# CSE 397-497:
# *Computational Issues in Molecular Biology*

# Lecture 3

# Spring 2004

For your class lecture, you must give me a ranked list of your top 3 topics in order of preference by 5:00 pm today:

*order we will cover topics in course*

- sequence comparison & alignment (pairwise & multiple),
- sequencing and sequence assembly,
- physical mapping of DNA,
- phylogenetic trees,
- genome rearrangements,
- RNA and protein structure prediction,
- DNA microarrays,
- DNA computing.

Consider:

**AGTAGCATC**       *versus*       **AGTAGCATC**

↕                ↕

**AGTGCACC**                 **GACACGATT**

That the two DNA fragments on the left somehow seem more similar than the two on the right could be significant.

Question: how can we measure sequence similarity?

LEHIGH
U N I V E R S I T Y.

Why is this important?

- Given a new DNA sequence, one of the first things a biologist will want to do is search databases of known sequences to see if anyone has recorded something similar. (As we've seen, genetic sequences are long and the databases are enormous, so efficiency will be an issue.)

- Sequence similarity can provide clues about function.

- Similarity can provide clues about evolutionary relationships.

- Many other problems from computational biology incorporate some notion of sequence similarity as a basic premise.

LEHIGH
UNIVERSITY

A *sequence* is a linear string of symbols over a finite alphabet.

(Note: *string* and *sequence* are often used synonymously.)

Some basic sequence concepts:

*length*              number of symbols in *s* (written |*s*|).

*empty string*     sequence of length 0 (written $\varepsilon$).

*subsequence*    sequence that can be obtained from *s* by removing some symbols (**ACG** is a subsequence of **TATCTG**).

*supersequence*  if *t* is a subsequence of *s*, then *s* is a supersequence of *t*.

More basic sequence concepts:

*substring*        sequence of consecutive symbols appearing in *s* (**ACG** is <u>not</u> a substring of **TATCTG**, but **TCT** is).

(Observation: every substring is a subsequence of the string in question, but not vice versa.)

*superstring*      if *t* is a substring of *s*, then *s* is a superstring of *t*.

*interval*      set of consecutive indices of a string, e.g., [2..4] refers to 2*nd* through 4*th* symbols of string *s*, or *s*[2]*s*[3]*s*[4].

And lastly:

*prefix*     substring of s of the form $s[1..j]$ where $0 \leq j \leq |s|$ (when $j = 0$, prefix is the empty string).

*suffix*     substring of s of the form $s[i..|s|]$ where $1 \leq i \leq |s| + 1$ (when $i = |s| + 1$, suffix is the empty string).

**AT** is a prefix (and substring and subsequence) of **ATCCAG**.

**AG** is a suffix (and substring and subsequence) of **ATCCAG**.

The concept of a sequence is extremely broad.  In this course, we are concerned with genetic sequences.  However, there are other important kinds of sequence data:

- ASCII text,

- speech,

- handwriting (pen-strokes).

Likewise, the same algorithmic techniques turn up again and again, often under different names:

- approximate string matching,

- edit (or evolutionary) distance,

- dynamic time warping.

LEHIGH UNIVERSITY

An obvious idea that comes to mind is to line up each symbol and count the number that don't match:

A C G T G C

$= 2$

A A G A G C

This is known as Hamming distance and forms the basis for most error correcting codes.

But it doesn't work for the kinds of sequences we care about:

A C G T G C

$?$ $= 6$

C G T G C

Just one missing symbol at the start of the second sequence leads to a large distance.

# Genomes aren't static ...



... sequence comparison must account for this.

Different kinds of mutations can arise during DNA replication.



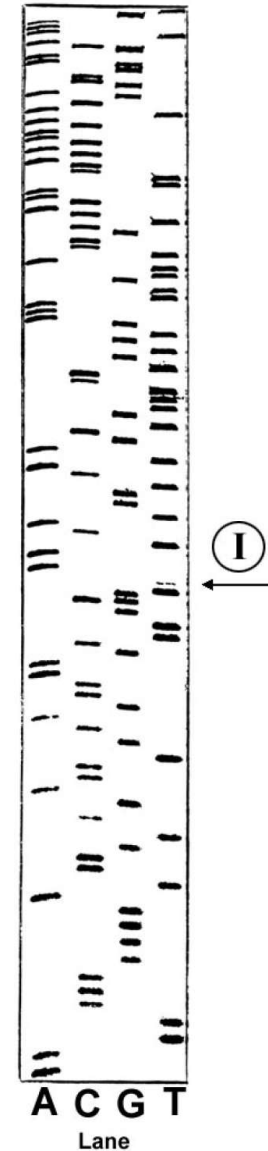http://www.accessexcellence.org/AB/GG/mutation.html

In addition, errors can arise during the sequencing process:

> "...the error rate is generally less than 1% over the first 650 bases and then rises significantly over the remaining sequence."

http://genome.med.harvard.edu/dnaseq.html

A hard-to-read gel (arrow marks location where bands of similar intensity appear in two different lanes):

http://hshgp.genome.washington.edu/teacher_resources/99-studentDNASequencingModule.pdf
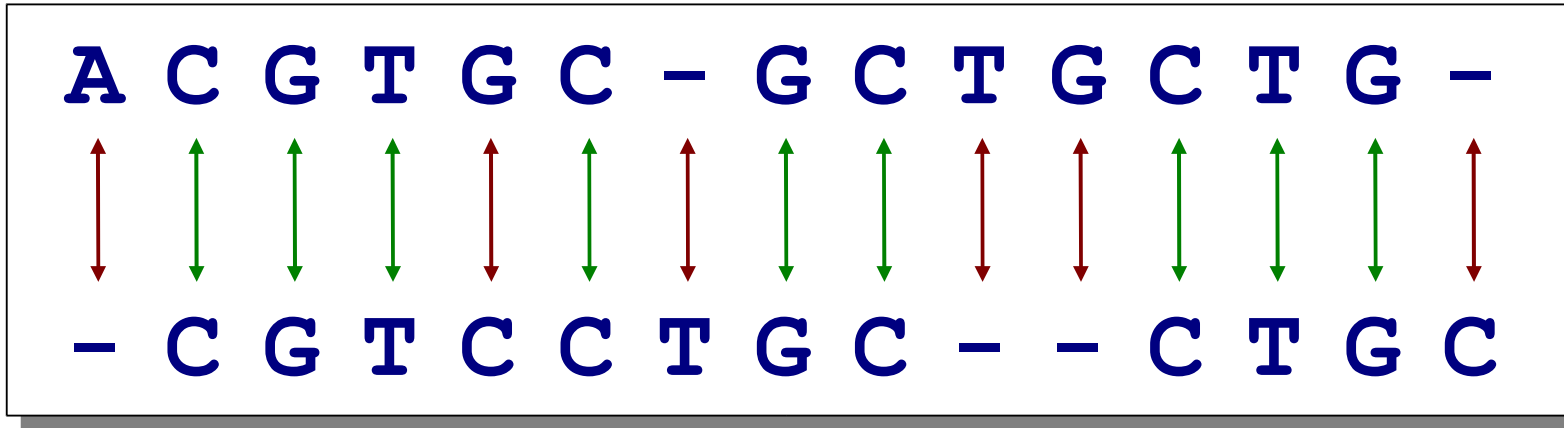


A C G T
Lane

LEHIGH
UNIVERSITY

As we have seen, the two sequences we wish to compare may have different lengths.  As a result, we need to allow for deletions and insertions.

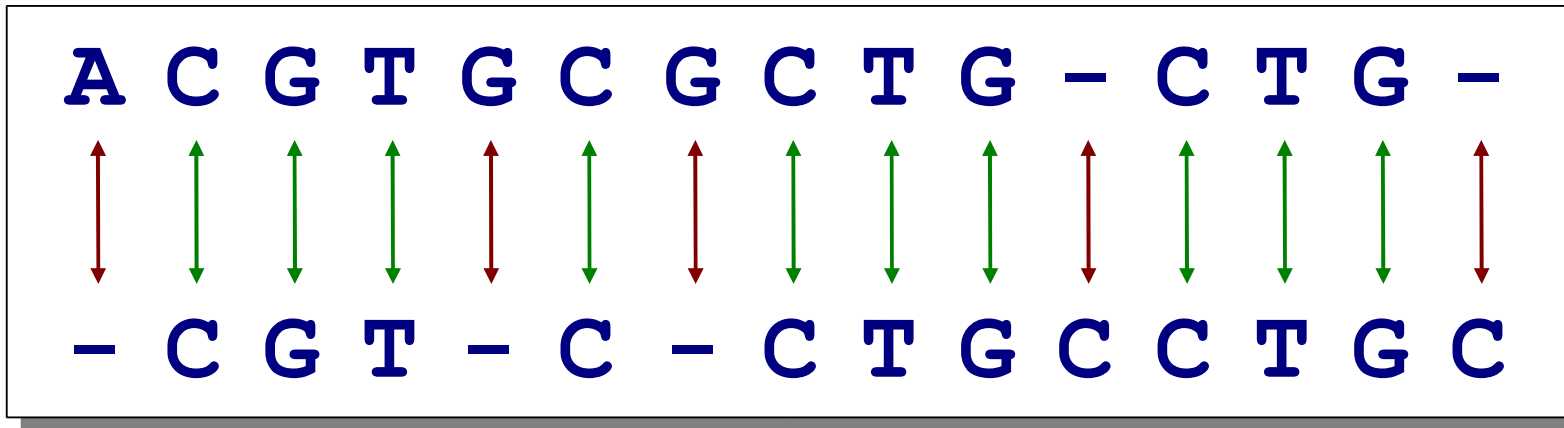The notion of an *alignment* helps us visualize this:



Alignment permits us to incorporate "spaces" (represented by a dash) in one or both sequences to make them the same length.

A C G T G C – G C T G C T G –
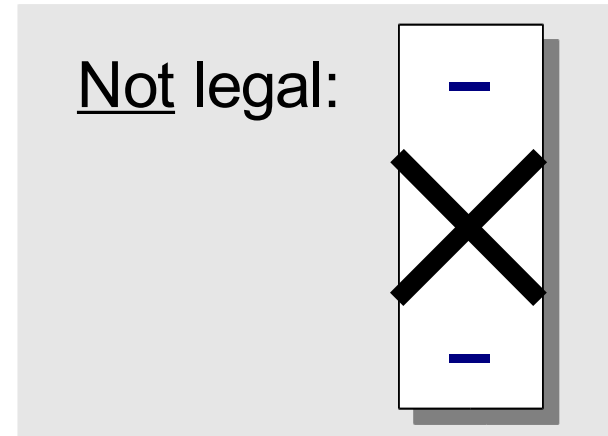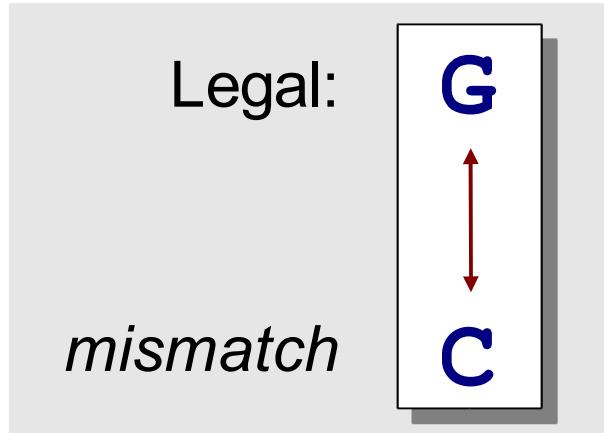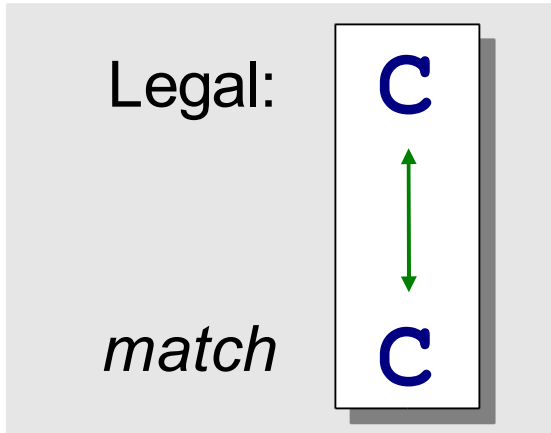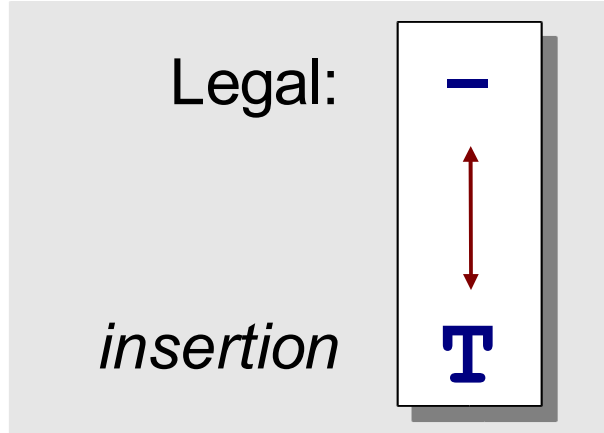
– C G T C C T G C – – C T G C

This alignment has 6 mismatches.  Is it the best possible?

A C G T G C G C T G – C T G –

– C G T – C – C T G C C T G C

How can we find the best alignment?

LEHIGH
UNIVERSITY

First some ground-rules.

Legal: **A** ↕ **–**

*deletion*

Legal: **–** ↕ **T**

*insertion*

<u>Not</u> legal: **–** ✕ **–**

Legal: **C** ↕ **C**

*match*

Legal: **G** ↕ **C**

*mismatch*

LEHIGH
U N I V E R S I T Y

We can't afford to enumerate all possible alignments looking for the best one – that would be an exponential search.

Fortunately we don't have to.  The optimal alignment can be found using a technique known as *dynamic programming*.

Dynamic programming is based on the premise of computing the solutions to smaller subproblems first and then using these to solve successively larger problems until we have our answer.
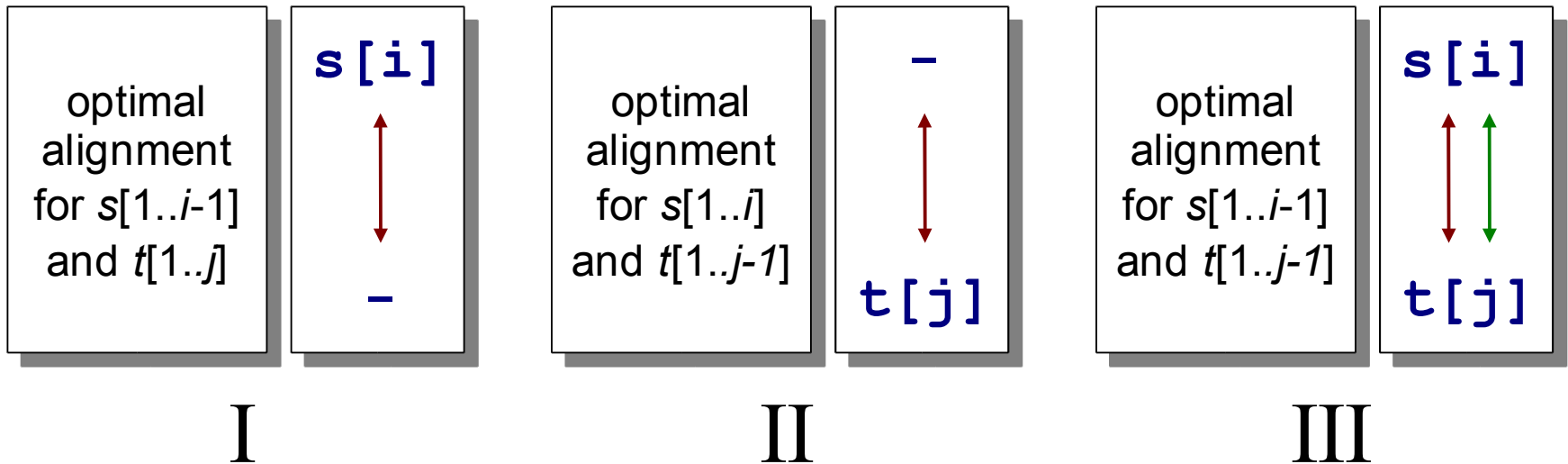
(Dynamic programming was invented by Richard Bellman in the 1950's.  Its application to sequence comparison came later, in the 1970's.)

http://fens.sabanciuniv.edu/msie/operations_research_50_years/anniversary/or50/1526-5463-2002-50-01-0048.pdf

Given two sequences *s* and *t*, consider what's required to compute optimal alignment for prefixes *s*[1..*i*] and *t*[1..*j*].  Based on our rules for alignments, there are three possible cases:

| | | | | | |
|---|---|---|---|---|---|
| optimal alignment for *s*[1..*i*-1] and *t*[1..*j*] | `s[i]`  `–` | optimal alignment for *s*[1..*i*] and *t*[1..*j-1*] | `–`  `t[j]` | optimal alignment for *s*[1..*i*-1] and *t*[1..*j-1*] | `s[i]`  `t[j]` |

<div align="center">

I        II        III

</div>

So, assuming we've already computed solutions for all shorter prefixes, we can compute the alignment for *s*[1..*i*] and *t*[1..*j*].

# *Sequence comparison: the basic algorithm*

Conceptually, this might look something like this:

$$\begin{array}{c}\text{optimal alignment at} \\ s[1..i] \text{ and } t[1..j]\end{array} = \max \left\{ \begin{array}{c}\text{optimal alignment at} \\ s[1..i\text{-}1] \text{ and } t[1..j] \\ + \\ \text{cost of deleting } s[i] \\ \\ \text{optimal alignment at} \\ s[1..i] \text{ and } t[1..j\text{-}1] \\ + \\ \text{cost of inserting } t[j] \\ \\ \text{optimal alignment at} \\ s[1..i\text{-}1] \text{ and } t[1..j\text{-}1] \\ + \\ \text{cost of matching } s[i] \text{ and } t[j]\end{array} \right.$$

Here we assume that deletions, insertions, and mismatches have negative costs, while matches have positive cost.

LEHIGH UNIVERSITY

# *Sequence comparison: the basic algorithm*

This computation can be viewed as building a 2-D matrix:

| ε | *s t r i n g   t* |
|---|---|
| 0 | ← cost of inserting t |

*string s* — cost of deleting s

$$a[i,j] \;=\; max \begin{cases} a[i\text{-}1,j] - 2 \\[6pt] a[i,j\text{-}1] - 2 \\[6pt] a[i\text{-}1,j\text{-}1] + p(i,j) \end{cases}$$

where -2 is the cost of an indel, and *p(i,j)* is the cost of a match/mismatch (using your book's notation).

LEHIGH UNIVERSITY

Stated more generally, say that our two sequences are:

$$s[1]s[2]s[3]...s[m] \qquad t[1]t[2]t[3]...t[n]$$

Then:

$$a[0,0] = 0$$

$$a[i,0] = a[i-1,j] + c_{del}(s[i]) \qquad 1 \le i \le m$$

$$a[0,j] = a[0,j-1] + c_{ins}(t[j]) \qquad 1 \le j \le n$$

And:

$$a[i,j] = max \begin{cases} a[i-1,j] + c_{del}(s[i]) \\ a[i,j-1] + c_{ins}(t[j]) \\ a[i-1,j-1] + c_{sub}(s[i],t[j]) \end{cases} \qquad 1 \le i \le m, \ 1 \le j \le n$$

Where $c_{del}$, $c_{ins}$, and $c_{sub}$ are the costs of a deletion, an insertion, and a substitution, respectively.

# *Sequence comparison: the basic algorithm*

Example:      say that $c_{del} = -1$, $c_{ins} = -1$,

$c_{sub} = -1$ if mismatch and $+1$ if match

|     | C   | A   | T   |
|-----|-----|-----|-----|
| **0** | **−1** | **−2** | **−3** |
| **A**   **−1** | **−1** | **0** | **−1** |
| **C**   **−2** | **0** | **−1** | **−1** |
| **G**   **−3** | **−1** | **−1** | **−2** |
| **T**   **−4** | **−2** | **−2** | **0** |

LEHIGH
UNIVERSITY

Computation can progress in a number of ways:



*or*

*or*

For sequences of length *m* and *n*,

$$s[1]s[2]s[3]...s[m] \qquad\qquad t[1]t[2]t[3]...t[n]$$

Computation time = $O(mn)$, space = $O(mn)$.

LEHIGH
UNIVERSITY.

We started with the notion of alignment. How do we get this?
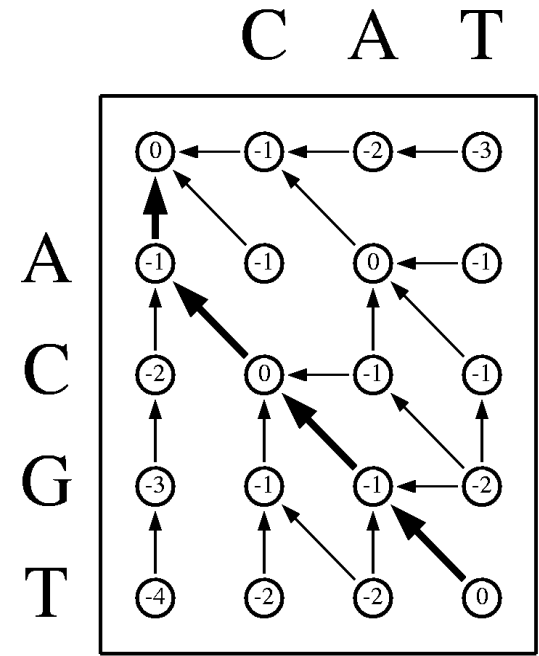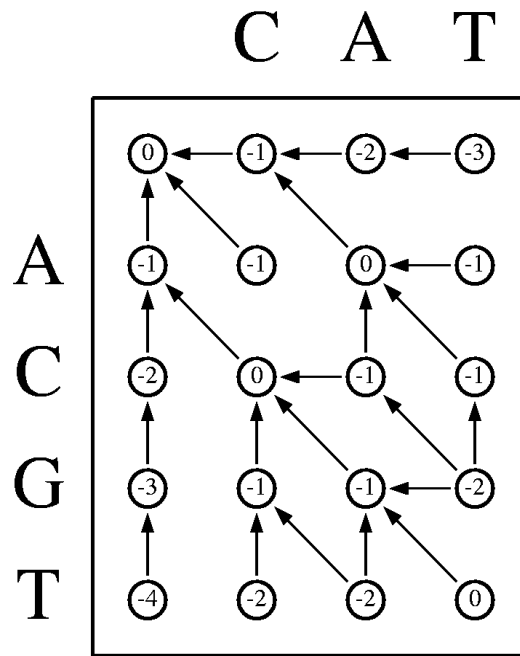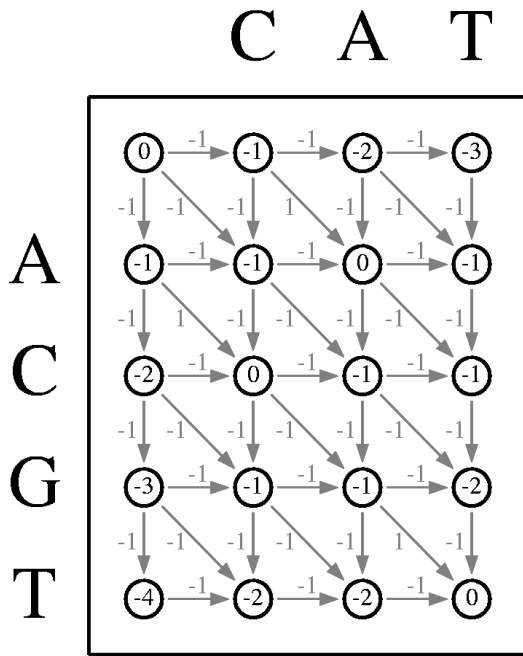
By keeping track of optimal decisions made during algorithm,
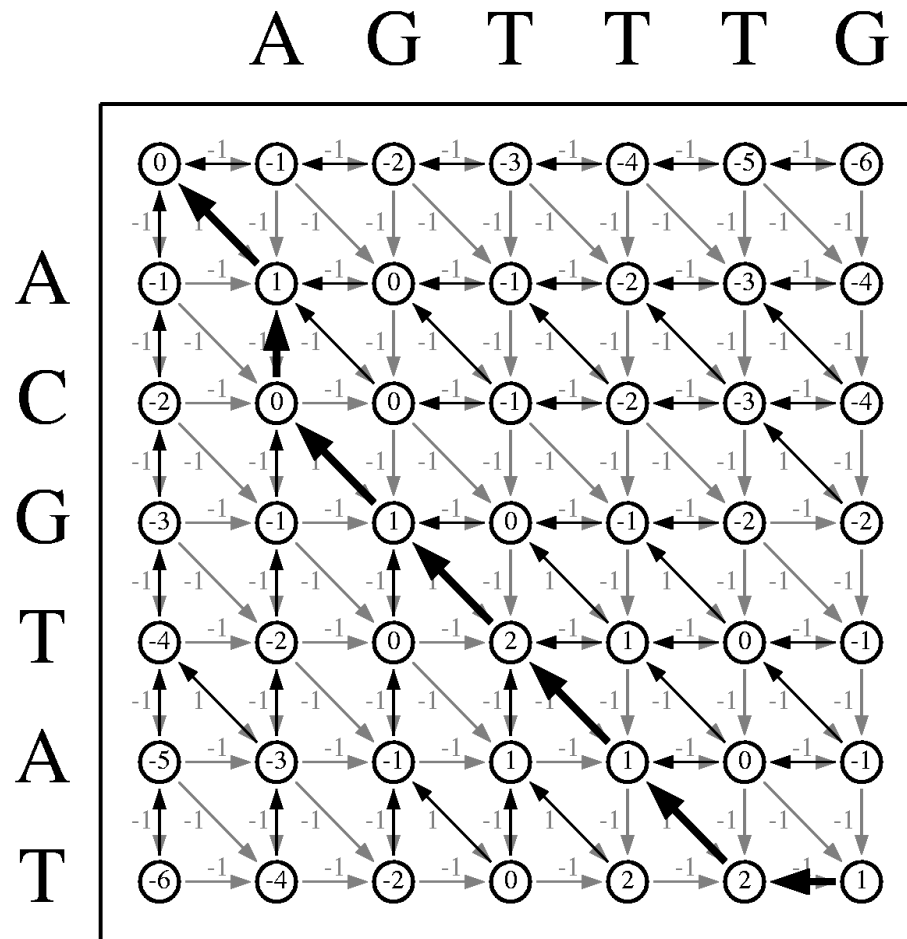


and then tracing back optimal path.

(May not be unique).

Comparing **ACGT** vs. **CAT**:
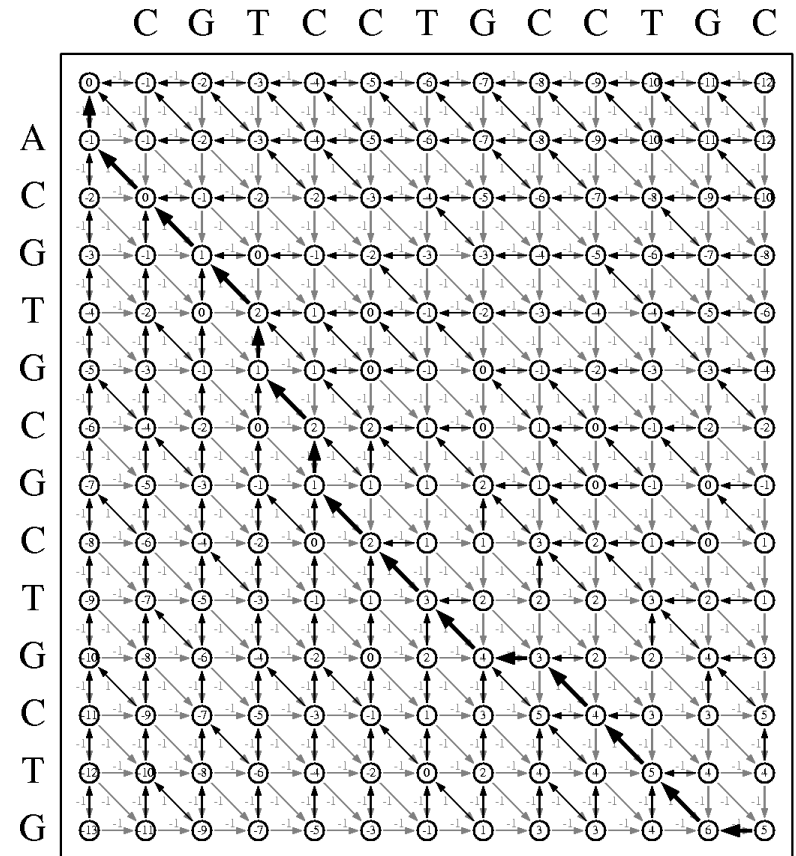
Comparing **ACGTAT** vs. **AGTTTG**:

Comparing **ACGTGCGCTGCTG**
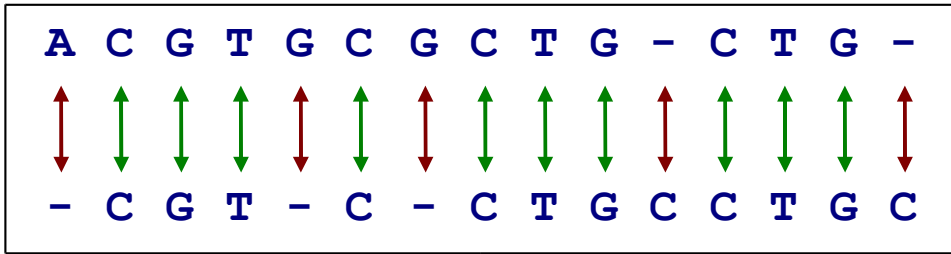vs. **CGTCCTGCCTGC**:

Recall the trace we saw earlier:

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

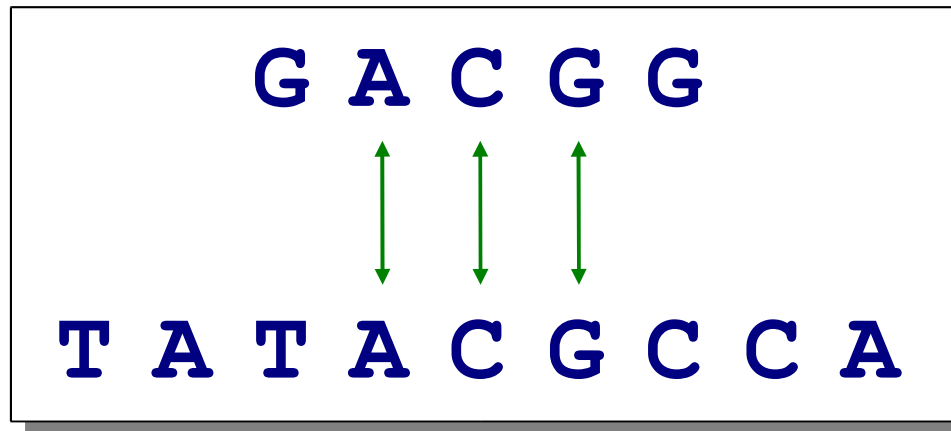How do we know this agorithm finds an optimal alignment?

Your textbook finesses this point, but basically it hinges on:

- Bellman's *principle of optimality* for dynamic programming.
- The cost functions behaving properly (i.e., they must satisfy the *triangle inequality*).

This could be an interesting topic for a lecture or final paper ...

So far, we have assumed the sequences must be matched in their entirety.  But this ignores interesting similarities that might be present at the subsequence level:



Fortunately, a slight modification of the original algorithm can handle this.

As before, say that our two sequences are:

$$s[1]s[2]s[3]...s[m] \qquad t[1]t[2]t[3]...t[n]$$

Then:
$$a[0,0] = 0$$
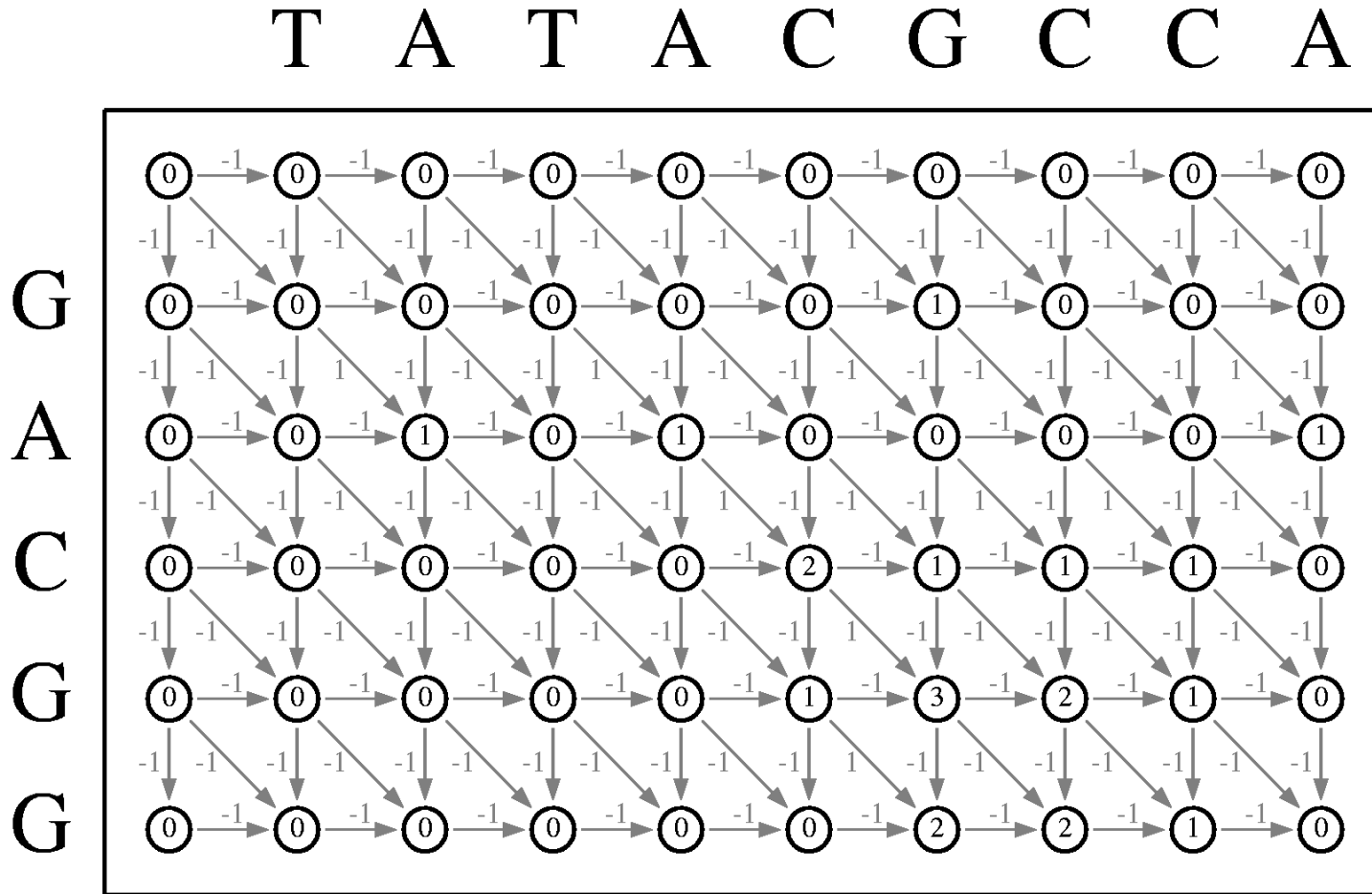$$a[i,0] = 0 \qquad 1 \leq i \leq m$$
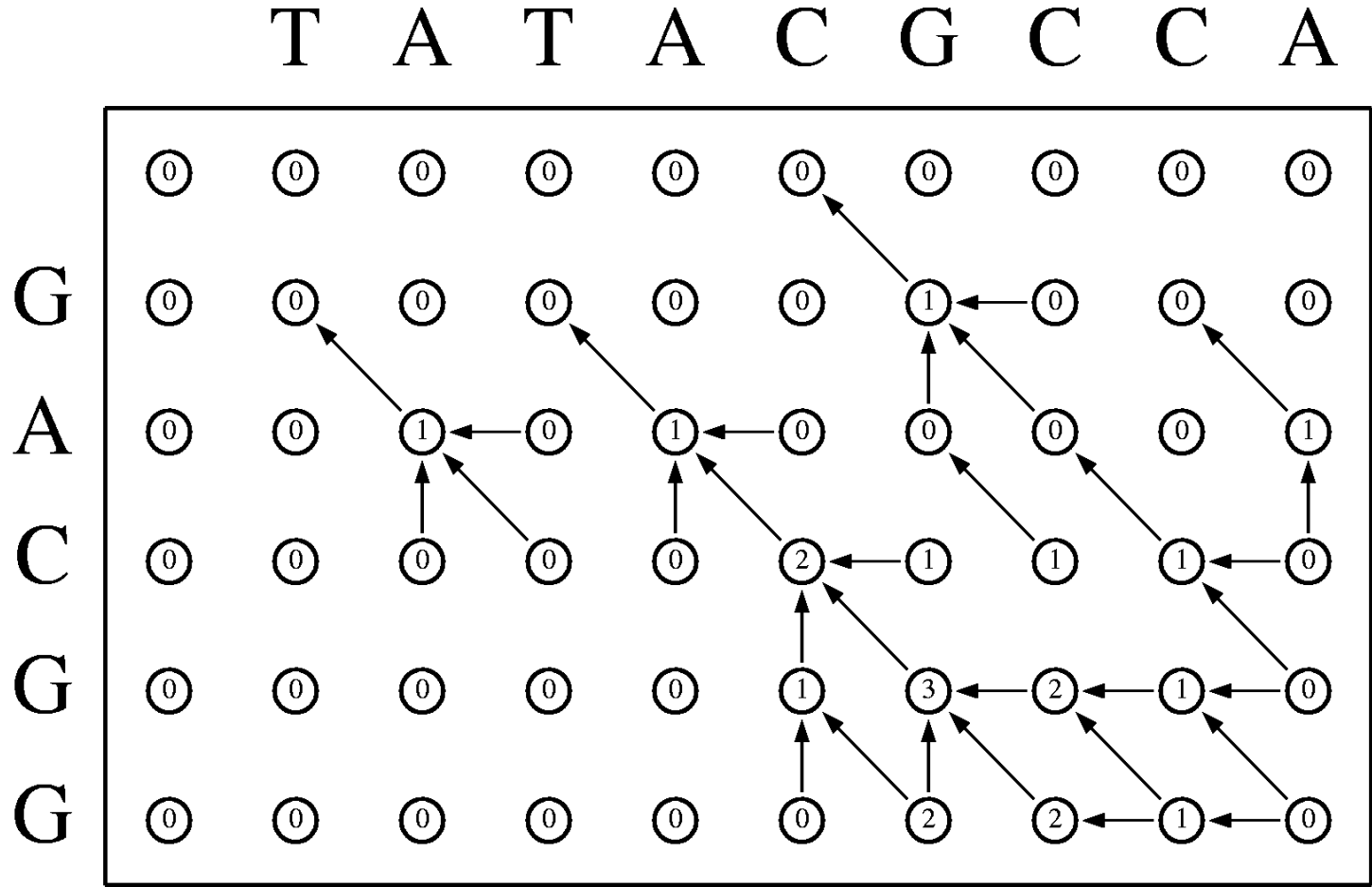$$a[0,j] = 0 \qquad 1 \leq j \leq n$$

And:

$$a[i,j] = max \begin{cases} a[i\text{-}1,j] + c_{del}(s[i]) \\ a[i,j\text{-}1] + c_{ins}(t[j]) \\ a[i\text{-}1,j\text{-}1] + c_{sub}(s[i],t[j]) \\ 0 \end{cases} \qquad 1 \leq i \leq m, \ 1 \leq j \leq n$$
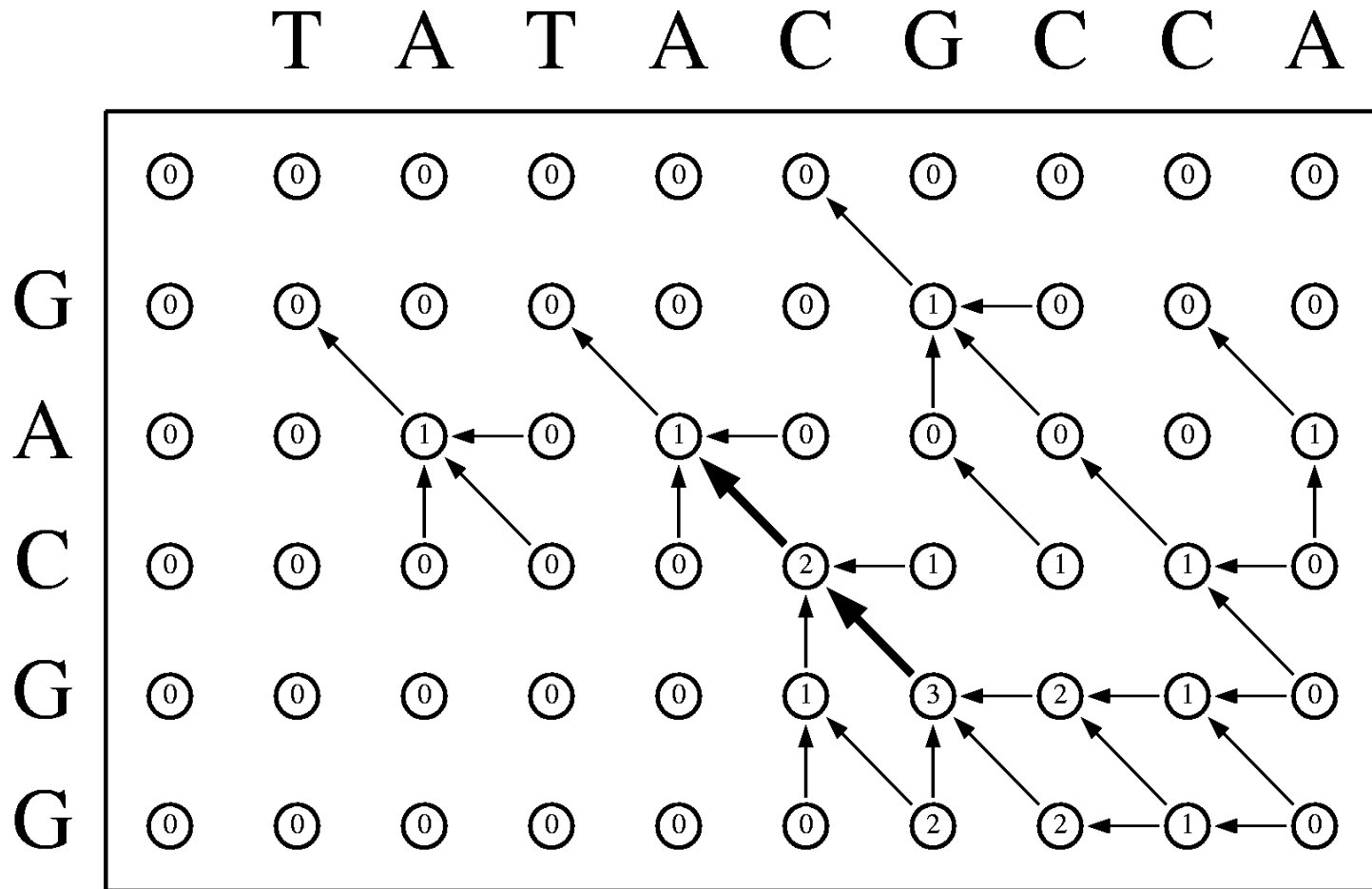
When done, we search the matrix for the largest value.

LEHIGH
U N I V E R S I T Y

LEHIGH
UNIVERSITY

LEHIGH
UNIVERSITY

Readings for next time:

- Section 3.3 in your textbook.
- "The emerging landscape of bioinformatics software systems" by L.S. Heath and N. Ramakrishnan, *IEEE Computer*, July 2002, pp. 41-45.  (Available online or in Blackboard.)

Remember:

- Come to class prepared to discuss what you have read.
- Check Blackboard regularly for updates.

LEHIGH
UNIVERSITY