest Common Subsequence

**Longest Common Subsequence Problem**:
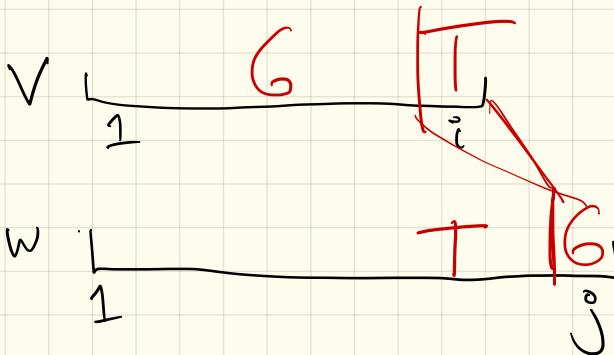*Find the longest subsequence common to two strings.*

AT G A A T
GT A A A C

**Input:** Two strings, **v** and **w**.

**Output:** The longest common subsequence of **v** and **w**.

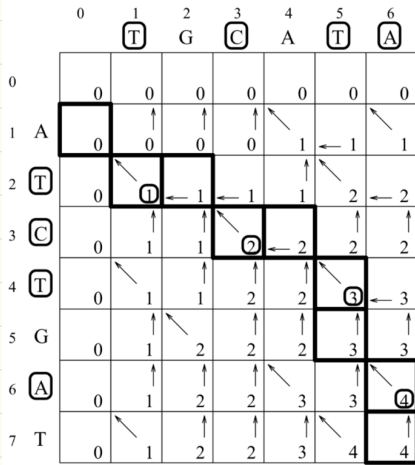Let $S_{i,j}$ = length of LCS between $v[1..i]$ and $w[1..j]$.
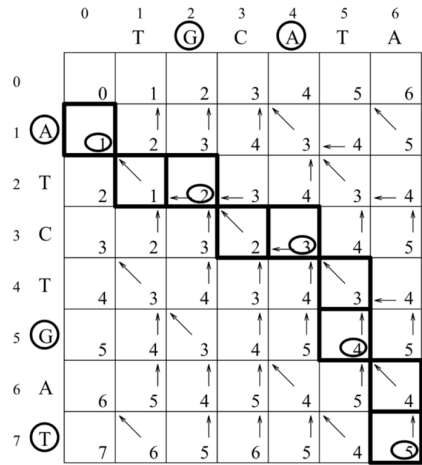
If I look at last letter(s), 2 options:

- last letter is common
- last letter isn't

# Aside: LCS + Edit distance are different!



Computing similarity s(V,W)=4
V and W have a subsequence TCTA in common

Computing distance d(V,W)=5
V can be transformed into W by deleting A,G,T and inserting G,A
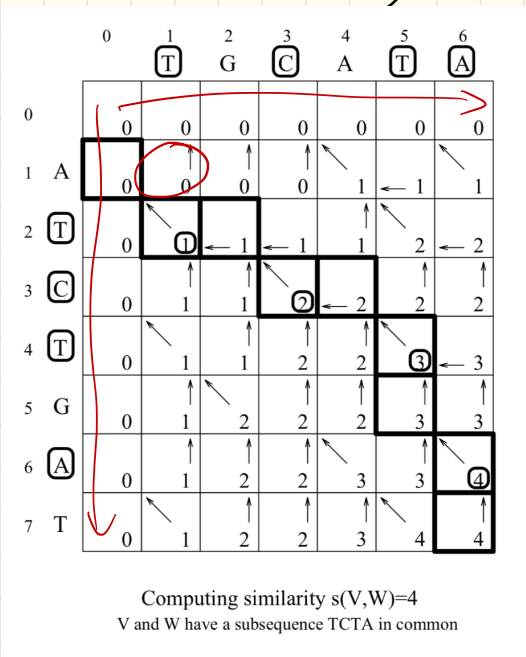
Alignment:
```
A T - C - T G A T
- T G C A T - A -
```

(But connected still...)
see HW..

# So : Try them all!

$$V[1..i], W[1..j]$$

$$s_{i,j} = \max \begin{cases} s_{i-1,j} \\ s_{i,j-1} \\ s_{i-1,j-1} + 1, & \text{if } v_i = w_j \end{cases}$$

↑ these are in LCS

Then think about not recomputing- so store in big table again



|   |   | 0 | 1 (T) | 2 G | 3 (C) | 4 A | 5 (T) | 6 (A) |
|---|---|---|---|---|---|---|---|---|
| 0 |   | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | A | 0 | 0 | 0 | 0 | 1 | ← 1 | 1 |
| 2 | (T) | 0 | ① | ← 1 | ← 1 | 1 | 2 | ← 2 |
| 3 | (C) | 0 | 1 | 1 | ② | ← 2 | 2 | 2 |
| 4 | (T) | 0 | 1 | 1 | 2 | 2 | ③ | ← 3 |
| 5 | G | 0 | 1 | 2 | 2 | 2 | 3 | 3 |
| 6 | (A) | 0 | 1 | 2 | 2 | 3 | 3 | ④ |
| 7 | T | 0 | 1 | 2 | 2 | 3 | 4 | 4 |

Computing similarity s(V,W)=4

V and W have a subsequence TCTA in common

Alignment:

```
A T - C - T G A T
- T G C A T - A -
```

# Pseudo code :

```
LCS(v, w)
1   for i ← 0 to n
2       s_{i,0} ← 0
3   for j ← 1 to m
4       s_{0,j} ← 0
5   for i ← 1 to n
6       for j ← 1 to m
7           s_{i,j} ← max { s_{i-1,j}
                            s_{i,j-1}
                            s_{i-1,j-1} + 1,   if v_i = w_j
8           b_{i,j} ← { " ↑ "    if s_{i,j} = s_{i-1,j}
                        " ← "    if s_{i,j} = s_{i,j-1}
                        " ↖ ",   if s_{i,j} = s_{i-1,j-1} + 1
9   return (s_{n,m}, b)
```

$$s_{i,j} \leftarrow \max \begin{cases} s_{i-1,j} \\ s_{i,j-1} \\ s_{i-1,j-1} + 1, & \text{if } v_i = w_j \end{cases}$$

$$b_{i,j} \leftarrow \begin{cases} \text{``}\uparrow\text{''} & \text{if } s_{i,j} = s_{i-1,j} \\ \text{``}\leftarrow\text{''} & \text{if } s_{i,j} = s_{i,j-1} \\ \text{``}\nwarrow\text{''}, & \text{if } s_{i,j} = s_{i-1,j-1} + 1 \end{cases}$$

```
PRINTLCS(b, v, i, j)
1    if i = 0 or j = 0
2        return
3    if b_{i,j} = " ↖ "
4        PRINTLCS(b, v, i - 1, j - 1)
5        print v_i
6    else
7        if b_{i,j} = " ↑ "
8            PRINTLCS(b, v, i - 1, j)
9        else
10           PRINTLCS(b, v, i, j - 1)
```

# Now: Back to some biology!

LCS is a way to score similarity:

- +1 for a match
- +0 for a mismatch (indel)

Edit distance is too!

- insertion, deletion + substitution all cost +1.

Biology changes are more complex...

Generalize:

Make a scoring matrix. $\delta$:

Ex: $\delta$:

|   | G | C | T | A | - |
|---|---|---|---|---|---|
| G |   | 1 |   |   |   |
| C |   |   | -3 |   |   |
| T |   |   | 1 |   |   |
| A |   |   |   |   |   |
| - |   |   | +½ |   |   |

C-T mutation charge

usually negative

T insertion / deletion charge

# New goal :

**Global Alignment Problem**:
*Find the best alignment between two strings under a given scoring matrix.*

**Input:** Strings **v**, **w** and a scoring matrix $\delta$.

**Output:** An alignment of **v** and **w** whose score (as defined by the matrix $\delta$) is maximal among all possible alignments of **v** and **w**.

Same type of recursion:
When looking at $v[1..i]$
and $w[1..j]$
- could match ~~them~~ $v[i] + w[j]$

Score
for $v[1..i-1]$
$w[1..j-1]$ $\longrightarrow T[i,j-1] + \delta(v[i], w[j])$

- could match $v[i]$ to $-$
$$T[i-1,j] + \delta(v[i], -)$$

- could match $w[j]$ to $-$
$$T[i,j-1] + \delta(-, w[j])$$

# End recurrence:

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \delta(v_i, -) \\ s_{i,j-1} + \delta(-, w_j) \\ s_{i-1,j-1} + \delta(v_i, w_j) \end{cases}$$

remove $V[i]$ + recurse

remove $w[j]$ + rec.

match last two

# Example:
- charge mismatches by $-\mu$
- indels by $-\sigma$
- and add $+1$ for matches

Get score =

$$\# matches - \mu(\# mismatches)$$
$$- \sigma(\# indels)$$

Then:

$$s_{i,j} = \max \begin{cases} s_{i-1,j} - \sigma \\ s_{i,j-1} - \sigma \\ s_{i-1,j-1} - \mu, \text{ if } v_i \neq w_j \\ s_{i-1,j-1} + 1, \text{ if } v_i = w_j \end{cases}$$

(Note: LCS is example of this!)
$$\sigma = 0$$
Runtime: Same alg! $O(mn)$

# How to get S?

- For DNA, usually just ask for $\mu$ & $\sigma$ / as part of input.

- For amino acids — harder!
    - Point Accepted Mutations (PAM)
    - Block Substitution (BLOSUM)

(Two most common ones)

Reason: probability that Ser mutates to Phe is ≈ 3 times higher than Trp to Phe.

But — gets complex!

(Read 6.7 for details — it's constructed iteratively.)

# Local alignment

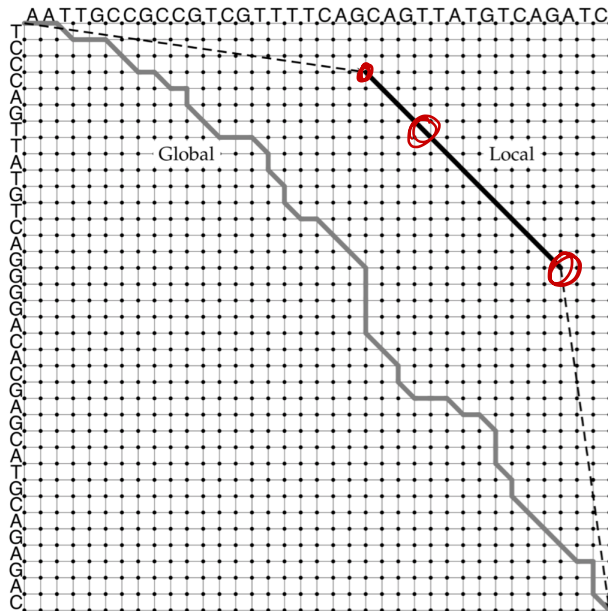Instead of entire string, might want to find substrings with good alignment.

(Book's example: homeobox genes)

High level issue:



```
--T--CC-C-AGT--TATGT-CAGGGGACACG--A-GCATGCAGA-GAC
  |  ||| || | | ||| |||   || | || | ||||  |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG--T-CAGAT--C


        tccCAGTTATGTCAGgggacacgagcatgcagagac
           ||||||||||||||
aattgccgccgtcgttttcagCAGTTATGTCAGatc
```

Essentially, since we charge:
"  $-\mu$ for mismatches
•  $-\sigma$ for indels,

Global alignment will likely
   miss this good substring
   entirely.

<u>So</u>:                          $S_{i,j} =$

---

**Local Alignment Problem**:
*Find the best local alignment between two strings.*

   **Input:** Strings v and w and a scoring matrix $\delta$.

   **Output:** Substrings of v and w whose global alignment, as
   defined by $\delta$, is maximal among all global alignments of all
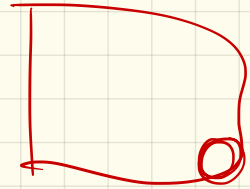   substrings of v and w.

---

However, this is still pretty close −
 ○ Instead of best $(0,0)$ to $(m,n)$
        path
 ○ Want best $(i,j)$ to $(i',j')$
      path, for <u>any</u> $i,j, i',j'$.

Really — just need to let
you start match at any
$(i,j)$ if matching$_v$ $v[1..i]$
$+$ $w[1..j]$ is bad

So :

$$s_{i,j} = \max \begin{cases} 0 \\ s_{i-1,j} + \delta(v_i, -) \\ s_{i,j-1} + \delta(-, w_j) \\ s_{i-1,j-1} + \delta(v_i, w_j) \end{cases}$$

Then:

- Build a matrix
- Fill it in
- Look at entire
→ m×n matrix for
   highest score

→ for $i \leftarrow 1$ to $n$
→ for $j \leftarrow 1$ to $m$
   4 lookups, 3 additions
   take max

Runtime:
$O(mn) + O(mn) \cdot \{ O(1) = \boxed{O(mn)}$
↑
also
$O(mn)$ space

## Next time :

- Alignment w/ Gap penalties
- Multiple alignment

Read 6.3
for next time