
Introduction to Algorithms

Kiyoko F. Aoki-Kinoshita

Computational problems

- A computational problem specifies an input-output relationship
 - What does the input look like?
 - What should the output be for each input?
 - Example:
 - Input: an integer number N
 - Output: Is the number prime?
 - Example:
 - Input: A list of names of people
 - Output: The same list sorted alphabetically
 - Example:
 - Input: A picture in digital format
 - Output: An English description of what the picture shows
-

Algorithms

- An algorithm is an exact specification of how to solve a computational problem
- An algorithm must specify every step completely, so a computer can implement it without any further “understanding”
- An algorithm must work for all possible inputs of the problem.
- Algorithms must be:
 - Correct: For each input, terminate and produce an appropriate output
 - Efficient: run as quickly as possible, and use as little memory as possible – more about this later
- There can be many different algorithms for each computational problem.

Describing Algorithms

- Algorithms can be implemented in any programming language
- Usually we use “pseudo-code” to describe algorithms

Testing whether input N is prime:

- ```
For j = 2 .. N-1
 If the remainder of j/N is 0
 Output "N is composite" and halt
Output "N is prime"
```

- In this course we will just describe algorithms in Perl and pseudocode
-

---

# Greatest Common Divisor

- The first algorithm “invented” in history was Euclid’s algorithm for finding the greatest common divisor (GCD) of two natural numbers
  - **Definition:** The GCD of two natural numbers  $x, y$  is the largest integer  $j$  that divides both evenly (with remainder 0).
  - **The GCD Problem:**
    - Input: natural numbers  $x, y$
    - Output:  $GCD(x,y)$  – their GCD
-

# Euclid's GCD Algorithm

```
sub gcd {
 my ($x, $y) = @_ ; // retrieve input x and y
 while ($y != 0) { // while y is not equal to 0
 $t = $x % $y; // get the modulus of x and y
 $x = $y; // replace x by y
 $y = $t; // replace y by t
 }
 return $x; // return the result (gcd of x and y)
}

print gcd(14, 21), "\n";
```

# Euclid's GCD Algorithm – sample

```
while ($y != 0) { // while y is not equal to 0
 $t = $x % $y; // get the modulus of x and y
 $x = $y; // replace x by y
 $y = $t; // replace y by t
}
```

## Example: Computing GCD(48,120)

|                | t  | x   | y   |
|----------------|----|-----|-----|
| After 0 rounds | -- | 72  | 120 |
| After 1 round  | 72 | 120 | 72  |
| After 2 rounds | 48 | 72  | 48  |
| After 3 rounds | 24 | 48  | 24  |
| After 4 rounds | 0  | 24  | 0   |

Output: 24

# Termination of Euclid's Algorithm

- Why does this algorithm terminate?
  - After any iteration we have that  $x > y$  since the new value of  $y$  is the remainder of the division by the new value of  $x$ .
  - In further iterations, we replace  $(x, y)$  with  $(y, x\%y)$ , and  $x\%y < x$ , thus the numbers decrease in each iteration.
  - Formally, the value of  $xy$  decreases at each iteration (except, maybe, the first one). When it reaches 0, the algorithm must terminate.

```
sub gcd {
 my ($x, $y) = @_; // retrieve input x and y
 while ($y != 0) { // while y is not equal to 0
 $t = $x % $y; // get the modulus of x and y
 $x = $y; // replace x by y
 $y = $t; // replace y by t
 }
 return $x; // return the result (gcd of x and y)
}
```

---

# Introduction to Algorithms

---

Running Time Analysis

---

# How fast will your program run?

- The running time of your program will depend upon:
    - ❑ The algorithm
    - ❑ The input
    - ❑ Your implementation of the algorithm in a programming language
    - ❑ The compiler you use
    - ❑ The operating system (OS) on your computer
    - ❑ Your computer hardware
    - ❑ Maybe other things: temperature outside; other programs on your computer; ...
  - Our Motivation: analyze the running time of an algorithm as a function of only simple parameters of the input.
-

---

# Basic idea: counting operations

- Each algorithm performs a sequence of basic operations:
    - Arithmetic:  $(\text{low} + \text{high})/2$
    - Comparison:  $\text{if} ( x > 0 ) \dots$
    - Assignment:  $\text{temp} = x$
    - Branching:  $\text{while} ( y \neq 0 ) \{ \dots \}$
    - ...
  - Idea: count the number of basic operations performed on the input.
  - Difficulties:
    - Which operations are basic?
    - Not all operations take the same amount of time.
    - Operations take different times with different hardware or compilers
-

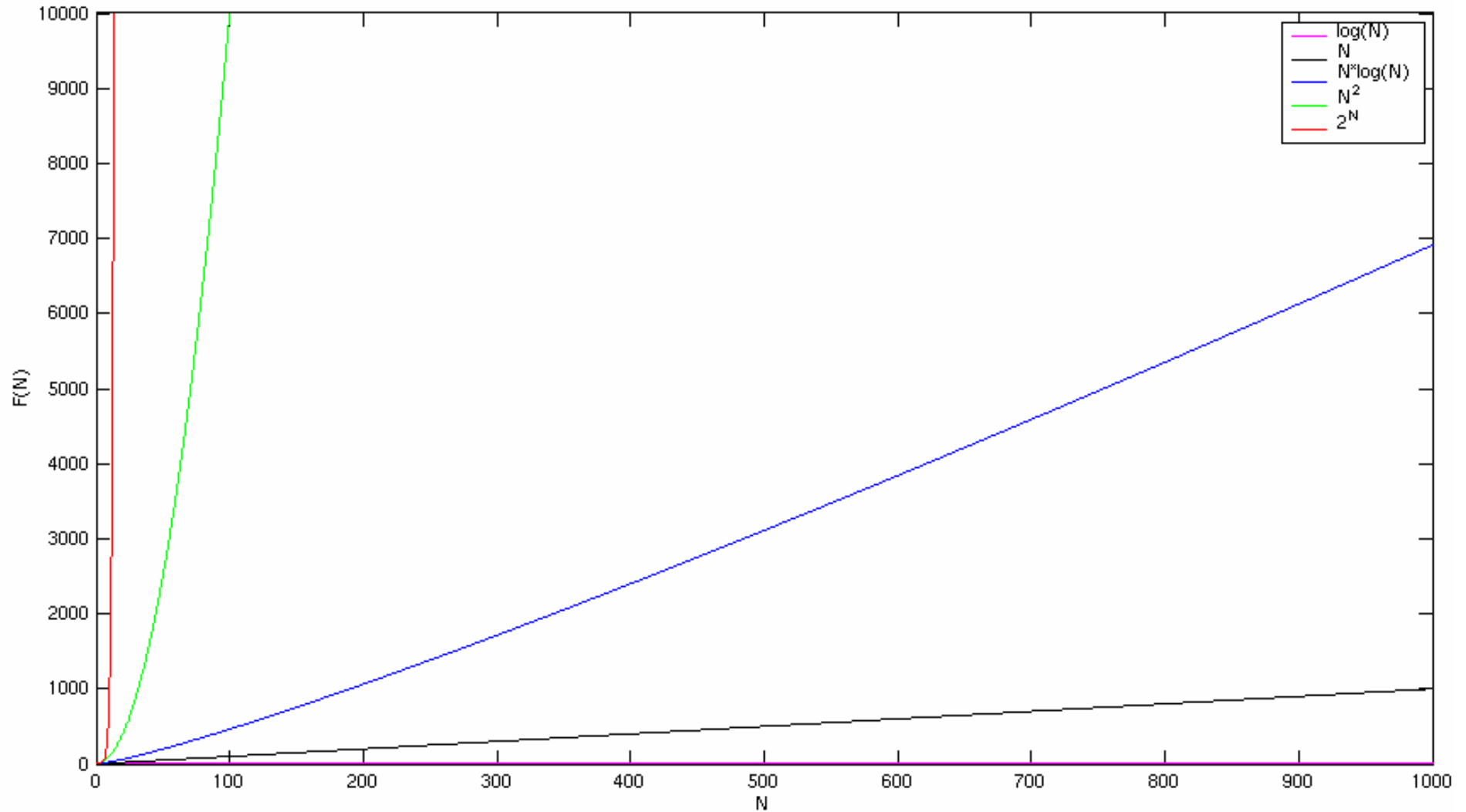
# Asymptotic running times

- Operation counts are only problematic in terms of constant factors.
- The general form of the function describing the running time is invariant over hardware, languages or compilers!

```
sub myMethod{
 my $N = shift @_;
 my $sq = 0;
 for($j=0; $j <$N ; $j++)
 for($k=0; $k<$N ; $k++)
 $sq++;
 return $sq;
}
```

- Running time is “about”  $N^2$ .
- We use “Big-O” notation, and say that the running time is  $O(N^2)$

# Asymptotic behavior of functions



# Mathematical Formalization

- Definition: Let  $f$  and  $g$  be functions from the natural numbers to the natural numbers. We write  $f=O(g)$  if there exists a constant  $c$  such that for all  $n$ :  $f(n) \leq cg(n)$ .

$$f=O(g) \Leftrightarrow \exists c \forall n: f(n) \leq cg(n)$$

- This is a mathematically formal way of ignoring constant factors, and looking only at the “shape” of the function.
- $f=O(g)$  should be considered as saying that “ $f$  is at most  $g$ , up to constant factors”.
- We usually will have  $f$  be the running time of an algorithm and  $g$  a nicely written function. E.g. The running time of the previous algorithm was  $O(N^2)$ .

---

# Asymptotic analysis of algorithms

- We usually embark on an *asymptotic worst case* analysis of the running time of the algorithm.
  - Asymptotic:
    - Formal, exact, depends only on the algorithm
    - Ignores constants
    - Applicable mostly for large input sizes
  - Worst Case:
    - Bounds on running time must hold for *all* inputs.
    - Thus the analysis considers the worst-case input.
    - Sometimes the “average” performance can be much better
    - Real-life inputs are rarely “average” in any formal sense
-

# The running time of Euclid's GCD Algorithm

- How fast does Euclid's algorithm terminate?
  - After the first iteration we have that  $x > y$ . In each iteration, we replace  $(x, y)$  with  $(y, x \% y)$ .
  - In an iteration where  $x > 1.5y$  then  $x \% y < y < 2x/3$ .
  - In an iteration where  $x \leq 1.5y$  then  $x \% y \leq y/2 < 2x/3$ .
  - Thus, the value of  $xy$  decreases by a factor of at least  $2/3$  each iteration (except, maybe, the first one).

```
sub gcd {
 my ($x, $y) = @_; // retrieve input x and y
 while ($y != 0) { // while y is not equal to 0
 $t = $x % $y; // get the modulus of x and y
 $x = $y; // replace x by y
 $y = $t; // replace y by t
 }
 return $x; // return the result (gcd of x and y)
}
```

# The running time of Euclid's Algorithm

- Theorem: Euclid's GCD algorithm runs in time  $O(N)$ , where  $N$  is the input length ( $N = \log_2 x + \log_2 y$ ).
- Proof:
  - Every iteration of the loop (except maybe the first) the value of  $xy$  decreases by a factor of at least  $2/3$ . Thus after  $k+1$  iterations the value of  $xy$  is at most  $(2/3)^k$  the original value.
  - Thus the algorithm must terminate when  $k$  satisfies:  $xy(2/3)^k < 1$  (for the original values of  $x, y$ ).
  - Thus the algorithm runs for at most  $1 + \log_{3/2} xy$  iterations.
  - Each iteration has only a constant  $L$  number of operations, thus the total number of operations is at most  $(1 + \log_{3/2} xy)L$
  - Formally,  $(1 + \log_{3/2} xy)L \leq L(1 + 2\log_2 x + 2\log_2 y) \leq 3LN$
  - Thus the running time is  $O(N)$ .

---

# Introduction to Algorithms

---

Recursion

---

# Designing Algorithms

- There is no single recipe for inventing algorithms
- There are basic rules:
  - Understand your problem well – may require much mathematical analysis!
  - Use existing algorithms (reduction) or algorithmic ideas
- There is a single basic algorithmic technique:

## Divide and Conquer

- In its simplest (and most useful) form it is simple induction
    - In order to solve a problem, solve a similar problem of smaller size
  - The key conceptual idea:
    - Think only about how to use the smaller solution to get the larger one
    - Do not worry about how to solve the smaller problem (it will be solved using an even smaller one)
-

---

# Recursion

- A recursive method is a method that contains a call to itself
  - Technically:
    - All modern computing languages allow writing methods that call themselves
    - We will discuss how this is implemented later
  - Conceptually:
    - This allows programming in a style that reflects divide-n-conquer algorithmic thinking
    - At the beginning recursive programs are confusing – after a while they become clearer than non-recursive variants
-

# Factorial

```
sub factorial {
 my $n = shift @_; // retrieve input
 if ($n == 0) {
 return 1; // if input is 0, return 1
 } else {
 // otherwise, compute the factorial of $n-1,
 // multiply it by $n and return the product
 return $n * factorial($n-1);
 }
}

print "5! = ", factorial(5), "\n";
```

---

# Elements of a recursive program

- **Basis:** a case that can be answered without using further recursive calls
    - In our case: `if ($n==0) { return 1; }`
  - **Creating the smaller problem, and invoking a recursive call on it**
    - In our case: `factorial($n-1)`
  - **Finishing to solve the original problem**
    - In our case: `return $n; //solution of recursive call`
-

# Tracing the factorial method

```
print "5! = ", factorial(5), "\n";

 5 * factorial(4)
 4 * factorial(3)
 3 * factorial(2)
 2 * factorial(1)
 1 * factorial(0)
 return 1
 return 1
 return 2
 return 6
 return 24
 return 120
```

# Correctness of factorial method

- Theorem: For every positive integer  $n$ , factorial ( $n$ ) returns the value  $n!$ .
- Proof: By induction on  $n$ :
- Basis: for  $n=0$ , factorial ( $0$ ) returns  $1=0!$ .
- Induction step: When called on  $n>1$ , factorial calls factorial ( $n-1$ ), which by the induction hypothesis returns  $(n-1)!$ . The returned value is thus  $n*(n-1)!=n!$ .

# Raising to power – take 1

```
sub power {
 my ($x, $n) = @_; // retrieve the input
 if ($n == 0) { // if $n is 0, return 1
 return 1.0;
 }
 // otherwise, return $x multiplied by the
 // result of power of x to the (n-1)th
 return $x * power($x, $n-1);
}

print "3^9 = ", power(3, 9), "\n";
```

---

# Running time analysis

- Simplest way to calculate the running time of a recursive program is to add up the running times of the separate levels of recursion.
  - In the case of the power method:
    - There are  $n+1$  levels of recursion
      - $\text{power}(x,n), \text{power}(x,n-1), \text{power}(x, n-2), \dots \text{power}(x,0)$
    - Each level takes  $O(1)$  steps
    - Total time =  $O(n)$
-

# Raising to power – take 2

```
sub power2 {
 my ($x, $n) = @_;
 if ($n == 0) {
 return 1.0;
 }
 if ($n%2 == 0) {
 my $t = power2($x, $n/2);
 return $t*$t;
 }
 return $x * power2($x, $n-1);
}
```

# Analysis

- Theorem: For any  $x$  and positive integer  $n$ , the power method returns  $x^n$ .
- Proof: by complete induction on  $n$ .
  - Basis: For  $n=0$ , we return 1.
  - If  $n$  is even, we return  $\text{power}(x, n/2) * \text{power}(x, n/2)$ . By the induction hypothesis  $\text{power}(x, n/2)$  returns  $x^{n/2}$ , so we return  $(x^{n/2})^2 = x^n$ .
  - If  $n$  is odd, we return  $x * \text{power}(x, n-1)$ . By the induction hypothesis  $\text{power}(x, n-1)$  returns  $x^{n-1}$ , so we return  $x \cdot x^{n-1} = x^n$ .
- The running time is now  $O(\log n)$ :
  - After 2 levels of recursion  $n$  has decreased by a factor of at least two (since either  $n$  or  $n-1$  is even, in which case the recursive call is with  $n/2$ )
  - Thus we reach  $n==0$  after at most  $2\log_2 n$  levels of recursion
  - Each level still takes  $O(1)$  time.

---

# Introduction to Algorithms

---

Algorithms for bioinformatics

---

# Bring in the Bioinformaticians

- Gene similarities between two genes with known and unknown function alert biologists to some possibilities
  - Computing a similarity score between two genes tells how likely it is that they have similar functions
  - **Dynamic programming** is a technique for revealing similarities between genes
  - The ***Change Problem*** is a good problem to introduce the idea of dynamic programming
-

# The Change Problem

**Goal**: Convert some amount of money  $M$  into given denominations, using the fewest possible number of coins

**Input**: An amount of money  $M$ , and an array of  $d$  denominations  $\mathbf{c} = (c_1, c_2, \dots, c_d)$ , in a decreasing order of value ( $c_1 > c_2 > \dots > c_d$ )

**Output**: A list of  $d$  integers  $i_1, i_2, \dots, i_d$  such that

$$c_1 i_1 + c_2 i_2 + \dots + c_d i_d = M$$

and  $i_1 + i_2 + \dots + i_d$  is minimal

# Change Problem: Example

Given the denominations 1, 3, and 5, what is the minimum number of coins needed to make change for a given value?

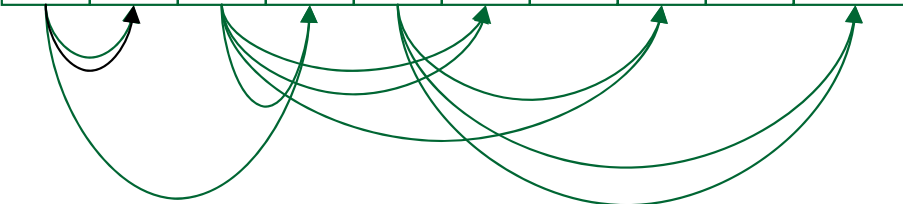
| Value          | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------------|---|---|---|---|---|---|---|---|---|----|
| Min # of coins | 1 |   | 1 |   | 1 |   |   |   |   |    |

Only one coin is needed to make change for the values 1, 3, and 5

# Change Problem: Example (cont'd)

Given the denominations 1, 3, and 5, what is the minimum number of coins needed to make change for a given value?

| Value          | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------------|---|---|---|---|---|---|---|---|---|----|
| Min # of coins | 1 | 2 | 1 | 2 | 1 | 2 |   | 2 |   | 2  |



However, two coins are needed to make change for the values 2, 4, 6, 8, and 10.

# Change Problem: Example (cont'd)

Given the denominations 1, 3, and 5, what is the minimum number of coins needed to make change for a given value?

| Value          | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------------|---|---|---|---|---|---|---|---|---|----|
| Min # of coins | 1 | 2 | 1 | 2 | 1 | 2 | 3 | 2 | 3 | 2  |

The diagram shows a table with two rows. The top row is labeled 'Value' and contains the numbers 1 through 10. The bottom row is labeled 'Min # of coins' and contains the numbers 1, 2, 1, 2, 1, 2, 3, 2, 3, 2. Arrows point from the 'Min # of coins' row to the 'Value' row, showing the relationship between the minimum number of coins and the value they represent. Specifically, arrows point from the '1' coin count to values 1, 3, and 5; from the '2' coin count to values 2, 4, 6, 8, and 10; and from the '3' coin count to values 7 and 9.

Lastly, three coins are needed to make change for the values 7 and 9

# Change Problem: Recurrence

This example is expressed by the following recurrence relation:

$$\mathit{minNumCoins}(M) = \mathbf{\min\ of} \left\{ \begin{array}{l} \mathit{minNumCoins}(M-1) + 1 \\ \mathit{minNumCoins}(M-3) + 1 \\ \mathit{minNumCoins}(M-5) + 1 \end{array} \right.$$

Given the denominations  $\mathbf{c}$ :  $c_1, c_2, \dots, c_d$ , the recurrence relation is:

$$\mathit{minNumCoins}(M) = \mathbf{\min\ of} \left\{ \begin{array}{l} \mathit{minNumCoins}(M-c_1) + 1 \\ \mathit{minNumCoins}(M-c_2) + 1 \\ \dots \\ \mathit{minNumCoins}(M-c_d) + 1 \end{array} \right.$$

# Change Problem: A Recursive Algorithm

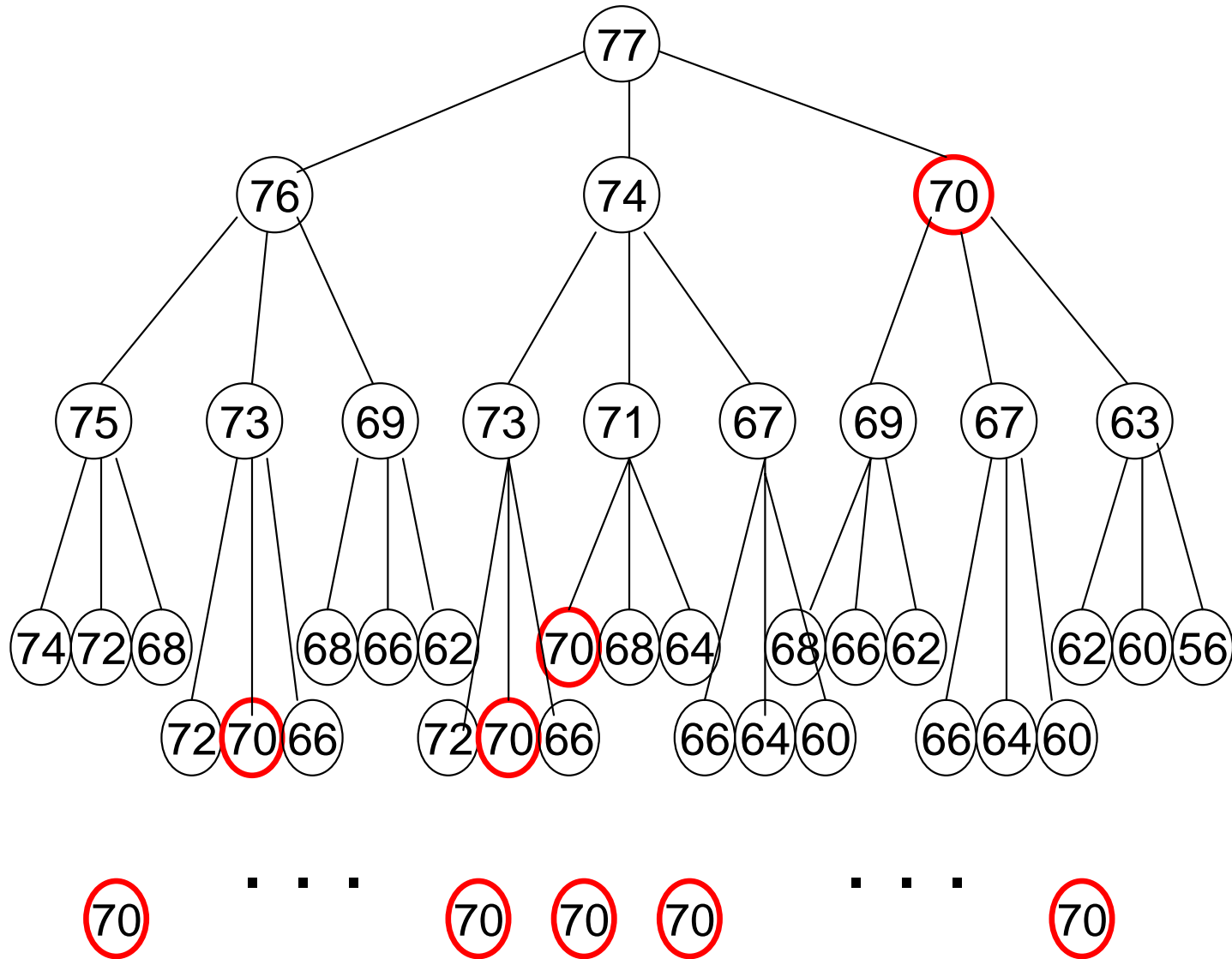
1. RecursiveChange( $M, c, d$ )
2.     if  $M = 0$
3.         return 0
4.      $bestNumCoins = \text{infinity}$
5.     for  $i = 1$  to  $d$
6.         if  $M \geq c_i$
7.              $numCoins = \text{RecursiveChange}(M - c_i, c, d)$
8.             if  $numCoins + 1 < bestNumCoins$
9.                  $bestNumCoins = numCoins + 1$
10.     return  $bestNumCoins$

---

# RecursiveChange Is Not Efficient

- It recalculates the optimal coin combination for a given amount of money repeatedly
  - i.e.,  $M = 77$ ,  $c = (1,3,7)$ :
    - Optimal coin combo for 70 cents is computed **9** times!
-

# The RecursiveChange Tree



---

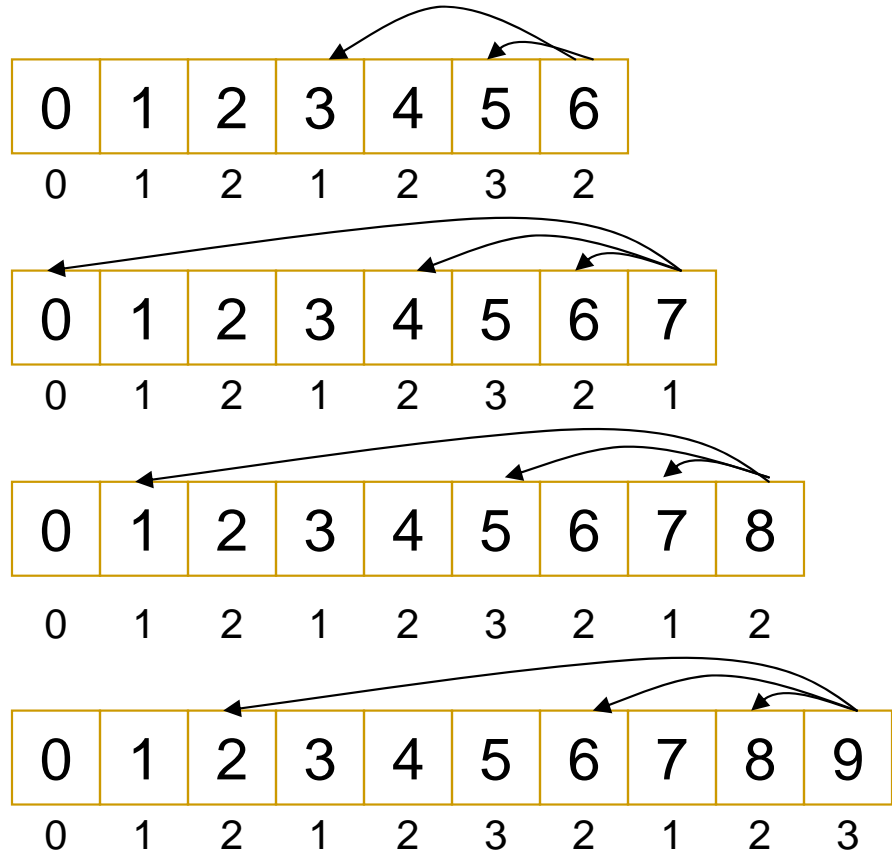
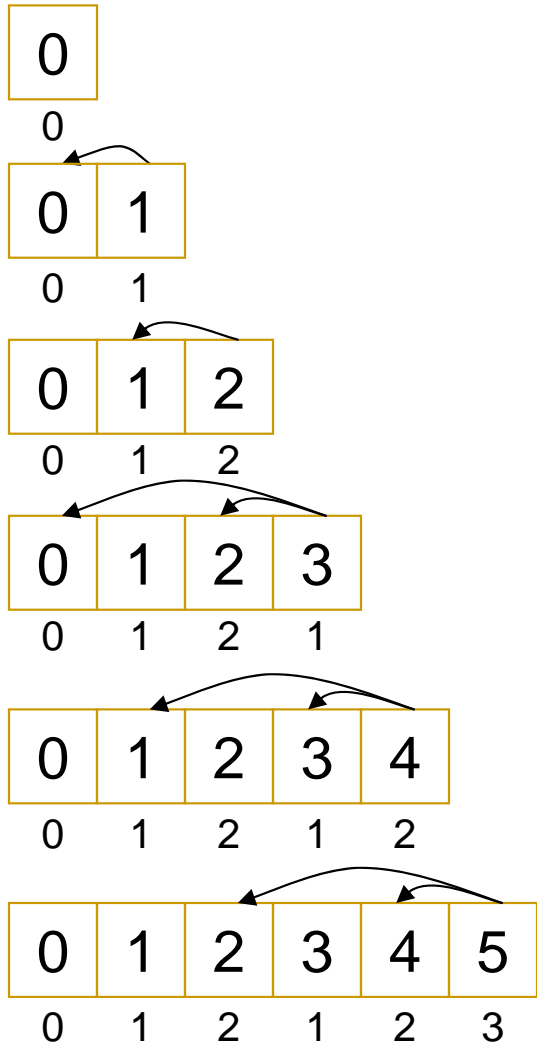
# We Can Do Better

- We're re-computing values in our algorithm more than once
  - Save results of each computation for 0 to  $M$
  - This way, we can do a reference call to find an already computed value, instead of re-computing each time
  - Running time becomes  $M^* d$ , where  $M$  is the value of money and  $d$  is the number of denominations
-

# The Change Problem: Dynamic Programming

1. DPChange( $M, c, d$ )
2.  $bestNumCoins_0 = 0$
3. for  $m = 1$  to  $M$
4.      $bestNumCoins_m = \text{infinity}$
5.     for  $i = 1$  to  $d$
6.         if  $m \geq c_i$
7.             if  $bestNumCoins_{m - c_i} + 1 < bestNumCoins_m$
8.                  $bestNumCoins_m = bestNumCoins_{m - c_i} + 1$
9.     return  $bestNumCoins_M$

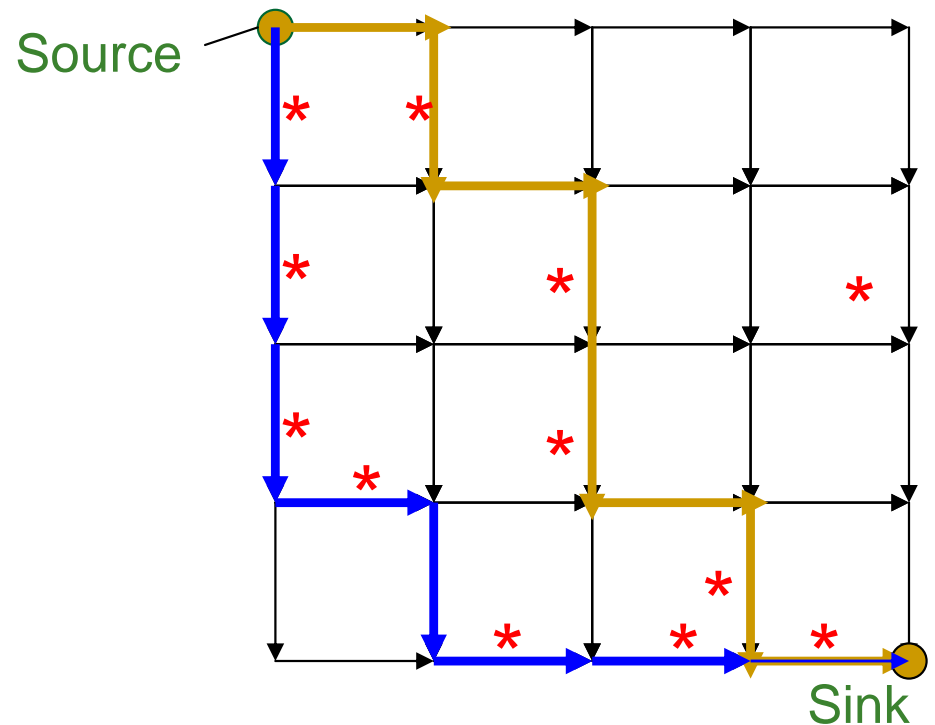
# DPChange: Example



$$\mathbf{c} = (1, 3, 7)$$
$$\mathbf{M} = 9$$

# Manhattan Tourist Problem (MTP)

Imagine seeking a path (from source to sink) to travel (only eastward and southward) with the most number of attractions (\*) in the Manhattan grid



---

## Manhattan Tourist Problem: Formulation

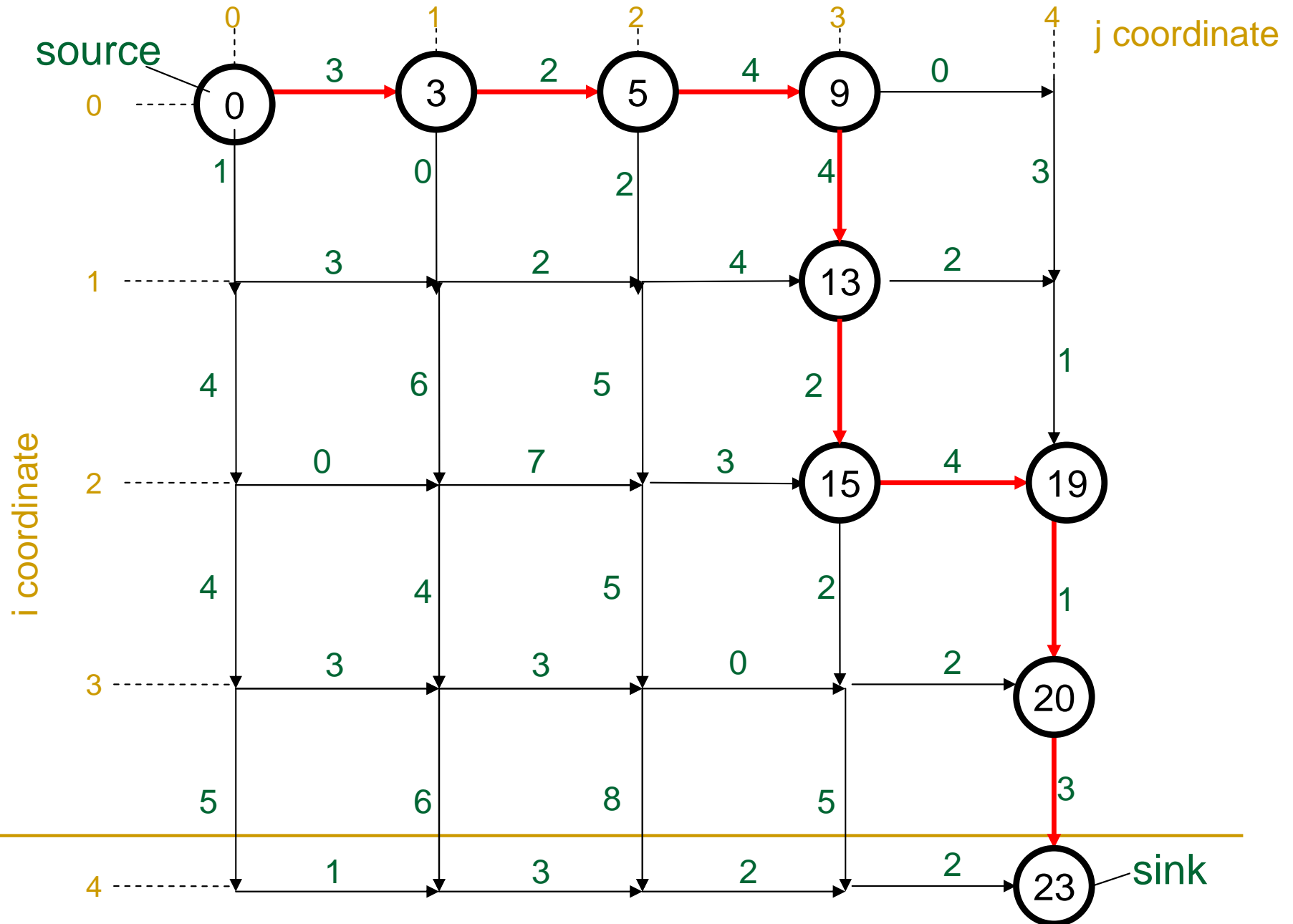
Goal: Find the longest path in a weighted grid.

Input: A weighted grid  $\mathbf{G}$  with two distinct vertices, one labeled “*source*” and the other labeled “*sink*”

Output: A longest path in  $\mathbf{G}$  from “*source*” to “*sink*”

---

# MTP: An Example



# MTP: Simple Recursive Program

MT( $n, m$ )

if  $n=0$  or  $m=0$

return  $MT(n, m)$

$x = MT(n-1, m) +$

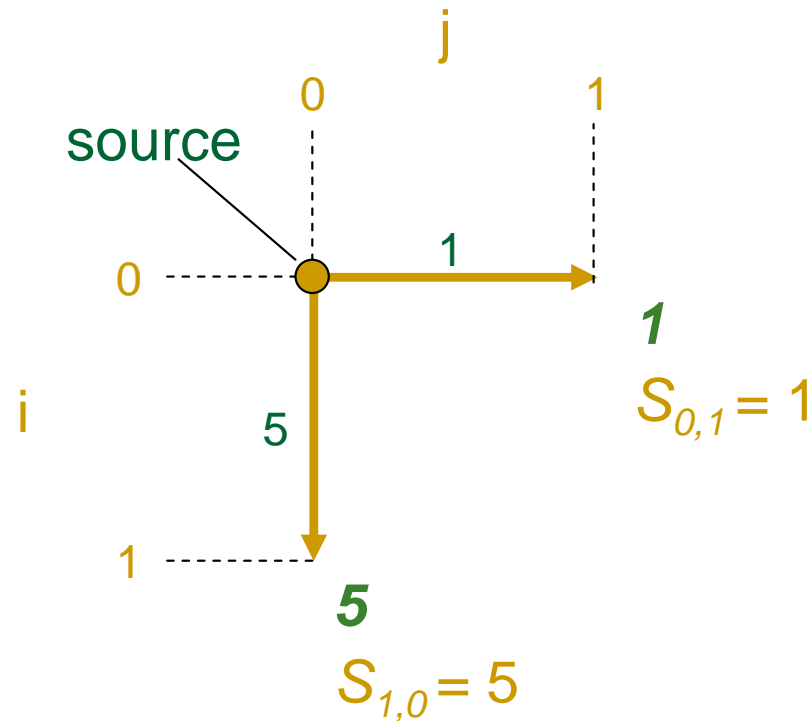
length of the edge from  $(n-1, m)$  to  $(n, m)$

$y = MT(n, m-1) +$

length of the edge from  $(n, m-1)$  to  $(n, m)$

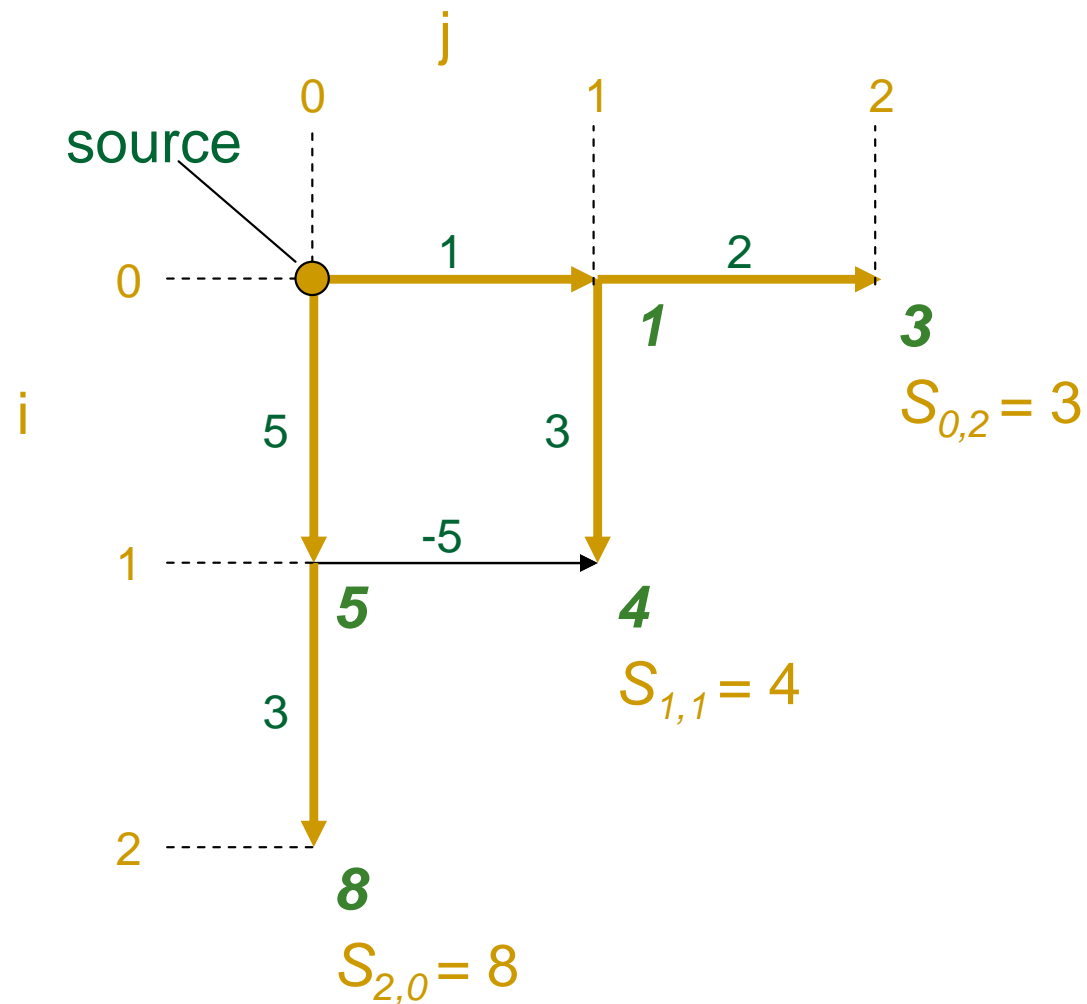
return  $\max\{x, y\}$

# MTP: Dynamic Programming

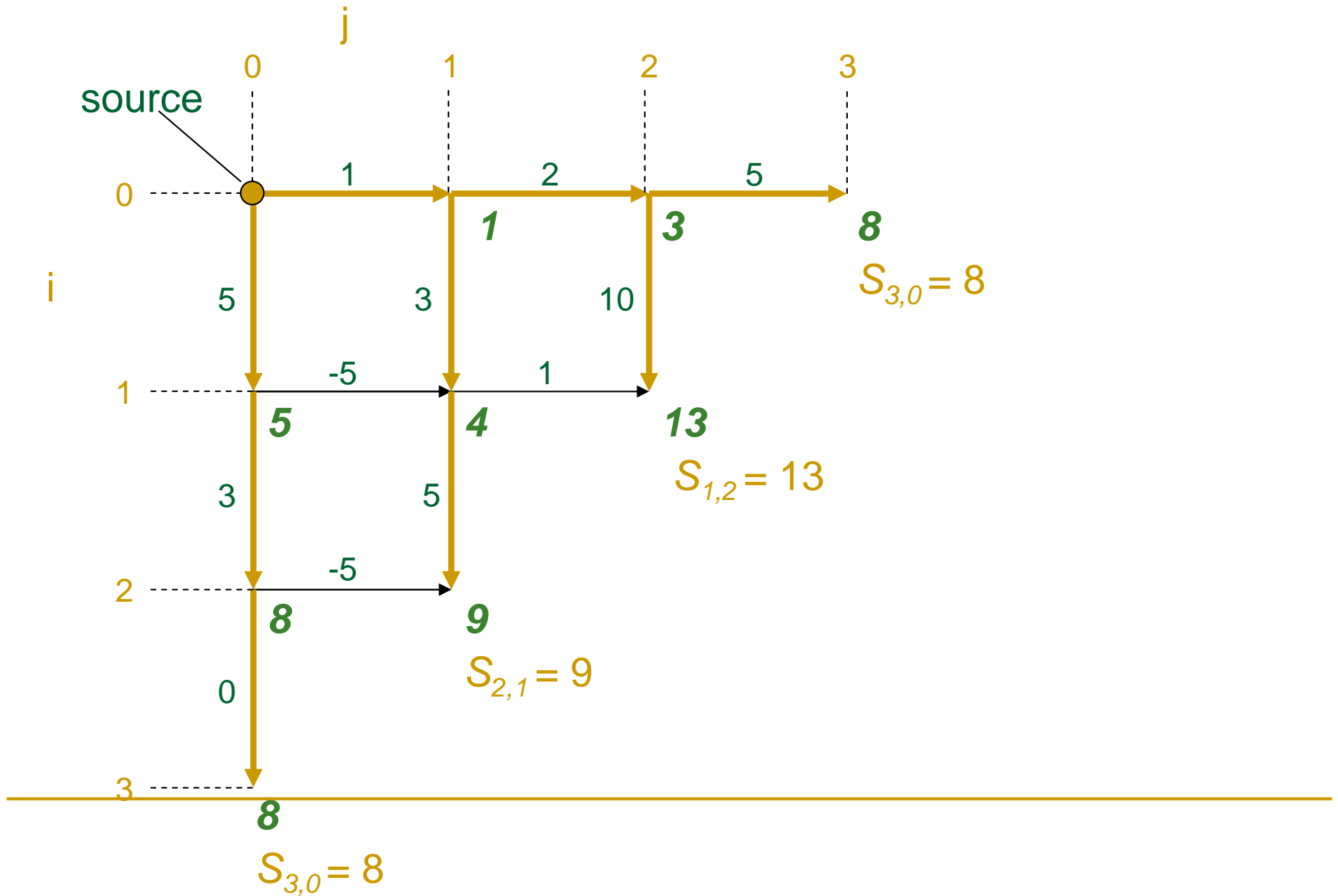


- Calculate optimal path score for each vertex in the graph
- Each vertex's score is the maximum of the prior vertices score plus the weight of the respective edge in between

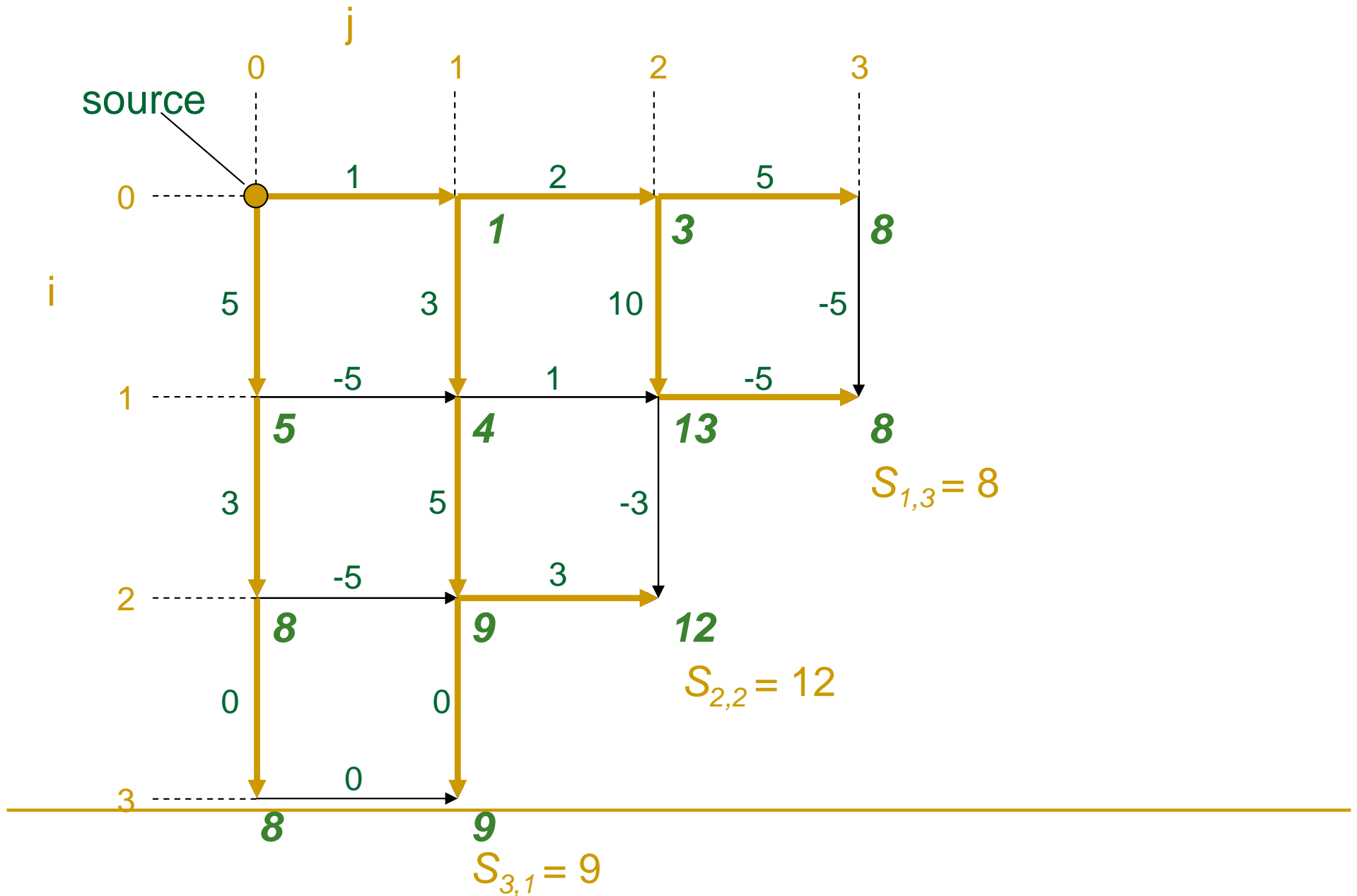
# MTP: Dynamic Programming (cont'd)



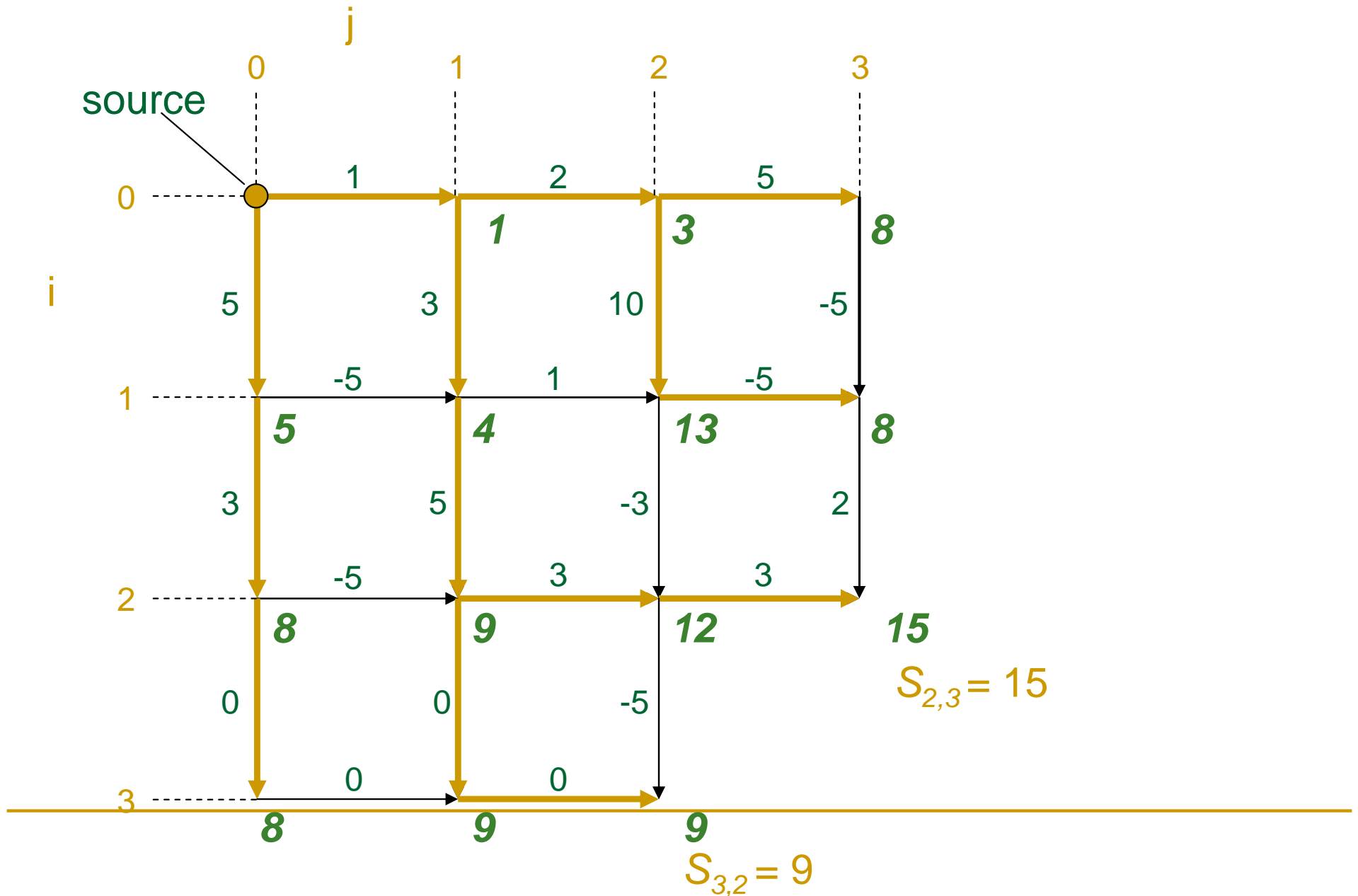
# MTP: Dynamic Programming (cont'd)



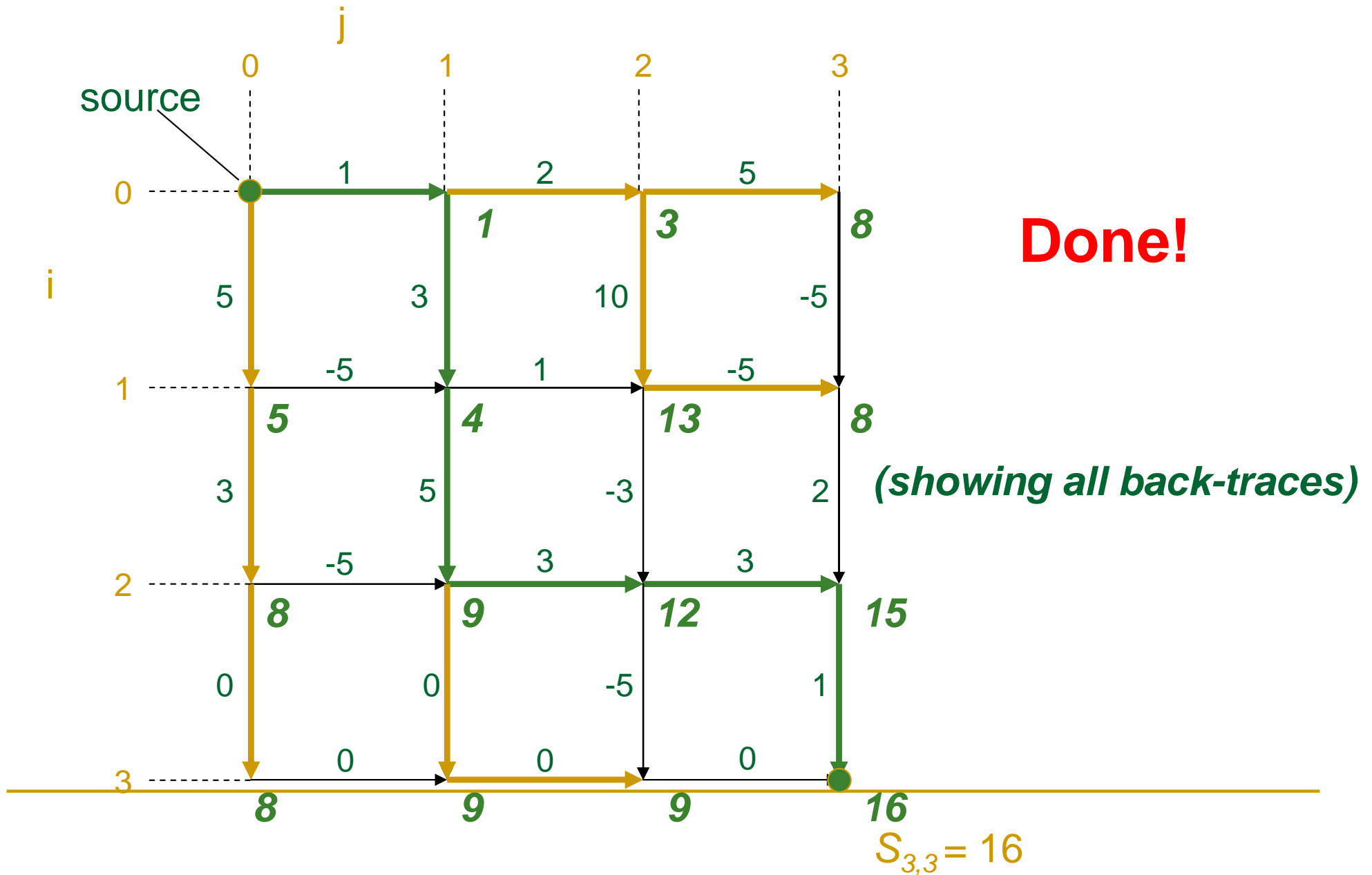
# MTP: Dynamic Programming (cont'd)



# MTP: Dynamic Programming (cont'd)



# MTP: Dynamic Programming (cont'd)



---

# MTP: Recurrence

Computing the score for a point  $(i,j)$  by the recurrence relation:

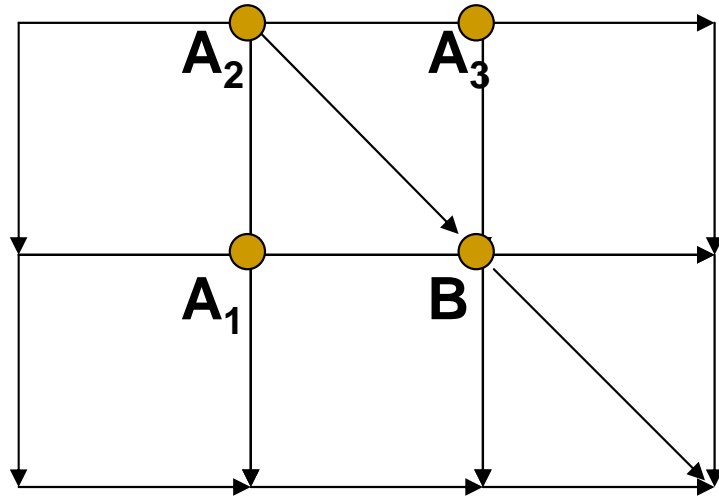
$$s_{i,j} = \max \left\{ \begin{array}{l} s_{i-1,j} + \text{weight of the edge between } (i-1, j) \text{ and } (i, j) \\ s_{i,j-1} + \text{weight of the edge between } (i, j-1) \text{ and } (i, j) \end{array} \right.$$

The running time is  $n \times m$  for a  $n$  by  $m$  grid

( $n = \#$  of rows,  $m = \#$  of columns)

---

# Manhattan Is Not A Perfect Grid



What about diagonals?

- The score at point B is then given by:

$$s_B = \max \text{ of } \begin{cases} s_{A_1} + \text{weight of the edge } (A_1, B) \\ s_{A_2} + \text{weight of the edge } (A_2, B) \\ s_{A_3} + \text{weight of the edge } (A_3, B) \end{cases}$$

## Manhattan Is Not A Perfect Grid (cont'd)

Computing the score for point  $x$  is given by the recurrence relation:

$$s_x = \max_{\text{of}} \left\{ s_y + \text{weight of vertex } (y, x) \text{ where } y \in \text{Predecessors}(x) \right.$$

- Predecessors ( $x$ ) = set of vertices that have edges leading to  $x$
- The running time for a graph  $G(\mathbf{V}, \mathbf{E})$  ( $\mathbf{V}$  is the set of all vertices and  $\mathbf{E}$  is the set of all edges) is  $O(\mathbf{E})$  since each edge is evaluated once

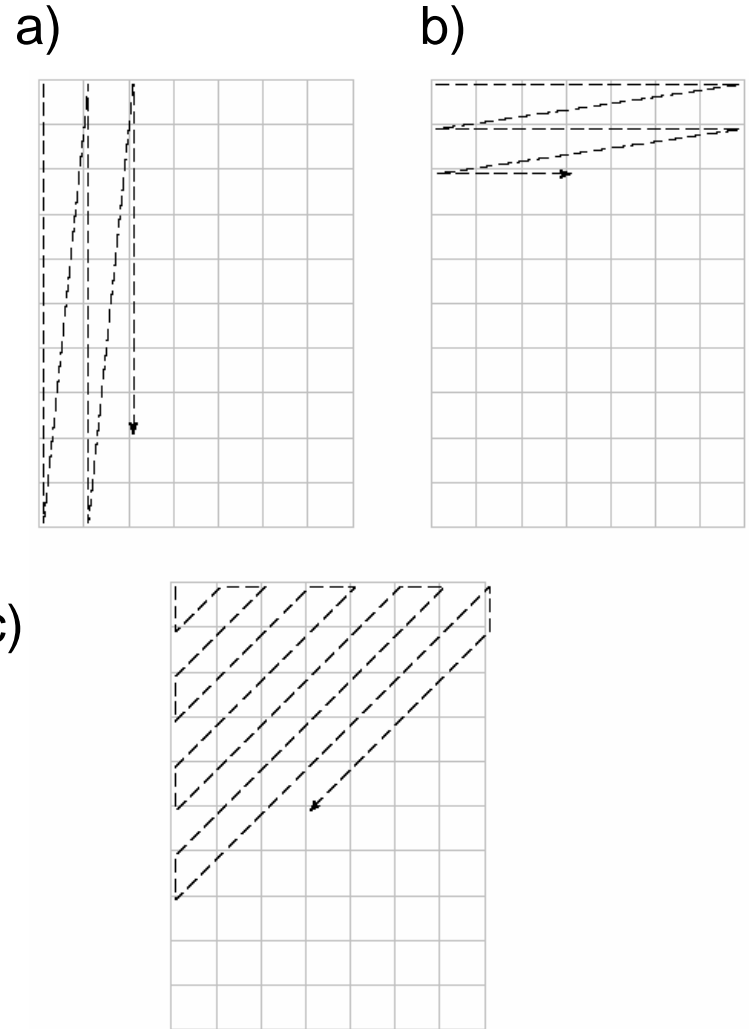
---

## Traveling in the Grid

- The only hitch is that one must decide on the order in which to visit the vertices
  - By the time the vertex  $x$  is analyzed, the values  $s_y$  for all its predecessors  $y$  should be computed – otherwise we are in trouble.
  - We need to traverse the vertices in some order
-

# Traversing the Manhattan Grid

- 3 different strategies:
  - a) Column by column
  - b) Row by row
  - c) Along diagonals



# Alignment: 2 row representation

Given 2 DNA sequences  $v$  and  $w$ :

$v$  : **A** **T** **C** **T** **G** **A** **T**      $m = 7$   
 $w$  : **T** **G** **C** **A** **T** **A**      $n = 6$

Alignment :  $2 * k$  matrix (  $k \geq \max(m, n)$  )

|                |   |   |    |   |   |    |   |    |    |
|----------------|---|---|----|---|---|----|---|----|----|
| letters of $v$ | A | T | -- | G | T | T  | A | T  | -- |
| letters of $w$ | A | T | C  | G | T | -- | A | -- | C  |

4 matches

2 insertions

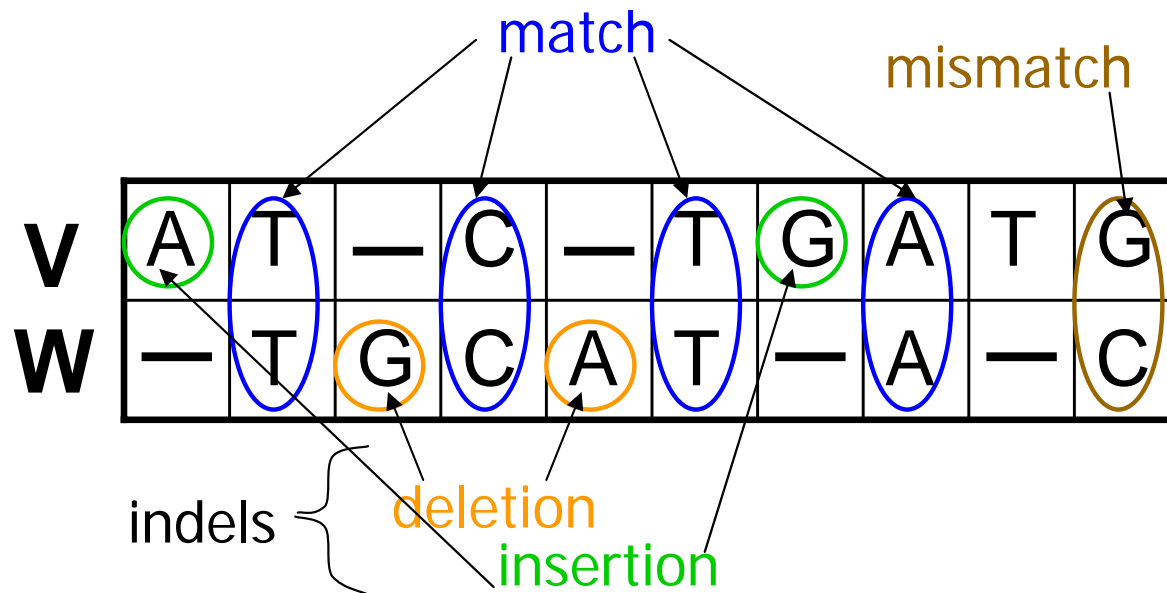
2 deletions

# Aligning DNA Sequences

**V** = ATCTGATG       $n = 8$

**W** = TGCATAC       $m = 7$

4 matches  
1 mismatches  
2 insertions  
2 deletions



Note:  
insertions and  
deletions are  
together  
called **indels**

# Longest Common Subsequence (LCS) – Alignment without Mismatches

- Given two sequences

$$\mathbf{v} = v_1 v_2 \dots v_m \text{ and } \mathbf{w} = w_1 w_2 \dots w_n$$

- The LCS of  $\mathbf{v}$  and  $\mathbf{w}$  is a sequence of positions in

$$\mathbf{v}: 1 \leq i_1 < i_2 < \dots < i_t \leq m$$

and a sequence of positions in

$$\mathbf{w}: 1 \leq j_1 < j_2 < \dots < j_t \leq n$$

such that  $i_t$ -th letter of  $\mathbf{v}$  equals to  $j_t$ -letter of  $\mathbf{w}$  and  $t$  is maximal

# LCS: Example

|                      |    |   |    |   |    |   |    |   |    |   |   |
|----------------------|----|---|----|---|----|---|----|---|----|---|---|
| <i>i</i> coords:     | 0  | 1 | 2  | 2 | 3  | 3 | 4  | 5 | 6  | 7 | 8 |
| elements of <i>v</i> | A  | T | -- | C | -- | T | G  | A | T  | C |   |
| elements of <i>w</i> | -- | T | G  | C | A  | T | -- | A | -- | C |   |
| <i>j</i> coords:     | 0  | 0 | 1  | 2 | 3  | 4 | 5  | 5 | 6  | 6 | 7 |

$(0,0) \rightarrow (1,0) \rightarrow (2,1) \rightarrow (2,2) \rightarrow (3,3) \rightarrow (3,4) \rightarrow (4,5) \rightarrow (5,5) \rightarrow (6,6) \rightarrow (7,6) \rightarrow (8,7)$

Matches shown in red      positions in *v*:  $2 < 3 < 4 < 6 < 8$   
                                         positions in *w*:  $1 < 3 < 5 < 6 < 7$

Every common subsequence is a path in 2-D grid

# LCS: Dynamic Programming

- Find the LCS of two strings

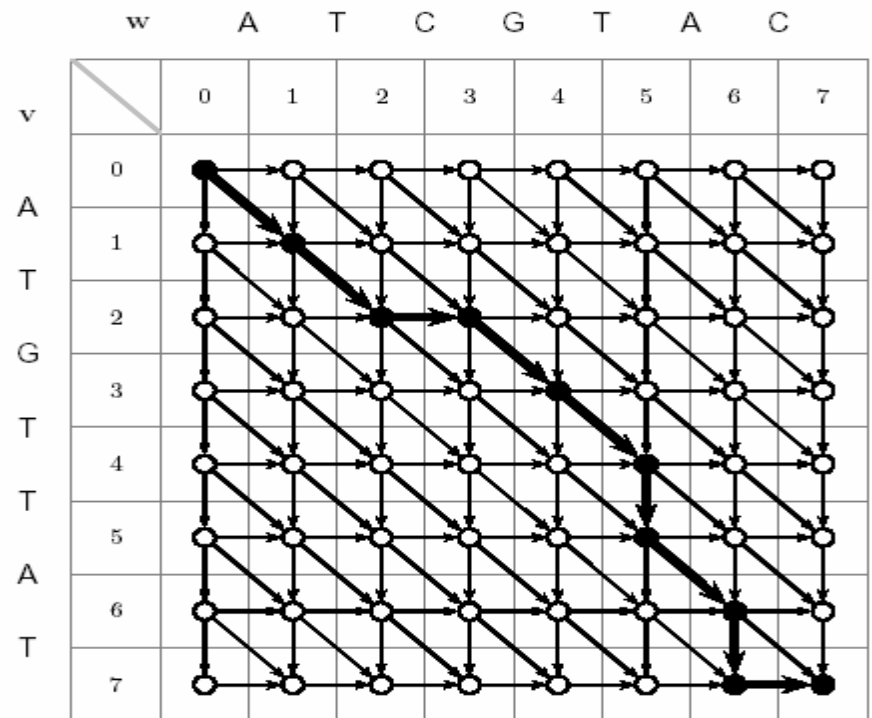
Input: A weighted graph  $G$  with two distinct vertices, one labeled “*source*” one labeled “*sink*”

Output: A longest path in  $G$  from “*source*” to “*sink*”

```

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

```

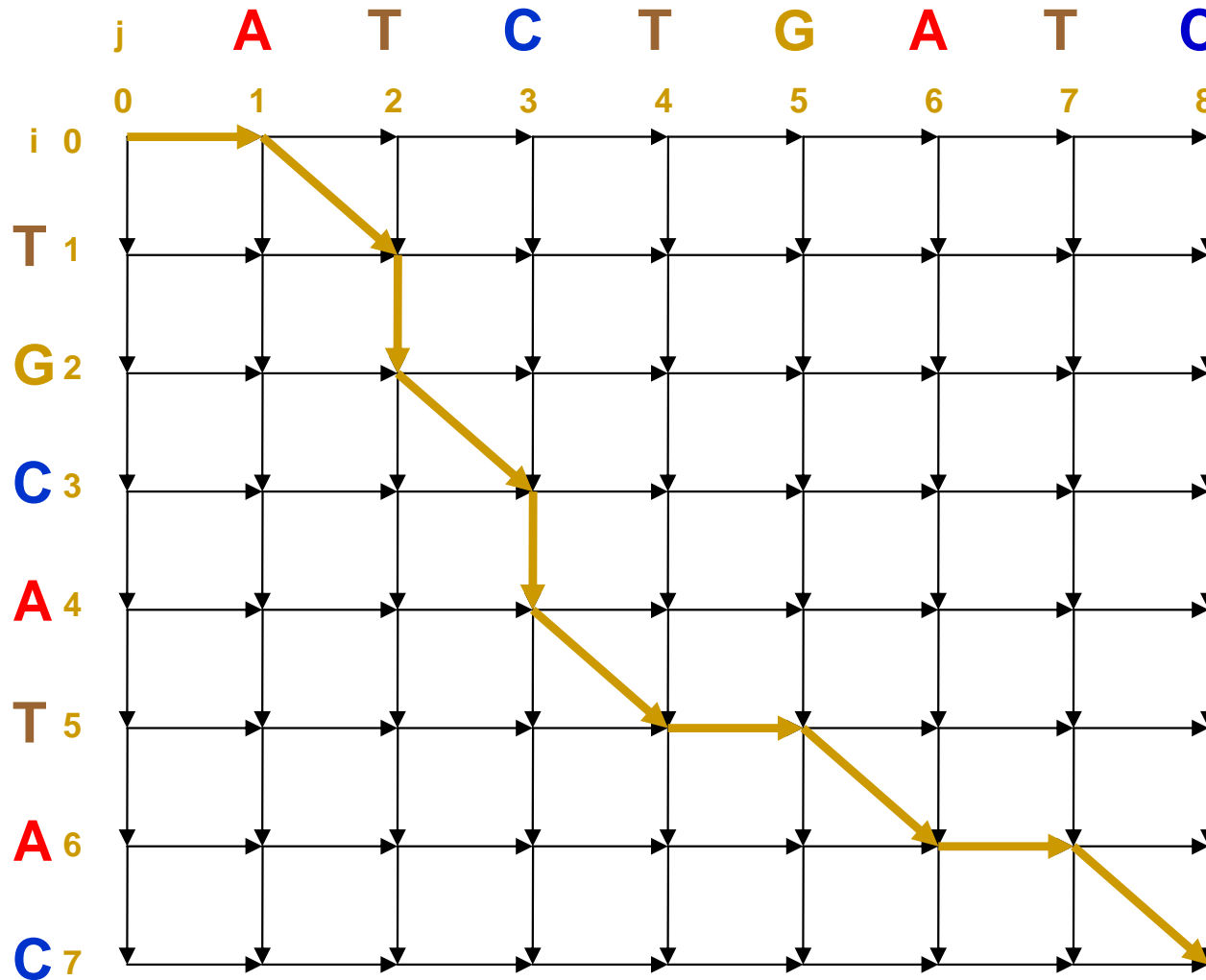


```

↘ ↘ → ↘ ↘ ↓ ↘ ↓ →
A T - G T T A T -
A T C G T - A - C

```

# LCS Problem as Manhattan Tourist Problem



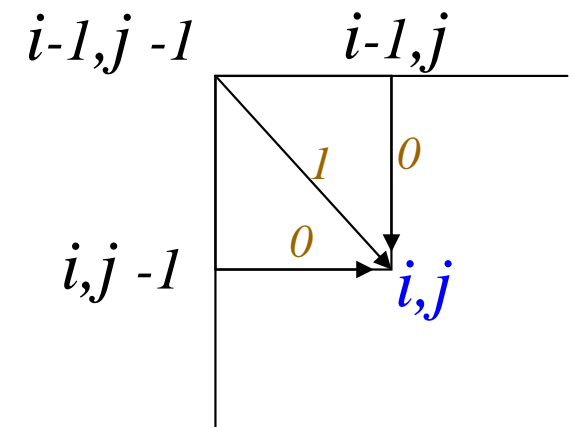
# Computing LCS

Let  $\mathbf{v}_i =$  prefix of  $\mathbf{v}$  of length  $i$ :  $v_1 \dots v_i$

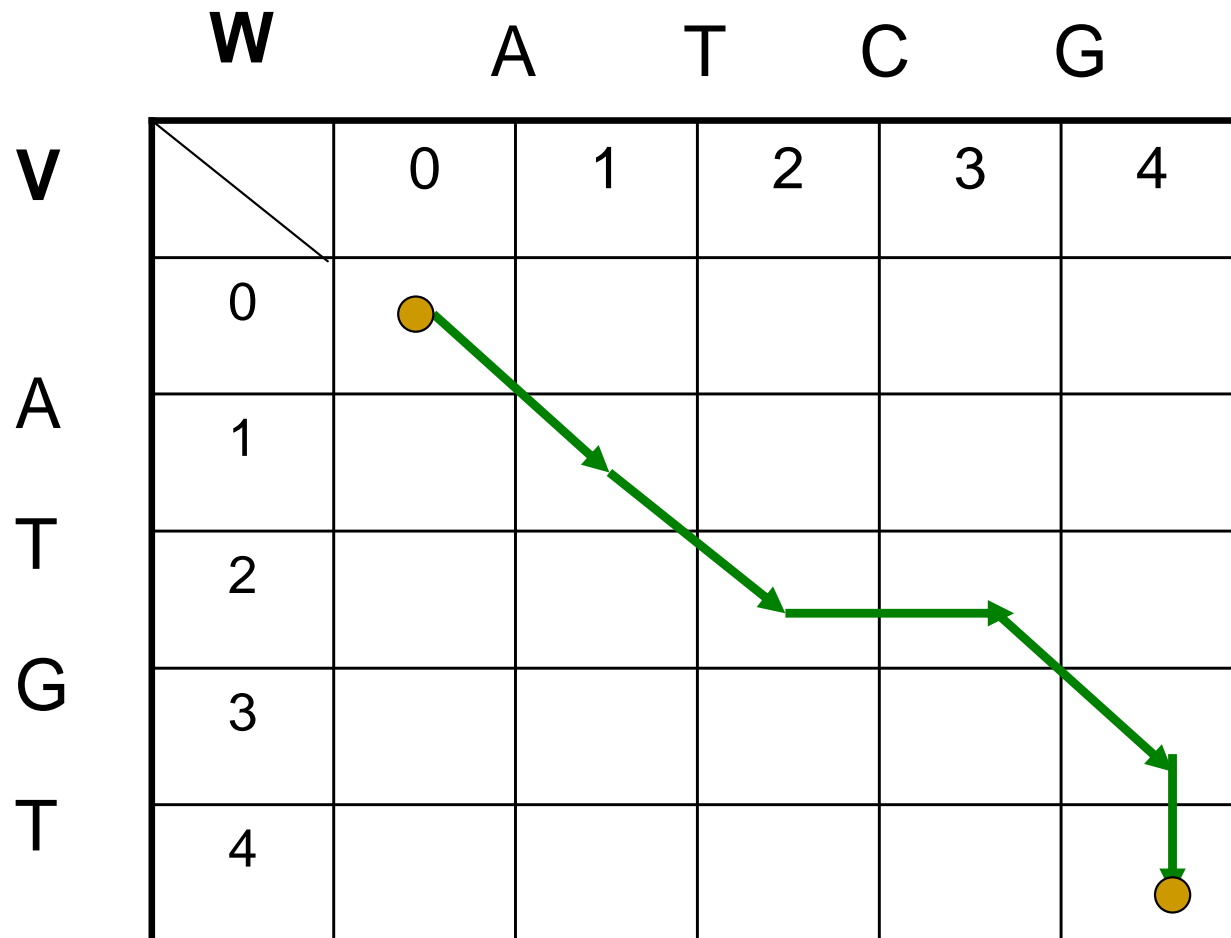
and  $\mathbf{w}_j =$  prefix of  $\mathbf{w}$  of length  $j$ :  $w_1 \dots w_j$

The length of  $\text{LCS}(\mathbf{v}_i, \mathbf{w}_j)$  is computed by:

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



# Every Path in the Grid Corresponds to an Alignment



$\swarrow \swarrow \rightarrow \swarrow \downarrow$   
 0 1 2 2 3 4  
 V = A T - G T  
 | | |  
 W = A T C G -  
 0 1 2 3 4 4

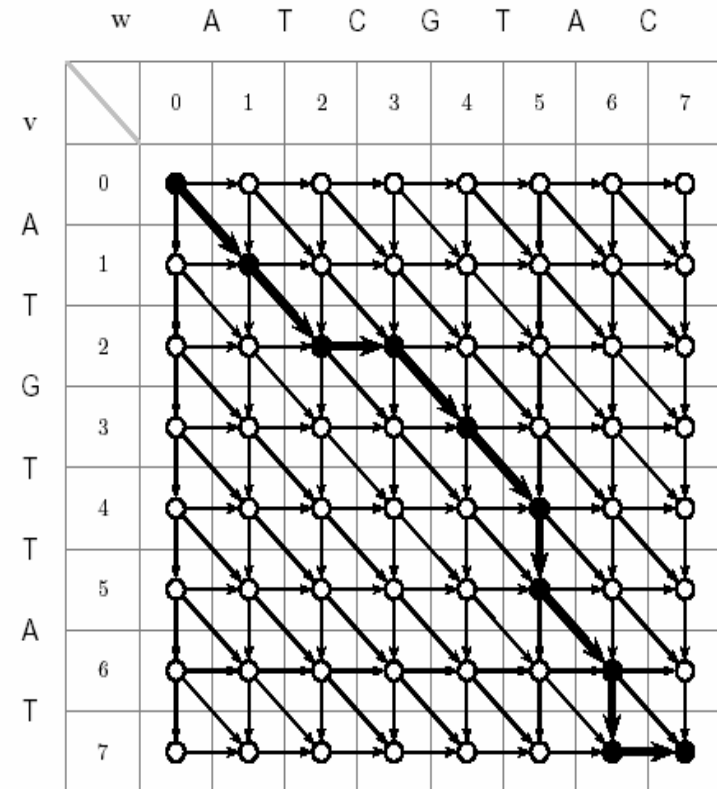
# The Alignment Grid

```

v = 0 1 2 2 3 4 5 6 7 7
 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

```

- Every alignment path is from source to sink

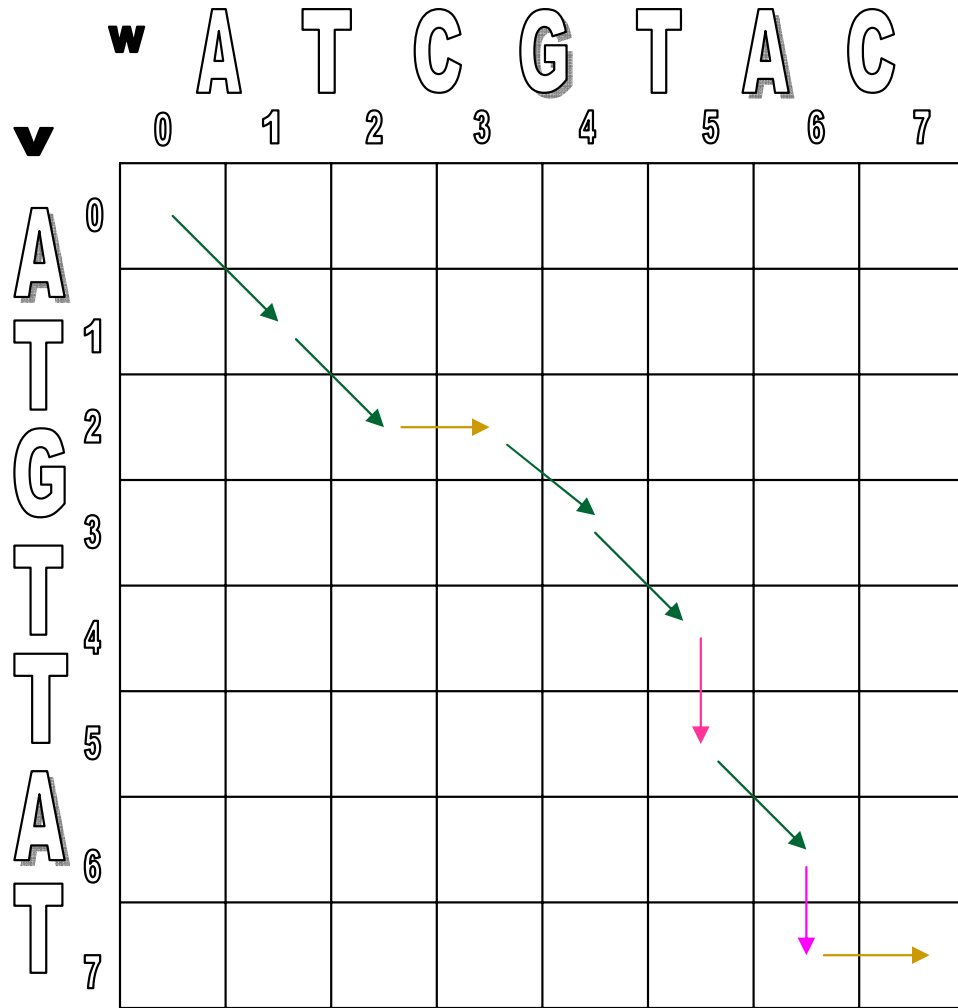


```

 \ \ → \ \ ↓ \ ↓ →
 A T - G T T A T -
 A T C G T - A - C

```

# Alignments in Edit Graph (cont'd)



↓ and → represent indels in **v** and **w** with score 0.

↘ represent matches with score 1.

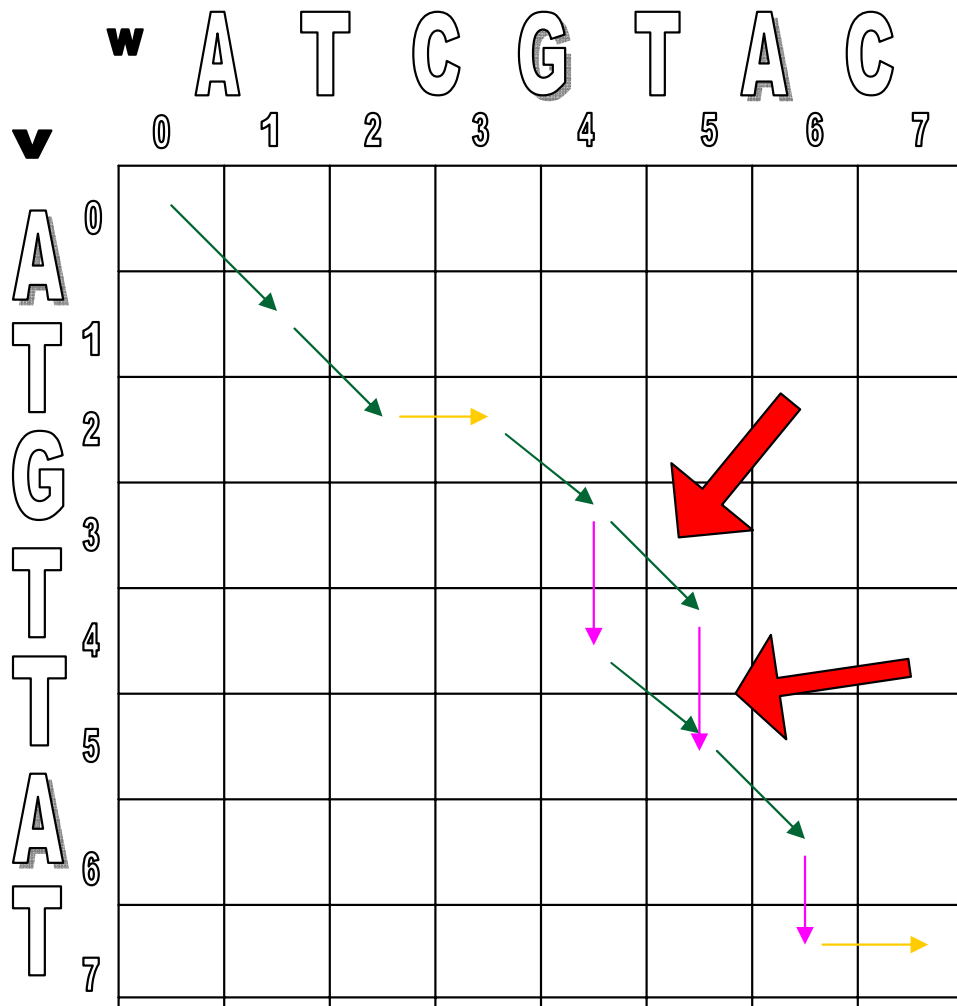
- The score of the alignment path is 5.

Every path in the edit graph corresponds to an alignment:



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

# Alignment as a Path in the Edit Graph



## Old Alignment

0122345677  
 v= AT\_GTTAT\_  
 w= ATCGT\_A\_C  
 0123455667

## New Alignment

0122345677  
 v= AT\_GTTAT\_  
 w= ATCG\_TA\_C  
 0123445667

# Dynamic Programming Example

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


Initialize  $1^{st}$  row and  $1^{st}$  column to be all zeroes.

Or, to be more precise, initialize  $0^{th}$  row and  $0^{th}$  column to be all zeroes.

# Dynamic Programming Example

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

$$S_{i,j} = \max \begin{cases} S_{i-1,j-1} & \leftarrow \text{value from NW} + 1, \text{ if } v_i = w_j \\ S_{i-1,j} & \leftarrow \text{value from North (top)} \\ S_{i,j-1} & \leftarrow \text{value from West (left)} \end{cases}$$

Arrows  show where the score originated from.

 if from the top

 if from the left

 if  $v_i = w_j$

# Backtracking Example

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

Find a match in row and column 2.

$i=2, j=2,5$  is a match (T).

$j=2, i=4,5,7$  is a match (T).

Since  $v_i = w_j$ ,  $s_{i,j} = s_{i-1,j-1} + 1$

$$s_{2,2} = [s_{1,1} = 1] + 1$$

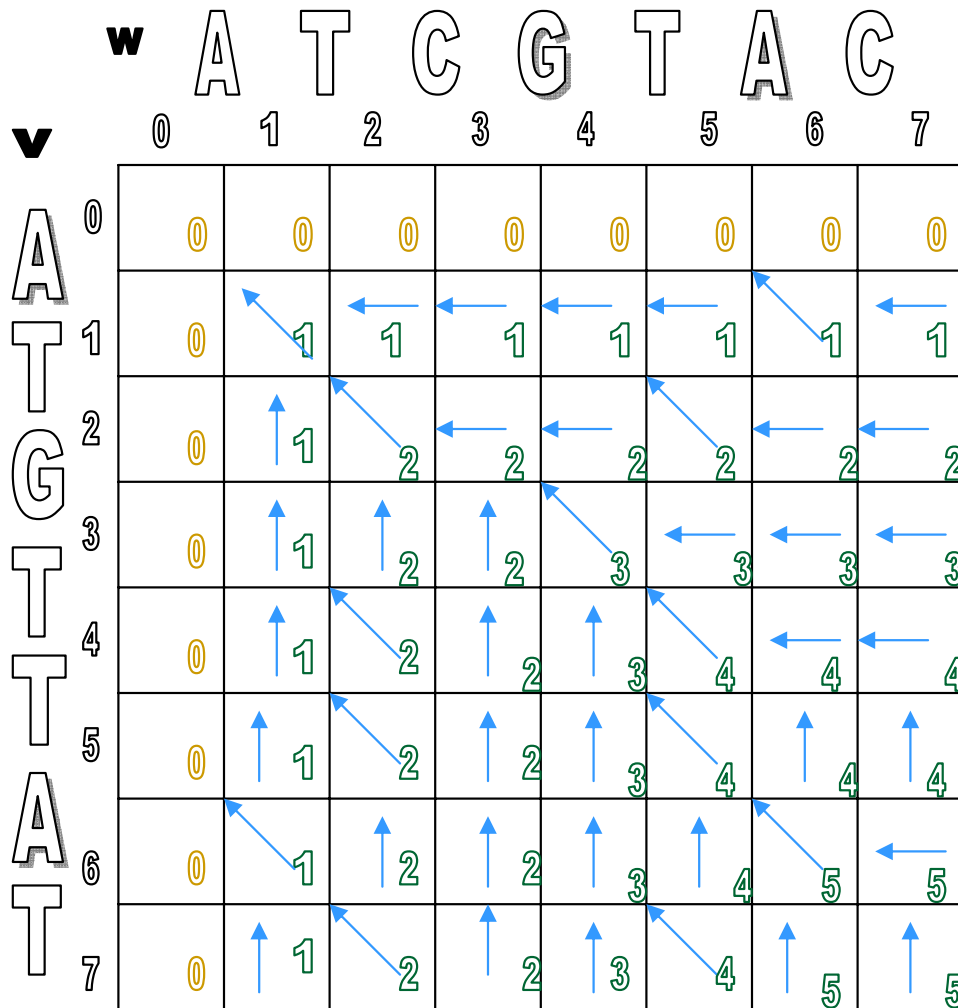
$$s_{2,5} = [s_{1,4} = 1] + 1$$

$$s_{4,2} = [s_{3,1} = 1] + 1$$

$$s_{5,2} = [s_{4,1} = 1] + 1$$

$$s_{7,2} = [s_{6,1} = 1] + 1$$

# Backtracking Example



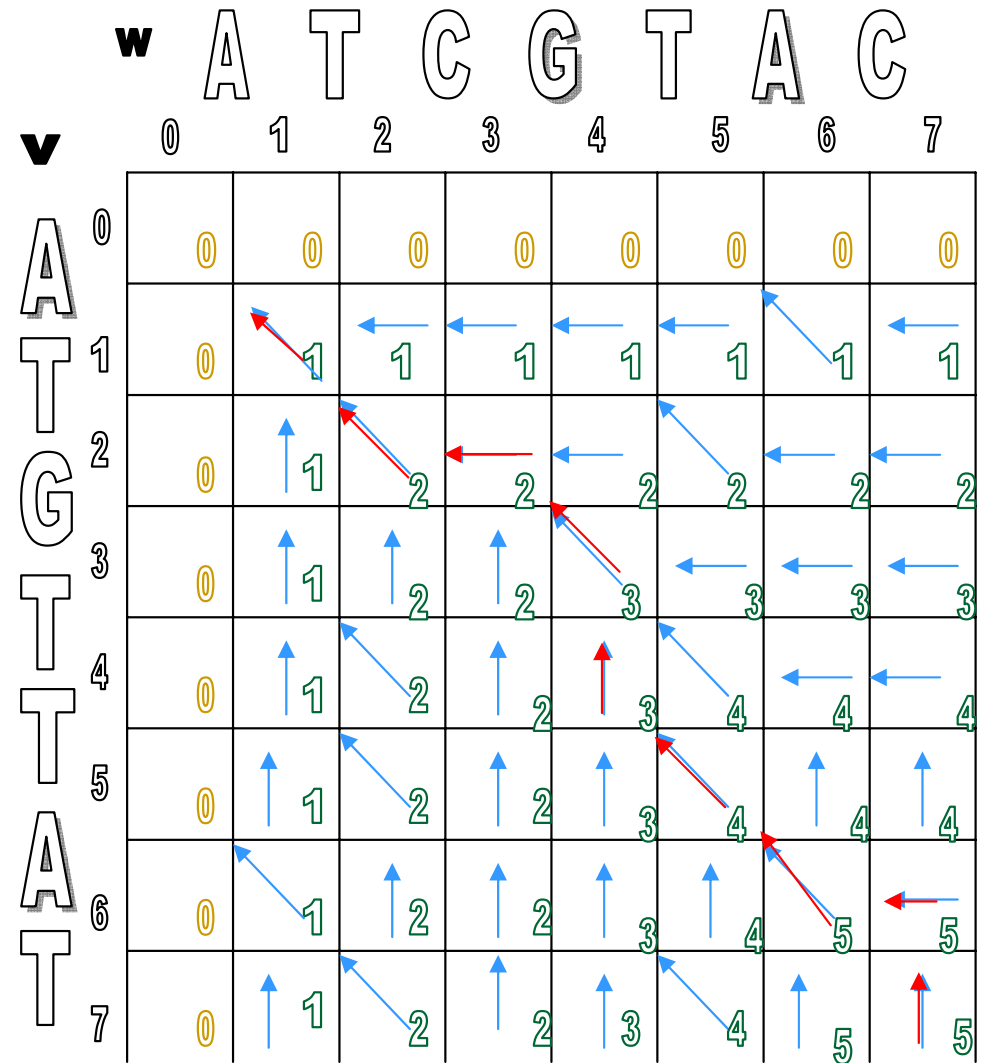
Continuing with the dynamic programming algorithm gives this result.

# LCS Algorithm

```
1. LCS(v, w)
2. for $i = 1$ to n
3. $s_{i,0} = 0$
4. for $j = 1$ to m
5. $s_{0,j} = 0$
6. for $i = 1$ to n
7. for $j = 1$ to m
8. $s_{i,j} = \max \begin{cases} s_{i-1,j} \\ s_{i,j-1} \\ s_{i-1,j-1} + 1, \text{ if } v_i = w_j \end{cases}$
9. $b_{i,j} = \begin{cases} \uparrow & \text{if } s_{i,j} = s_{i-1,j} \\ \leftarrow & \text{if } s_{i,j} = s_{i,j-1} \\ \swarrow & \text{if } s_{i,j} = s_{i-1,j-1} + 1 \end{cases}$
10.
11.
■ return $(s_{n,m}, b)$
```

# Now What?

- $LCS(v,w)$  created the alignment grid
- Now we need a way to read the best alignment of  $v$  and  $w$
- Follow the arrows backwards from sink



# Printing LCS: Backtracking

1. **PrintLCS(b,v,i,j)**
2.     **if**  $i = 0$  or  $j = 0$
3.         **return**
4.     **if**  $b_{i,j} = \nwarrow$
5.         **PrintLCS(b,v,i-1,j-1)**
6.         **print**  $v_i$
7.     **else**
8.         **if**  $b_{i,j} = \uparrow$
9.             **PrintLCS(b,v,i-1,j)**
10.         **else**
11.             **PrintLCS(b,v,i,j-1)**

---

# LCS Runtime

- It takes  $O(nm)$  time to fill in the  $n \times m$  dynamic programming matrix.
  - Why  $O(nm)$ ? The pseudocode consists of a nested “for” loop inside of another “for” loop to set up a  $n \times m$  matrix.
-

---

# Summary

- The running times of algorithms is important!
    - If it doesn't scale up, it won't be useful, especially in bioinformatics
  - Recursion is a basic technique which is useful for breaking down problems into simpler ones
  - Dynamic programming, which uses recursion, is often used in bioinformatics as well
    - Shown to be mathematically accurate
    - However, it can be inefficient for more than two sequences
      - BLAST and FASTA use heuristics (human-like techniques to speed up the computations)
-