

# Edit Distance

# Name: Dynamic Programming

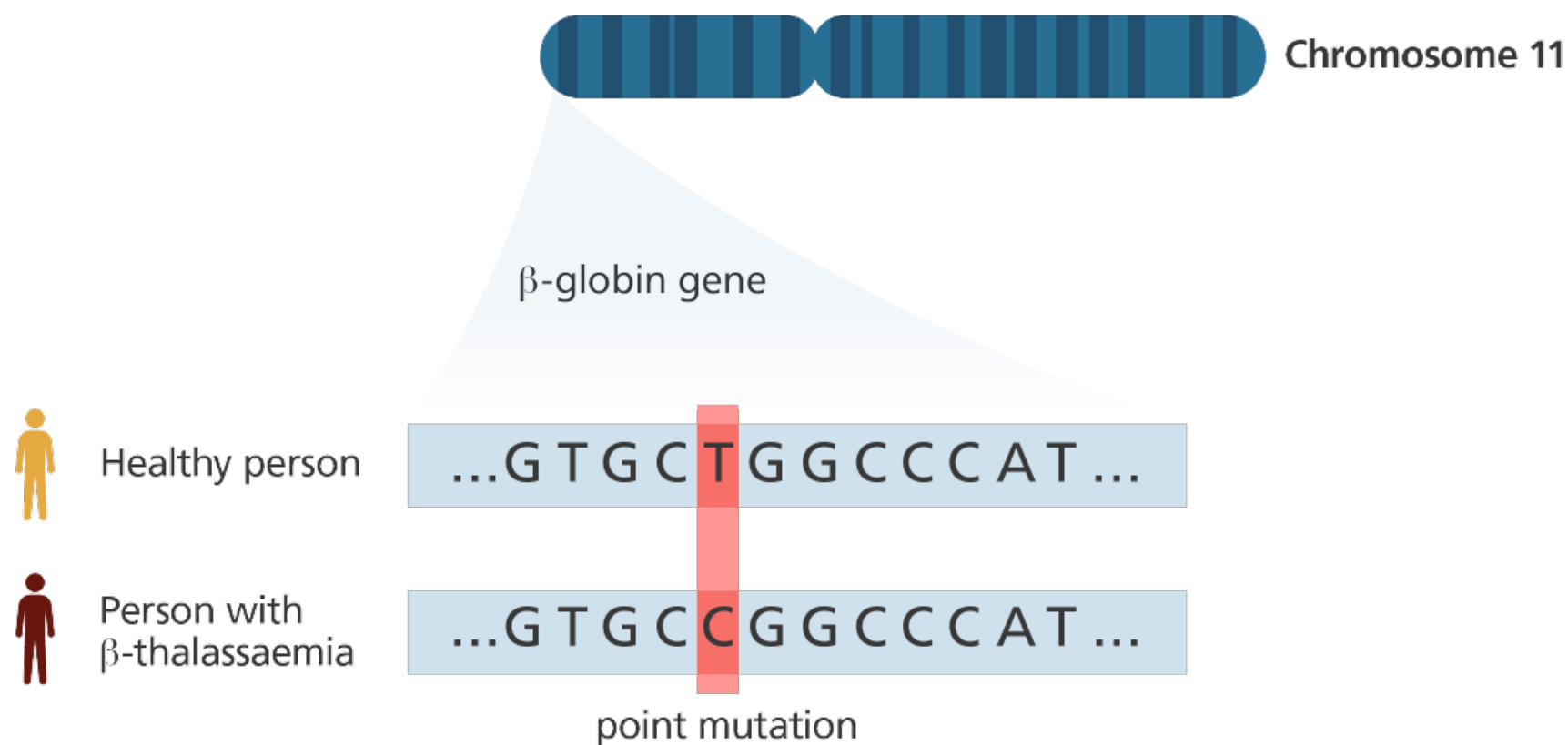
- Formalized by Richard Bellman in the 1950s

We had a very interesting gentleman in Washington named Wilson. He was secretary of Defense, and he actually had a pathological fear and hatred of the word “*research*”. I’m not using the term lightly; I’m using it precisely. His face would suffuse, he would turn red, and he would get violent if people used the term “*research*” in his presence. You can imagine how he felt, then, about the term “*mathematical*”.... I felt I had to do something to shield Wilson and the Air Force from the fact that I was really doing mathematics inside the RAND Corporation. What title, what name, could I choose?

- Chose the name “**dynamic programming**” to hide the mathematical nature of the work from military bosses

# Motivation

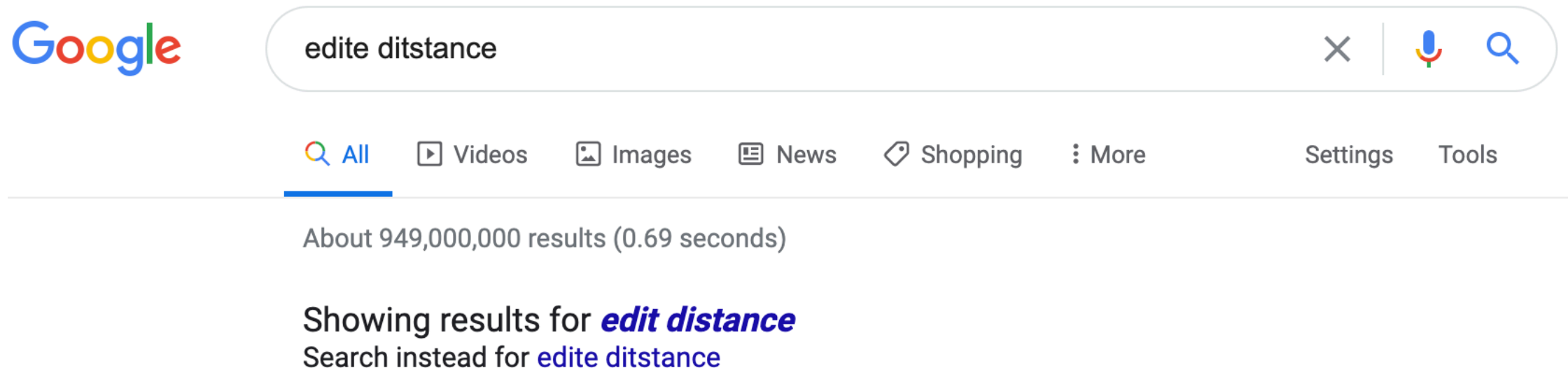
- **Edit distance:** is a metric that captures the similarity between two strings



DNA sequencing: finding similarities between two genome sequences

# Motivation

- **Edit distance:** is a metric that captures the similarity between two strings



# Problem Definition

- **Problem.** Given two strings  $A = a_1 \cdot a_2 \cdot \dots \cdot a_n$  and  $B = b_1 \cdot b_2 \cdot \dots \cdot b_m$  find the edit distance between them
- Edit distance between  $A$  and  $B$  is the smallest number of the following operations that are needed to transform  $A$  into  $B$ 
  - Replace a character (substitution)
  - Delete a character
  - Insert a character

riddle  $\xrightarrow{\text{delete}}$  ridle  $\xrightarrow{\text{substitute}}$  riple  $\xrightarrow{\text{insert}}$  triple

Edit distance(riddle, triple): 3

# Edit Distance

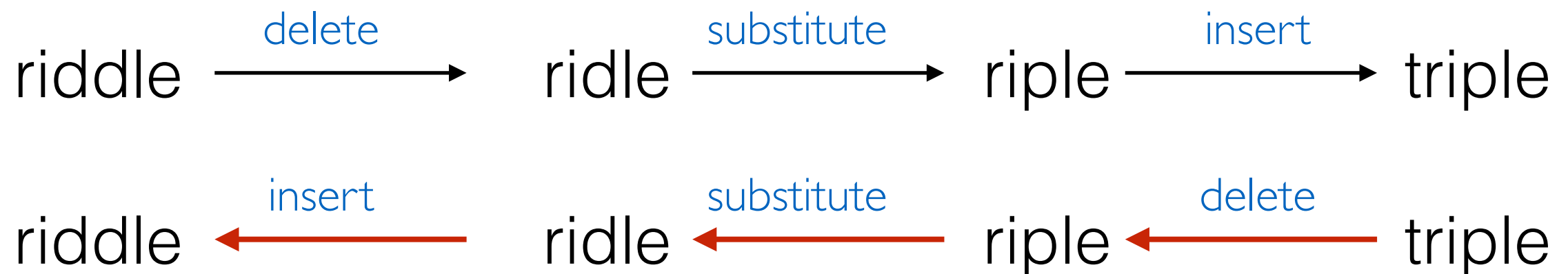
- **Problem.** Given two strings find the minimum number of edits (letter insertions, deletions and substitutions) that transform one string into the other
- Measure of similarity between strings
- For example, the edit distance between FOOD and MONEY is at most four:

FOOD → MOOD → MON<sup>^</sup>D → MONED → MONEY

- Can observe and conclude that 3 edits don't work
- Edit distance = 4 in this case

# Structure of the Problem

- **Problem.** Given two strings  $A = a_1 \cdot a_2 \cdot \dots \cdot a_n$  and  $B = b_1 \cdot b_2 \cdot \dots \cdot b_m$  find the edit distance between them
- Notice that the process of getting from string  $A$  to string  $B$  by doing substitutions, inserts and deletes is reversible
- **Inserts** in one string correspond to **deletes** in another



Edit distance(riddle, triple): 3

# Sequence Alignment

- We can visualize the problem of finding the edit distance as an the problem of finding the best alignment between two strings
- **Gaps** in alignment represent inserts/deletes
- **Mismatches** in alignment represent substitutes
- Cost of an alignment = **number of gaps + mismatches**
- Edit distance: minimum cost alignment

r	i	d	d	l	e	
	t	r	i	p	l	e

cost = 7

	r	i	d	d	l	e
t	r	i		p	l	e

cost = 3



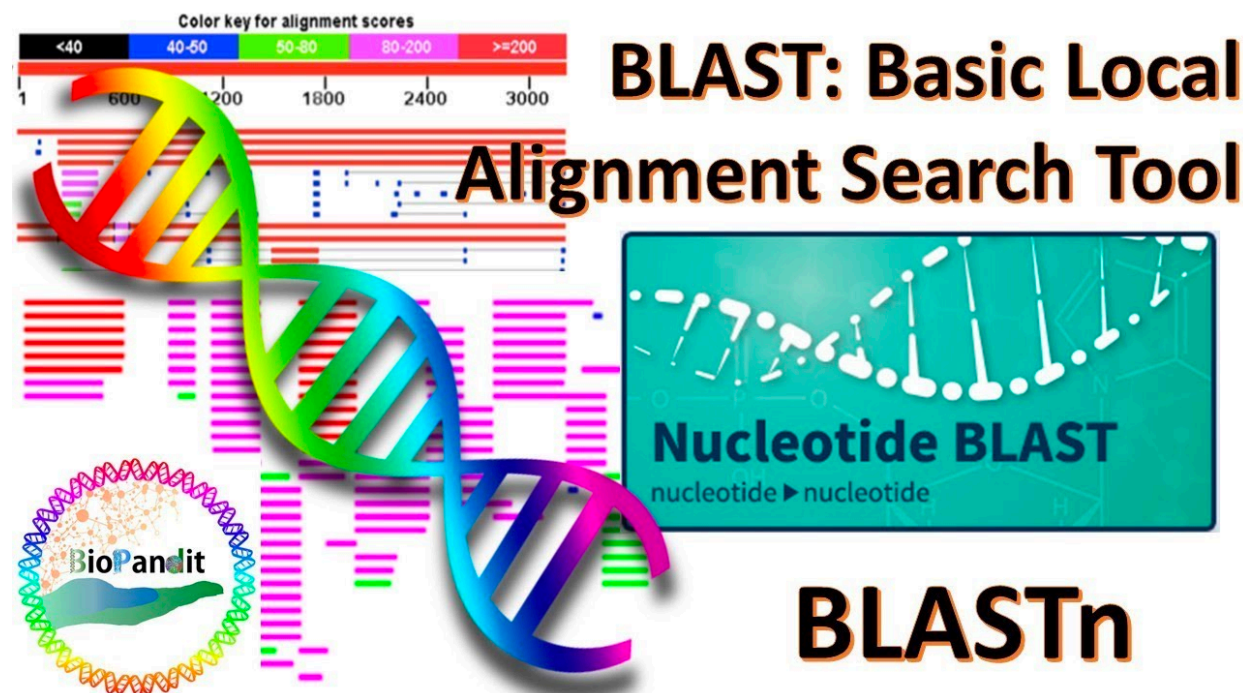
# Sequence Alignment

>gb|AC115706.7| Mus musculus chromosome 8, clone RP23-382B3, complete sequence

```
Query 1650 gtgtgtgtgggtgcacatttgtgtgtgtgtgcgcctgtgtgtgtgggtgcctgtgtgtgt 1709
          ||||| |  || | ||||| | ||||| ||| || |||||
Sbjct 56838 GTGTGTGTGGAAGTGAGTTCATCTGTGTGTGCACATGTGTGTGCA--TGCATGCATGTGT 56895

Query 1710 gtg-gggcacatttgtgtgtgtgtgtgtgcctgtgtgtgggtgcacatttgtgtgtgtgc 1768
          || ||||| || ||| ||||| ||||| ||| ||| ||||| |||
Sbjct 56896 GTCCGGGCA-----TGCATGTCTGTGTGCATGTGTGTGTGTGTGCAT--GTGTGAGTAC 56947

Query 1769 ctgtgtgtgtgtgcctgtgtgtgggggtgcacatttgtgtgtgtgtgtgcctgtgtgtgg 1828
          ||||| ||| ||| |||| | ||| ||| ||||| ||||| ||||| |
Sbjct 56948 CTGTGTGTGTATGCTTGTATGTGTGTGTGTGCATGTGTGTAGGTGTGTATATGTGTAAGT 57007
```



# Sequence Alignment

prin-ciple  
|||| |||xx  
prinncipal  
(1 gap, 2 mm)

prin-cip-le  
|||| ||| |  
prinncipal-  
(3 gaps, 0 mm)

misspell  
||| ||||  
mis-pell  
(1 gap)

prehistoric  
||| |||||  
---historic  
(3 gaps)

aa-bb-ccaabb  
|x || | | |  
ababbbbc-a-b-  
(5 gaps, 1 mm)

al-go-rithm-  
|| xx ||x |  
alKhwariz-mi  
(4 gaps, 3 mm)

# Sequence Alignment

- These alignments are a way of **visualizing** the edit distance between two strings
- For every sequence of edits between two strings, we can **draw** a sequence alignment

riddle  $\xrightarrow{\text{delete}}$  ridle  $\xrightarrow{\text{substitute}}$  riple  $\xrightarrow{\text{insert}}$  triple

	r	i	d	d	l	e
t	r	i		p	l	e

# Sequence Alignment Problem

- Find an alignment of the two strings  $A, B$  where
  - each character  $a_i$  in  $A$  is matched to a string  $b_j$  in  $B$  or unmatched
  - each character  $b_j$  in  $B$  is matched to a string  $a_i$  in  $A$  or unmatched
- $\text{cost}(a_i, b_j) = 0$  if  $a_i = b_j$ , else  $\text{cost}(a_i, b_j) = 1$
- cost of an unmatched letter (gap) = 1
- Total cost = # unmatched (gaps) +  $\sum_{a_i, b_j} \text{cost}(a_i, b_j)$
- **Goal.** Compute edit distance by finding an alignment of the minimum total cost

# Sequence Alignment

- The problem of finding the smallest edit distance *is* the problem of finding the best alignment
- It's just drawn differently
- Any questions about this?

riddle  $\xrightarrow{\text{delete}}$  ridle  $\xrightarrow{\text{substitute}}$  riple  $\xrightarrow{\text{insert}}$  triple

	r	i	d	d	l	e
t	r	i		p	l	e

# Recipe for a Dynamic Program

- **Formulate the right subproblem.** The subproblem must have an optimal substructure
- **Formulate the recurrence.** Identify how the result of the smaller subproblems can lead to that of a larger subproblem
- **State the base case(s).** The subproblem that's so small we know the answer to it!
- **State the final answer.** (In terms of the subproblem)
- **Choose a memoization data structure.** Where are you going to store already computed results? (Usually a table)
- **Identify evaluation order.** Identify the dependencies: which subproblems depend on which ones? Using these dependencies, identify an evaluation order
- **Analyze space and running time.** As always!



# How to Come Up with a DP

Ask yourself **two questions**:

- **What subproblem** should I use?
- How can I **recursively find the solution** to a subproblem (using the solution to smaller subproblems)?
- **Splitting into “cases”** is a common way to answer these questions



# Recursive Structure

- **Imagine** for a second that we have the mismatch/gap representation of the shortest edit sequence of two strings
- How can I reduce this to a smaller subproblem?
- If we remove the last column, the remaining columns must represent the shortest edit sequence of the remaining prefixes!

A	L	G	O	R		I		T	H	M	
A	L			T	R	U	I	S	T	I	C



# Recursive Structure

- Suppose we have the mismatch/gap representation of the shortest edit sequence of two strings
- How can I reduce this to a smaller subproblem?
- If we remove the last column, the remaining columns must represent the shortest edit sequence of the remaining prefixes!

$$\begin{array}{cccccc} \text{(cost 2)} & & & \text{(cost 2)} & & \text{(cost 0)} \\ \begin{array}{cccccc} r & i & d & d & l & e \\ | & | & | & | & | & | \\ t & r & i & & p & l & e \end{array} & = & \begin{array}{cccccc} r & i & d & d & l & \\ | & | & | & | & | & \\ t & r & i & & p & l \end{array} & + & \begin{array}{cc} e \\ | \\ e \end{array} \end{array}$$

# Recursive Structure

- Suppose we have the mismatch/gap representation of the shortest edit sequence of two strings
- How can I reduce this to a smaller subproblem?
- If we remove the last column, the remaining columns must represent the shortest edit sequence of the remaining prefixes!

$$\begin{array}{cccccc} \text{(cost 3)} & & & \text{(cost 2)} & & \text{(cost 1)} \\ \begin{array}{cccccc} r & i & d & d & l & t \\ | & | & | & | & | & | \\ t & r & i & & p & l & e \end{array} & = & \begin{array}{cccccc} r & i & d & d & l & \\ | & | & | & | & | & \\ t & r & i & & p & l \end{array} & + & \begin{array}{cccccc} & & & & & t \\ & & & & & | \\ & & & & & e \end{array} \end{array}$$

# Recursive Structure

- Suppose we have the mismatch/gap representation of the shortest edit sequence of two strings
- How can I reduce this to a smaller subproblem?
- If we remove the last column, the remaining columns must represent the shortest edit sequence of the remaining prefixes!

$$\begin{array}{cccccc} \text{(cost 4)} & & & \text{(cost 3)} & & \text{(cost 1)} \\ \begin{array}{cccccc} \mathbf{d} & \mathbf{u} & \mathbf{c} & \mathbf{k} & \mathbf{i} & \mathbf{e} \\ | & | & | & | & | & | \\ \mathbf{t} & \mathbf{r} & \mathbf{u} & \mathbf{c} & \mathbf{k} & \end{array} & = & \begin{array}{ccccc} \mathbf{d} & \mathbf{u} & \mathbf{c} & \mathbf{k} & \mathbf{i} \\ | & | & | & | & | \\ \mathbf{t} & \mathbf{r} & \mathbf{u} & \mathbf{c} & \mathbf{k} \end{array} & + & \begin{array}{c} \mathbf{e} \\ | \end{array} \end{array}$$

# How to Come Up with a DP

Ask yourself **two questions**:

- **What subproblem** should I use?
- How can I **recursively find the solution** to a subproblem (using the solution to smaller subproblems)?



**Recursion**: “cut off” last column, recurse on remaining strings.  
Of course, we don’t know what the last column looks like.  
What are the possibilities? **[on board]**

# Recurrence

- Imagine the optimal alignment between the two strings
- What are the possibilities for the last column?
  - It could be that both letters match: cost 0
  - It could be that both letters do not match: cost 1
  - It could be that there is an unmatched character (gap): either from *A* or from *B*: cost 1

ALGO	R
ALT	R

ALGO	R
ALTR	U

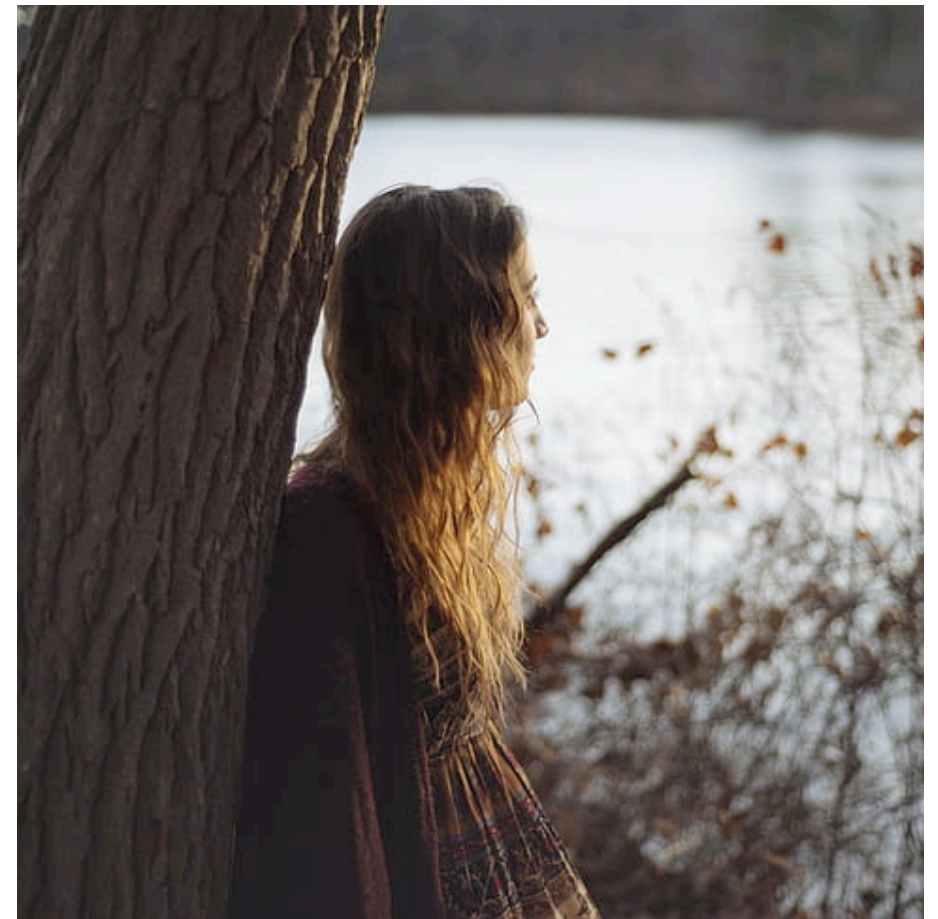
ALGO	R
ALTRU	

ALGOR	
ALTR	U

# How to Come Up with a DP

Ask yourself **two questions**:

- **What subproblem** should I use?
- How can I **recursively find the solution** to a subproblem



Recursion “cuts off” the last column. How should we define the subproblem to allow us to do this?

How does “cutting off” the last column affect the strings?

# Subproblem

- **Subproblem**

Edit( $i, j$ ): edit distance between  
the strings  $a_1 \cdot a_2 \cdot \dots \cdot a_i$  and  $b_1 \cdot b_2 \cdot \dots \cdot b_j$ ,  
where  $0 \leq i \leq n$  and  $0 \leq j \leq m$

- **Final answer**

Edit( $n, m$ )

# Base Cases

- We have to fill out a **two-dimensional array** to memoize our recursive dynamic program
- Let us think about which rows/columns can we fill initially
- $\text{Edit}(i, 0)$ : Min number of edits to transform a string of length  $i$  to an empty string

$$\text{Edit}(i, 0) = i \text{ for } 0 \leq i \leq n$$

$$\text{Edit}(0, j) = j \text{ for } 0 \leq j \leq m$$



# Recurrence

- Three possibilities for the last column in the optimal alignment of  $a_1 \cdot a_2 \cdot \dots \cdot a_i$  and  $b_1 \cdot b_2 \cdot \dots \cdot b_j$ :
- **Case 1.** Only one row has a character:
  - Case 1a. Letter  $a_i$  is unmatched  
 $\text{Edit}(i, j) = \text{Edit}(i - 1, j) + 1$
  - Case 1b. Letter  $b_j$  is unmatched  
 $\text{Edit}(i, j) = \text{Edit}(i, j - 1) + 1$
- **Case 2:** Both rows have characters:
  - Case 2a. Same characters:  
 $\text{Edit}(i, j) = \text{Edit}(i - 1, j - 1)$
  - Case 2b. Different characters:  
 $\text{Edit}(i, j) = \text{Edit}(i - 1, j - 1) + 1$

ALGO	R
ALTRU	

ALGOR	
ALTR	U

ALGO	R
ALT	R

ALGO	R
ALTR	U

# Final Recurrence

- For  $1 \leq i \leq n$  and  $1 \leq j \leq m$ , we have:

$$\text{Edit}(i, j) = \min \begin{cases} \text{Edit}(i, j - 1) + 1 \\ \text{Edit}(i - 1, j) + 1 \\ \text{Edit}(i - 1, j - 1) + (a_i \neq b_j) \end{cases}$$

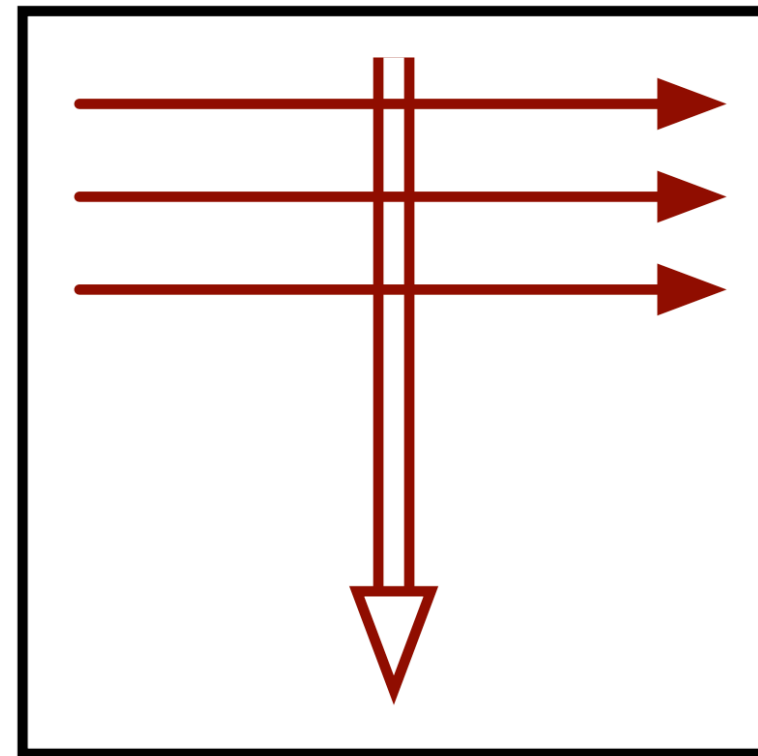
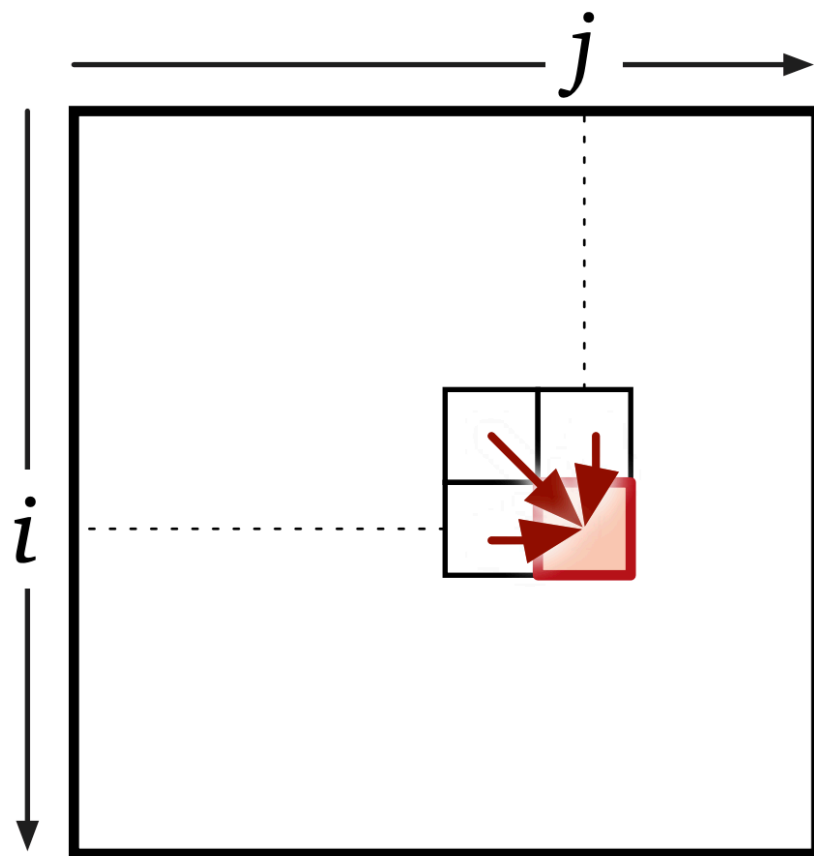
- Uses the shorthand:  $(a_i \neq b_j)$  which is 1 if it is true (and they mismatch), and zero otherwise

# From Recurrence to DP

- We can now transform it into a dynamic program
- **Memoization Structure**: We can memoize all possible values of  $\text{Edit}(i, j)$  in a table/ two-dimensional array of size  $O(nm)$ :
  - Store  $\text{Edit}[i, j]$  in a 2D array;  $0 \leq i \leq n$  and  $0 \leq j \leq m$
- **Evaluation order**:
  - Is interesting for a 2D problem
  - Based on dependencies between subproblems
  - We want values required to be already computed

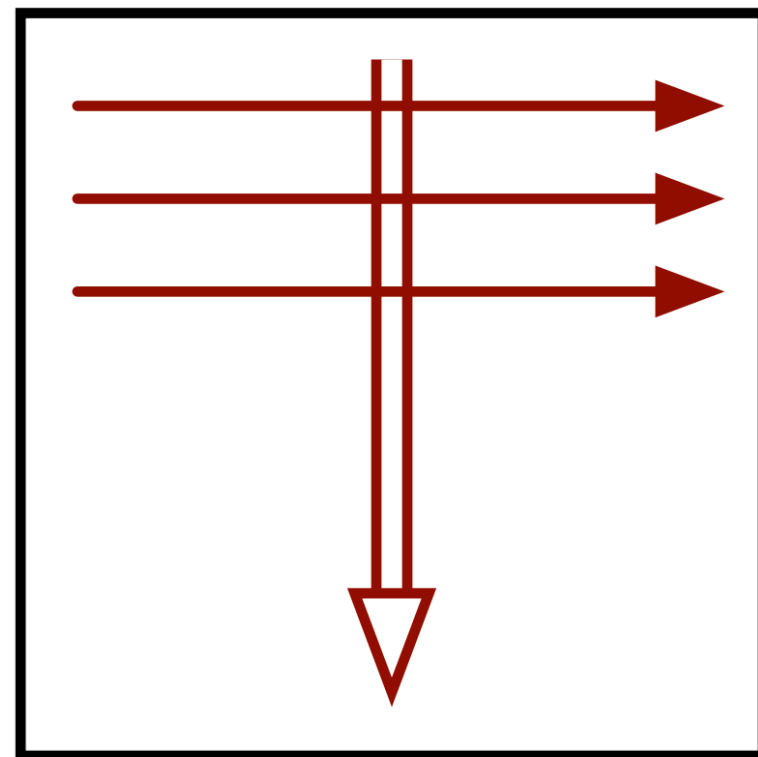
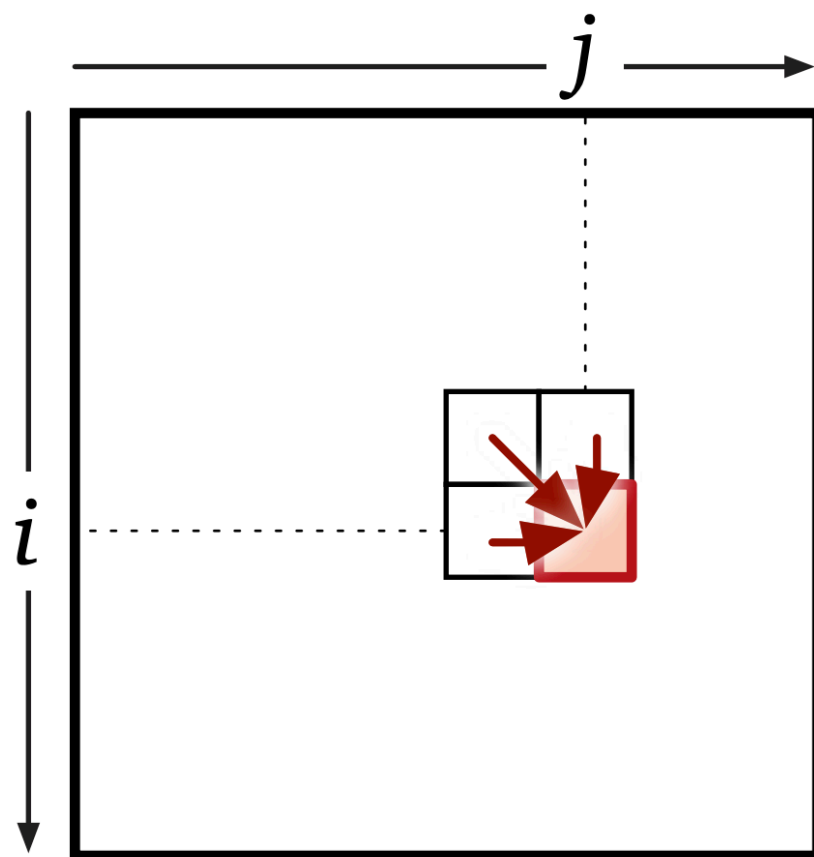
# From Recurrence to DP

- **Evaluation order**
  - We can fill in **row major order**, which is row by row from top down, each row from left to right: when we reach an entry in the table, it depends only on filled-in entries



# Space and Time

- The memoization uses  $O(nm)$  space
- We can compute each  $\text{Edit}[i, j]$  in  $O(1)$  time
- Overall running time:  $O(nm)$



# Memoization Table: Example

- Memoization table for **ALGORITHM** and **ALTRUISTIC**
- Bold numbers indicate where characters are same
- Horizontal arrow: deletion in *A*
- Vertical arrow: insertion in *A*
- Diagonal: substitution
- Bold red: free substitution
- Only draw an arrow if used in DP
- Any directed path of arrows from top left to bottom right represents an optimal edit distance sequence

		A	L	G	O	R	I	T	H	M
	0	1	2	3	4	5	6	7	8	9
A	1	<b>0</b>	1	2	3	4	5	6	7	8
L	2	1	<b>0</b>	1	2	3	4	5	6	7
T	3	2	1	1	2	3	4	<b>4</b>	5	6
R	4	3	2	2	2	<b>2</b>	3	4	5	6
U	5	4	3	3	3	3	3	4	5	6
I	6	5	4	4	4	4	<b>3</b>	4	5	6
S	7	6	5	5	5	5	4	4	5	6
T	8	7	6	6	6	6	5	<b>4</b>	5	6
I	9	8	7	7	7	7	<b>6</b>	5	5	6
C	10	9	8	8	8	8	7	6	6	6

# Reconstructing the Edits

- We don't need to store the arrow!
- Can be reconstructed on the fly in  $O(1)$  time using the numerical values
- Once the table is built, we can construct the shortest edit distance sequence in  $O(n + m)$  time

		A	L	G	O	R	I	T	H	M
	0	1	2	3	4	5	6	7	8	9
A	1	0	1	2	3	4	5	6	7	8
L	2	1	0	1	2	3	4	5	6	7
T	3	2	1	1	2	3	4	4	5	6
R	4	3	2	2	2	2	3	4	5	6
U	5	4	3	3	3	3	3	4	5	6
I	6	5	4	4	4	4	3	4	5	6
S	7	6	5	5	5	5	4	4	5	6
T	8	7	6	6	6	6	5	4	5	6
I	9	8	7	7	7	7	6	5	5	6
C	10	9	8	8	8	8	7	6	6	6

# In Pairs: Calculate ED

- What is the edit distance between SNOW and SOWS?
- $\text{Edit}(i, 0) = i$  and  $\text{Edit}(0, j) = j$
- For  $1 \leq i \leq n$  and  $1 \leq j \leq m$ , we have:

$$\text{Edit}(i, j) = \min \begin{cases} \text{Edit}(i, j - 1) + 1 \\ \text{Edit}(i - 1, j) + 1 \\ \text{Edit}(i - 1, j - 1) + (a_i \neq b_j) \end{cases}$$



# Edit Distance Fun Facts

- Can we compute edit distance using less space?
  - $O(nm)$  is huge for large genomic sequences
  - If we only care cost cost, we only need  $O(n + m)$  space (just keep row above current)
  - But this doesn't let us recreate the path
- **Hirschberg's algorithm:** Can compute the actual path (edits) in  $O(nm)$  time using  $O(n + m)$  space
  - Neat divide-and-conquer trick to save space

# Edit Distance Fun Facts

- Can we do better than  $O(n^2)$  if  $n = m$ ?
- Yes; can get  $O(n^2/\log^2 n)$  [Masek Paterson '80]
  - Uses “bit packing” trick called “Four Russians Technique” or “Indirection”
    - Not a great name; only one of them was actually Russian
- Can we get an algorithm for edit distance with runtime  $O(n^{2-\epsilon})$ , e.g.  $O(n^{1.9})$ ?
  - Probably not (unless a well-known conjecture breaks)

# Edit Distance Fun Facts

## Edit Distance Cannot Be Computed in Strongly Subquadratic Time (unless SETH is false)

Arturs Backurs  
MIT  
backurs@mit.edu

Piotr Indyk  
MIT  
indyk@mit.edu

### ABSTRACT

The edit distance (a.k.a. the Levenshtein distance) between two strings is defined as the minimum number of insertions, deletions or substitutions of symbols needed to transform one string into another. The problem of computing the edit distance between two strings is a classical computational task, with a well-known algorithm based on dynamic programming. Unfortunately, all known algorithms for this problem run in nearly quadratic time.

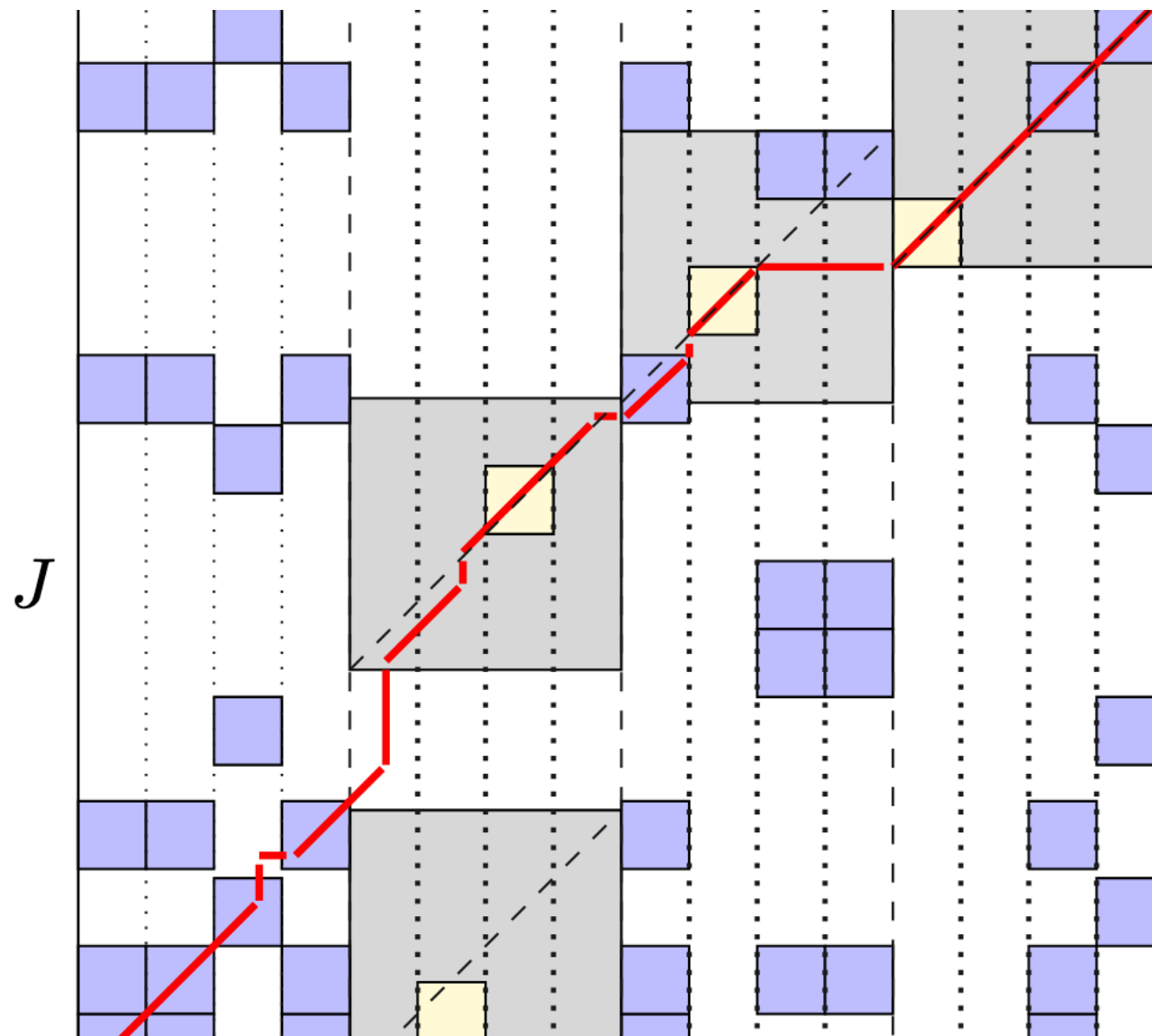
In this paper we provide evidence that the near-quadratic running time bounds known for the problem of computing edit distance might be tight. Specifically, we show that, if the edit distance can be computed in time  $O(n^{2-\delta})$  for some constant  $\delta > 0$ , then the satisfiability of conjunctive normal form formulas with  $N$  variables and  $M$  clauses can be solved in time  $M^{O(1)}2^{(1-\epsilon)N}$  for a constant  $\epsilon > 0$ . The latter result would violate the *Strong Exponential Time Hypothesis*, which postulates that such algorithms do not exist.

with many applications in computational biology, natural language processing and information theory. The problem of computing the edit distance between two strings is a classical computational task, with a well-known algorithm based on dynamic programming. Unfortunately, that algorithm runs in quadratic time, which is prohibitive for long sequences (e.g., the human genome consists of roughly 3 billions base pairs). A considerable effort has been invested into designing faster algorithms, either by assuming that the edit distance is bounded, by considering the average case or by resorting to approximation<sup>1</sup>. However, the fastest known exact algorithm, due to [MP80], has a running time of  $O(n^2/\log^2 n)$  for sequences of length  $n$ , which is still nearly quadratic.

In this paper we provide evidence that the (near)-quadratic running time bounds known for this problem might, in fact, be tight. Specifically, we show that if the edit distance can be computed in time  $O(n^{2-\delta})$  for some constant  $\delta > 0$ , then the satisfiability of conjunctive normal form (CNF) formulas with  $N$  variables and  $M$  clauses can be solved in time

# Edit Distance Fun Facts

- Can approximate to any  $1 + \epsilon$  factor in  $O(n)$  time!  
[Andoni Nosatski '20]



A figure from [CDGKS'18], the first approximation algorithm for edit distance.  
The idea: rule out large portions of the dynamic programming table