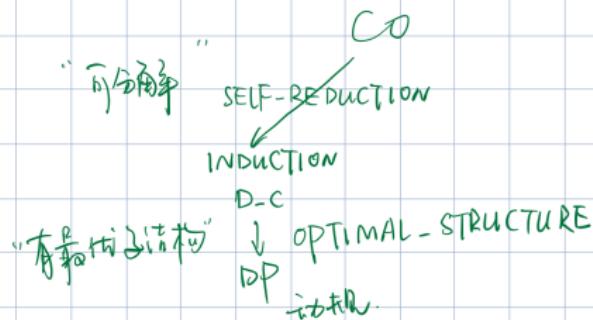


CS711008Z Algorithm Design and Analysis

Lecture 6. Basic algorithm design technique: Dynamic programming

Dongbo Bu



Institute of Computing Technology
Chinese Academy of Sciences, Beijing, China

“优化目标函数”

动机目的

- The first example: MATRIXCHAINMULTIPLICATION
- Elements of dynamic programming technique;
- Various ways to describe subproblems: SEGMENTED LEAST SQUARES, KNAPSACK, RNA SECONDARY STRUCTURE, SEQUENCE ALIGNMENT, and SHORTEST PATH;
- Connection with greedy technique: INTERVAL SCHEDULING, SHORTEST PATH.

Dynamic programming and its connection with divide-and-conquer

动态规划

- Dynamic programming typically applies to the **optimization** ^{问题} **problems with optimal-substructure property**; that is, the original problem can be divided into smaller subproblems, and the optimal solution to the original problem can be obtained by **combining the optimal solutions to subproblems**.
- To identify meaningful recursions, one of the key steps is to define **an appropriate general form of sub-problems**. For this aim, it is helpful to describe the solving process as **a multiple-stage decision process**.
- Unlike the general divide-and-conqueror framework, a dynamic programming algorithm usually **enumerate** all possible dividing strategies. In addition, **the repetition of computing the common subproblems** is avoided through “programing”.

1

动态规划与分治法

¹Here, “programming” means “tabular” rather than “coding”. The best example of programming might be the calculation of Fibonacci numbers.

Program: [Date: 1600-1700; Language: French; 'to write before']

- Suppose a problem is related to the following data structure, perhaps we can try to divide it into sub-problems.
 - An array with n elements
 - A matrix
 - A set of n elements
 - A tree
 - A graph

MATRIXCHAINMULTIPLICATION problem: recursion over
sequences

INPUT:

A sequence of n matrices A_1, A_2, \dots, A_n ; matrix A_i has dimension $p_{i-1} \times p_i$;

OUTPUT:

Fully parenthesizing the product $A_1A_2\dots A_n$ in a way to minimize the number of scalar multiplications.

Let's start from a simple example

$$A_1 = \begin{bmatrix} 1 & 2 \end{bmatrix} \quad A_2 = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} \quad A_3 = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

1×2 2×3 3×4

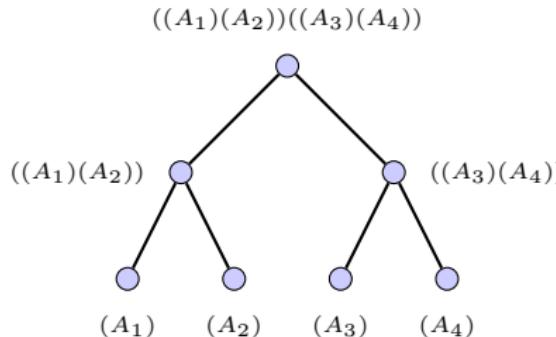
Solutions: $((A_1)(A_2))(A_3)$ $(A_1)((A_2)(A_3))$

#Multiplications:

$$\begin{array}{ll} 1 \times 2 \times 3 & 2 \times 3 \times 4 \\ +1 \times 3 \times 4 & +1 \times 2 \times 4 \\ = 18 & = 32 \end{array}$$

- Here, the calculation of $A_1 A_2$ needs $1 \times 2 \times 3$ scalar multiplications.
- The objective is to determine a calculation sequence such that the number of multiplications is minimized.

- Intuitively, a calculation sequence can be described as a binary tree, where each node corresponds to a subproblem.



- The total number of possible calculation sequences:
$$\binom{2n}{n} - \binom{2n}{n-1}$$
 (Catalan number)
- Thus, it takes exponential time to enumerate all possible calculation sequences.
- Question: can we design an efficient algorithm?

A dynamic programming algorithm (by S. S. Godbole, 1973?)

- ① It is not easy to solve the problem directly when n is large.
Let's investigate whether it is possible to reduce into smaller sub-problems.
- ② Solution: a full parentheses. Let's describe the solving process as a process of **multiple-stage decisions**, where each decision is to add parentheses at a position.
- ③ Suppose we have already worked out the optimal solution O , where the first **decision** adds two parentheses as $(A_1 \dots A_k)(A_{k+1} \dots A_n)$.
- ④ This decision decomposes the original problem into two **independent sub-problems**: to calculate $A_1 \dots A_k$ and $A_{k+1} \dots A_n$.
- ⑤ Summarizing these two cases, we define the general form of sub-problems as: to calculate $A_i \dots A_j$ with the minimal number of scalar multiplications.

General form of sub-problems and optimal substructure property

- The general form of sub-problems: to calculate $A_i \dots A_j$ with the minimal number of scalar multiplications. Let's denote the optimal solution value to the sub-problem as $OPT(i, j)$, thus the original problem can be solved via calculating $OPT(1, n)$.
- The optimal solution to the original problem can be obtained through combining the optimal solutions to sub-problems, implying the following **optimal substructure property**:
$$OPT(1, n) = OPT(1, k) + OPT(k + 1, n) + p_1 p_{k+1} p_{n+1}$$

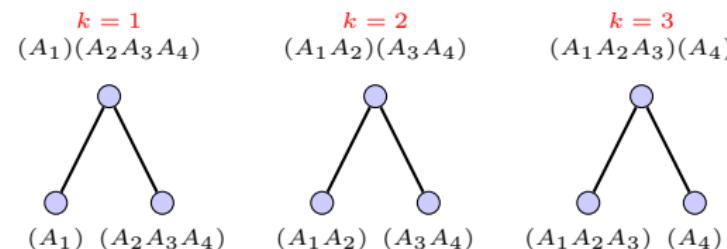
- “Cut-and-paste” proof:

- Suppose for $A_i \dots A_k$, there is another parentheses $OPT'(i, k)$ better than $OPT(i, k)$. Then the combination of $OPT'(i, k)$ and $OPT(k+1, j)$ leads to a new solution with lower cost than $OPT(i, j)$: a contradiction.
- Here, the independence between $A_i \dots A_k$ and $A_{k+1} \dots A_j$ guarantees that the substitution of $OPT(i, k)$ with $OPT'(i, k)$ does not affect solution to $A_{k+1} \dots A_j$.

问题往往是由错误引起的。

不知道 k 在哪
⇒ 枚举.

- So far so good! The only difficulty is that we have no idea of the first splitting position k in the optimal solution.
- How to overcome this difficulty? **Enumeration!** We **enumerate all possible options of the first decision**, i.e. for all k , $i \leq k < j$.



- Thus we have the following recursion:

$$OPT(i, j) = \begin{cases} 0 & i = j \\ \min_{i \leq k < j} \{OPT(i, k) + OPT(k+1, j) + p_i p_{k+1} p_{j+1}\} & \text{otherwise} \end{cases}$$

枚举 和 DC 不相同的地步.

Implementing the recursion: trial 1

Trial 1: Explore the recursion in the top-down manner

RECURSIVE_MATRIX_CHAIN(i, j)

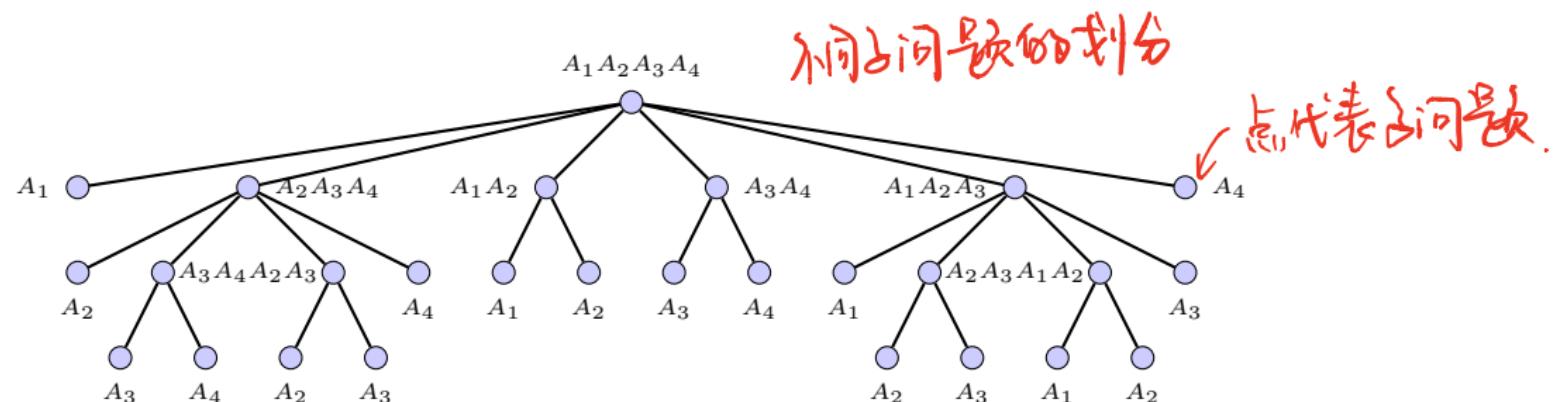
```
1: if  $i == j$  then
2:     return 0;
3: end if
4:  $OPT(i, j) = +\infty;$ 
5: for  $k = i$  to  $j - 1$  do  $\min_{i \leq k < j}$ 
6:      $q = \text{RECURSIVE\_MATRIX\_CHAIN}(i, k)$ 
7:     +  $\text{RECURSIVE\_MATRIX\_CHAIN}(k + 1, j)$ 
8:     +  $p_i p_{k+1} p_{j+1};$ 
9:     if  $q < OPT(i, j)$  then
10:         $OPT(i, j) = q;$ 
11:    end if
12: end for
13: return  $OPT(i, j);$ 
```

- Note: The optimal solution to the original problem can be obtained through calling
RECURSIVE_MATRIX_CHAIN(1, n).

An example

$$A_1 = \begin{bmatrix} 1 & 2 \end{bmatrix} \quad A_2 = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} \quad A_3 = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix} \quad A_4 = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 1 & 2 & 3 & 4 & 5 \\ 1 & 2 & 3 & 4 & 5 \\ 1 & 2 & 3 & 4 & 5 \end{bmatrix}$$

1×2 2×3 3×4 4×5



Note: each node of the recursion tree denotes a subproblem.

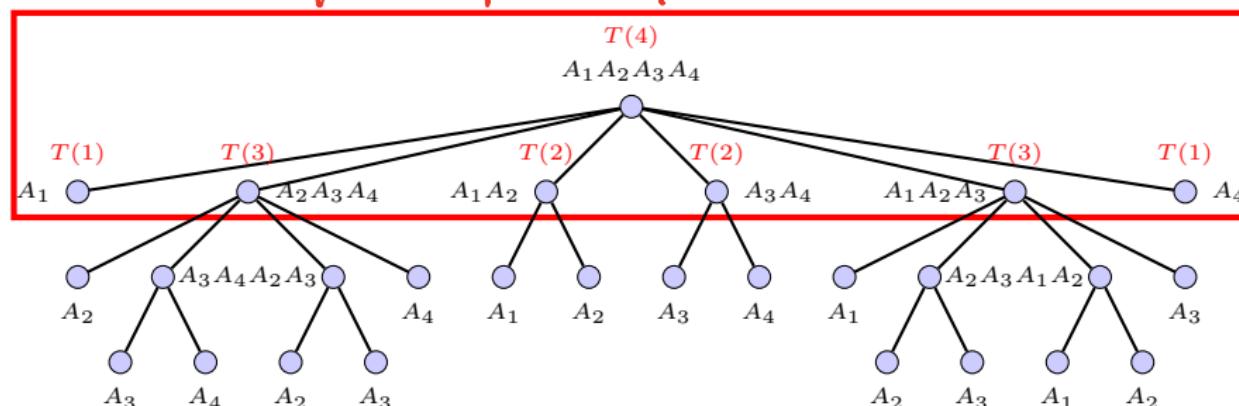
Theorem

Algorithm RECURSIVE-MATRIX-CHAIN costs exponential time.

Let $T(n)$ denote the time used to calculate product of n matrices.

Note that $T(n) \geq 1 + \sum_{k=1}^{n-1} (T(k) + T(n-k) + 1)$ for $n > 1$.

$$T(4) \geq 2T(1) + 2T(2) + 2T(3) \Rightarrow T(n) \geq 2^n$$



Theorem

Algorithm RECURSIVE-MATRIX-CHAIN costs exponential time.

Proof.

- We shall prove $T(n) \geq 2^{n-1}$ using the substitution technique.
 - Basis: $T(1) \geq 1 = 2^{1-1}$.
 - Induction:

$$T(n) \geq 1 + \sum_{k=1}^{n-1} (T(k) + T(n-k) + 1) \quad (1)$$

$$= n + 2 \sum_{k=1}^{n-1} T(k) \quad (2)$$

$$\geq n + 2 \sum_{k=1}^{n-1} 2^{k-1} \quad (3)$$

$$\geq n + 2(2^{n-1} - 1) \quad (4)$$

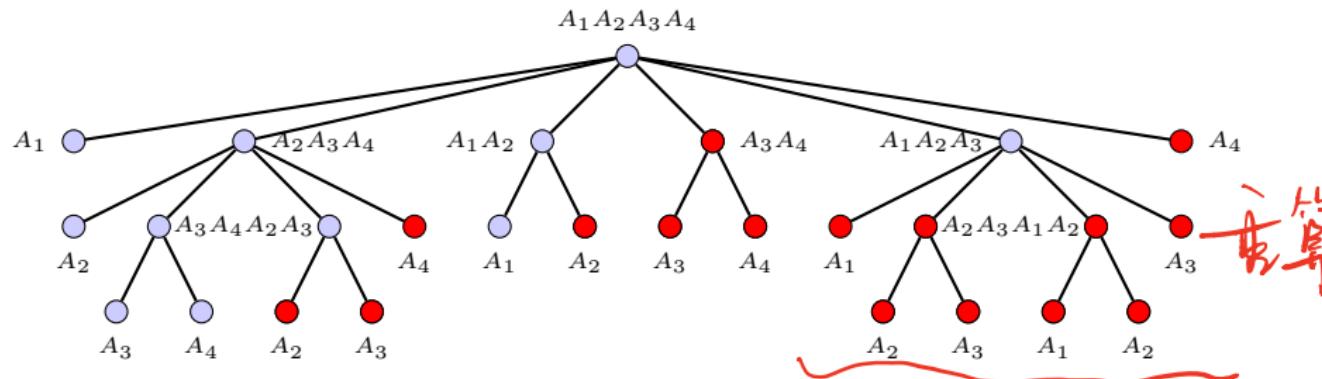
$$\geq n + 2^n - 2 \quad (5)$$

$$\geq 2^{n-1} \quad (6)$$

Implementing the recursion: trial 2

Why the first trial failed?

n^2 个问题，却花费了 2^n 时间 \Rightarrow 有问题重叠计算.



- Reason: There are only $O(n^2)$ subproblems. However, some subproblems (in red) were solved repeatedly. *记录所有算过的问題.*
- Solution: memorize the solutions to subproblems using an array $OPT[1..n; 1..n]$ for further look-up.

The “memorizing” technique

```
MEMORIZE_MATRIX_CHAIN( $i, j$ )
1: if  $OPT[i, j] \neq \text{NULL}$  then
2:   return  $OPT(i, j)$ ;
3: end if
4: if  $i == j$  then
5:    $OPT[i, j] = 0$ ;
6: else
7:   for  $k = i$  to  $j - 1$  do
8:      $q = \text{MEMORIZER\_MATRIX\_CHAIN}(i, k)$ 
9:     + $\text{MEMORIZER\_MATRIX\_CHAIN}(k + 1, j)$ 
10:    + $p_i p_{k+1} p_{j+1}$ ;
11:    if  $q < OPT[i, j]$  then
12:       $OPT[i, j] = q$ ;
13:    end if
14:   end for
15: end if
16: return  $OPT[i, j]$ ;
```

- The original problem can be solved by calling `MEMORIZE_MATRIX_CHAIN(1, n)` with all $OPT[i, j]$ initialized as `NULL`.
- Time-complexity: $O(n^3)$ (The calculation of each entry $OPT[i, j]$ makes $O(n)$ recursive calls in line 8.)
- Note: the calculation of Fibonacci number is a good example of the power of the “memorizing” technique.

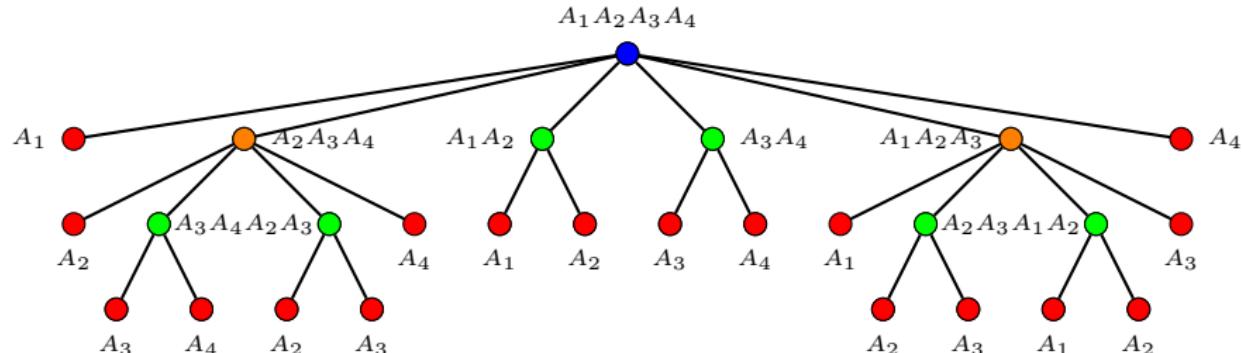
Implementing the recursion faster: trial 3

Trial 3: Faster implementation: unrolling the recursion in the bottom-up manner

MATRIX_CHAIN_MULTIPLICATION(P)

```
1: for  $i = 1$  to  $n$  do
2:    $OPT(i, i) = 0;$ 
3: end for
4: for  $l = 2$  to  $n$  do
5:   for  $i = 1$  to  $n - l + 1$  do
6:      $j = i + l - 1;$ 
7:      $OPT(i, j) = +\infty;$ 
8:     for  $k = i$  to  $j - 1$  do
9:        $q = OPT(i, k) + OPT(k + 1, j) + p_i p_{k+1} p_{j+1};$ 
10:      if  $q < OPT(i, j)$  then
11:         $OPT(i, j) = q;$ 
12:         $S(i, j) = k;$ 
13:      end if
14:    end for
15:  end for
16: end for
17: return  $OPT(1, n);$ 
```

Recursion tree: an intuitive view of the bottom-up calculation



Solving sub-problems in a **bottom-up** manner, i.e.

非递归方法.

- ① Solving the sub-problems in red first;
- ② Then solving the sub-problems in green;
- ③ Then solving the sub-problems in orange;
- ④ Finally we can solve the original problem in blue.

Step 1 of the bottom-up algorithm

所有区间嵌套优化解

OPT

1	2	3	4
0	6		
	0	24	
		0	60
			0

SPLITTER

1	2	3	4
	1		
		2	
			3
			4

Step 1:

$$OPT[1, 2] = p_0 \times p_1 \times p_2 = 1 \times 2 \times 3 = 6;$$

$$OPT[2, 3] = p_1 \times p_2 \times p_3 = 2 \times 3 \times 4 = 24;$$

$$OPT[3, 4] = p_2 \times p_3 \times p_4 = 3 \times 4 \times 5 = 60;$$

Step 2 of the bottom-up algorithm

OPT				SPLITTER			
1	2	3	4	1	2	3	4
0	6	18			1		
	0	24	64	2		2	
		0	60	3		3	
			0	4			4

Step 2: $A_1 A_2 A_3 \Rightarrow (A_1 A_2) A_3$ Combine

$$OPT[1, 3] = \min \begin{cases} OPT[1, 2] + OPT[3, 3] + p_0 \times p_2 \times p_3 (= 18) \\ OPT[1, 1] + OPT[2, 3] + p_0 \times p_1 \times p_3 (= 32) \end{cases}$$

Thus, $SPLITTER[1, 2] = 2$. $A_1 (A_2 A_3)$ combine

$$OPT[2, 4] = \min \begin{cases} OPT[2, 2] + OPT[3, 4] + p_1 \times p_2 \times p_4 (= 90) \\ OPT[2, 3] + OPT[4, 4] + p_1 \times p_3 \times p_4 (= 64) \end{cases}$$

Thus, $SPLITTER[2, 4] = 3$.

Step 3 of the bottom-up algorithm

OPT				SPLITTER			
1	2	3	4	1	2	3	4
0	6	18	38		1	2	3
	0	24	64	1		2	3
	0	60		2		3	
	0			3		3	
			0	4			

Step 3:

$$OPT[1, 4] = \min \begin{cases} OPT[1, 1] + OPT[2, 4] + p_0 \times p_1 \times p_4 (= 74) \\ OPT[1, 2] + OPT[3, 4] + p_0 \times p_2 \times p_4 (= 81) \\ OPT[1, 3] + OPT[4, 4] + p_0 \times p_3 \times p_4 (= 38) \end{cases}$$

Thus, $SPLITTER[1, 4] = 3$.

Question: We have calculated the optimal **value**, but how to get the optimal **calculation sequence**?

Final step: constructing an optimal solution through “backtracking” the optimal options

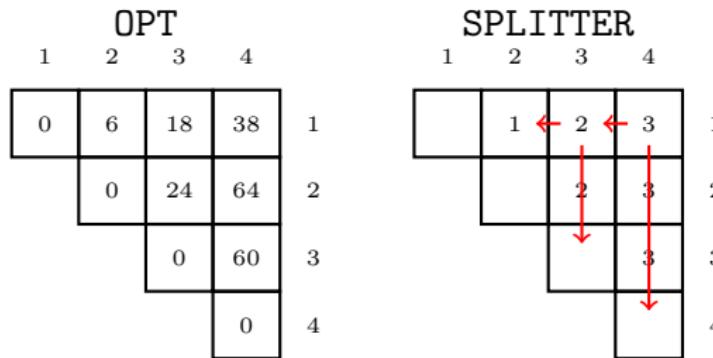
- Idea: **backtracking!** Starting from $OPT[1, n]$, we trace back the source of $OPT[1, n]$, i.e. which option we take at each decision stage.
- Specifically, an auxiliary array $S[1..n, 1..n]$ is used.
 - Each entry $S[i, j]$ records the optimal decision, i.e. the value of k such that the optimal parentheses of $A_i \dots A_j$ occurs between $A_k A_{k+1}$.
 - Thus, the optimal solution to the original problem $A_{1..n}$ is $A_{1..S[1,n]} A_{S[1,n]+1..n}$.
- Note: The optimal option cannot be determined before solving all subproblems.

Backtracking: step 1

OPT				SPLITTER			
1	2	3	4	1	2	3	4
0	6	18	38		1	2	3
	0	24	64			2	3
		0	60			3	3
			0			4	

Step 1: $(A_1 A_2 A_3) (A_4)$

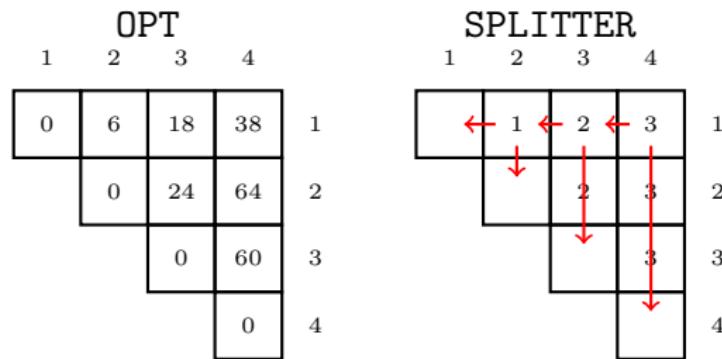
Backtracking: step 2



Step 1: $(A_1 A_2 A_3) (A_4)$

Step 2: $((A_1 A_2) (A_3)) (A_4)$

Backtracking: step 3



Step 1: $(A_1 A_2 A_3) (A_4)$

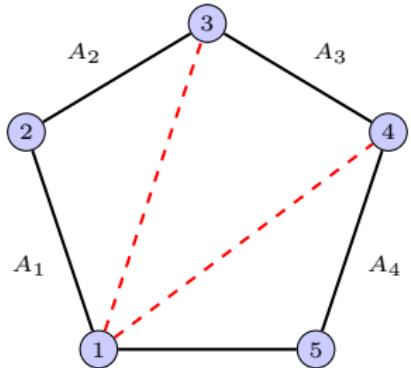
Step 2: $((A_1 A_2) (A_3)) (A_4)$

Step 3: $(((A_1) (A_2) (A_3)) (A_4)$

- ① It is usually not easy to solve a large problem directly. Let's consider whether the problem can be decomposed into smaller sub-problems.
How to define sub-problems? ²
 - Let's describe the solving process as a process of multiple-stage **decisions** first.
 - Suppose that we have already worked out the optimal solution. Let's consider the **first/final decision (in some order)** in the optimal solution. The **first/final decision** might have several options.
 - We enumerate all possible options for the decision, and observe the generated sub-problems. The general form of sub-problems can be defined via summarising all possible forms of sub-problems.
- ② Show the **optimal substructure property**, i.e. the optimal solution to the problem contains within it optimal solutions to subproblems.
- ③ Programming: if recursive algorithm solves the same subproblem over and over, “tabular” can be used to avoid the repetition of solving same sub-problems.

²Sometimes problem should be extended to identify meaningful recursion.

Question: is $O(n^3)$ the lower bound?



$O_1: 123 \Rightarrow A_1 A_2$

$O_2: 134 \Rightarrow (A_1 A_2) (A_3)$

$O_3: 145 \Rightarrow (A_1 A_2 A_3) (A_4)$

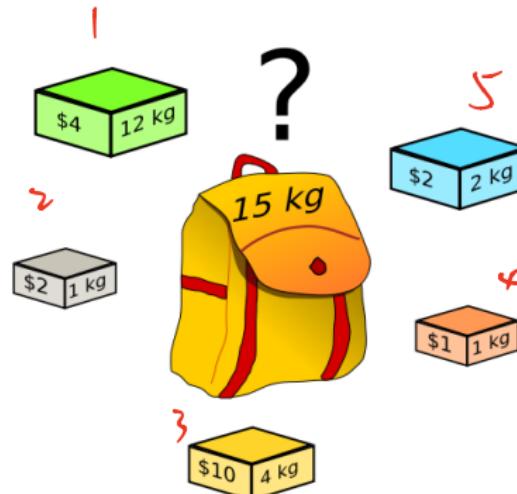
- One-to-one correspondence between parenthesis and partitioning a convex polygon into non-intersecting triangles.
 - Each node has a weight w_i , and a triangle corresponds to a product of the weight of its nodes.
 - The decomposition (red, dashed lines) has a weight sum of 38. In fact, it corresponds to the parenthesis $(((A_1) (A_2) (A_3)) (A_4)$.
- The optimal decomposition can be found in $O(n \log n)$ time.

(See Hu and Shing 1981 for details)

0/1 KNAPSACK problem: recursion over **sets**

A Knapsack instance

Given a set of items, each item has a weight and a value, to select a subset of items such that the total weight is less than a given limit and the total value is as large as possible.



$$\left\{ \begin{array}{l} 12a_1 + 1a_2 + 3a_3 + 1a_4 + 2a_5 \leq 15 \\ \max (4a_1 + 2a_2 + 10a_3 + 1a_4 + 2a_5) \\ \text{within } a_i \in \{0, 1\} \end{array} \right.$$

What's the best solution?

Formalized Definition:

- **Input:**

A set of items. Item i has weight w_i and value v_i , and a total weight limit W ;

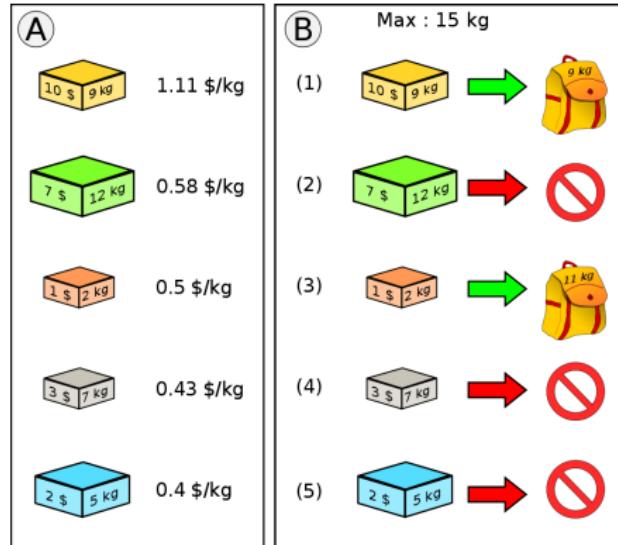
- **Output:**

A sub-set of items to maximize the total value with a total weight below W .

Note:

- ① Here, “0/1” means that we should select an item (1) or abandon it (0), and we cannot select parts of an item.
- ② In contrast, FRACTIONAL KNAPSACK problem allow one to select a fractional, say 0.5, of an item.

0/1 Knapsack problem: an intuitive algorithm



- Intuitive method: selecting “expensive” items first.
- But this is not the optimal solution.

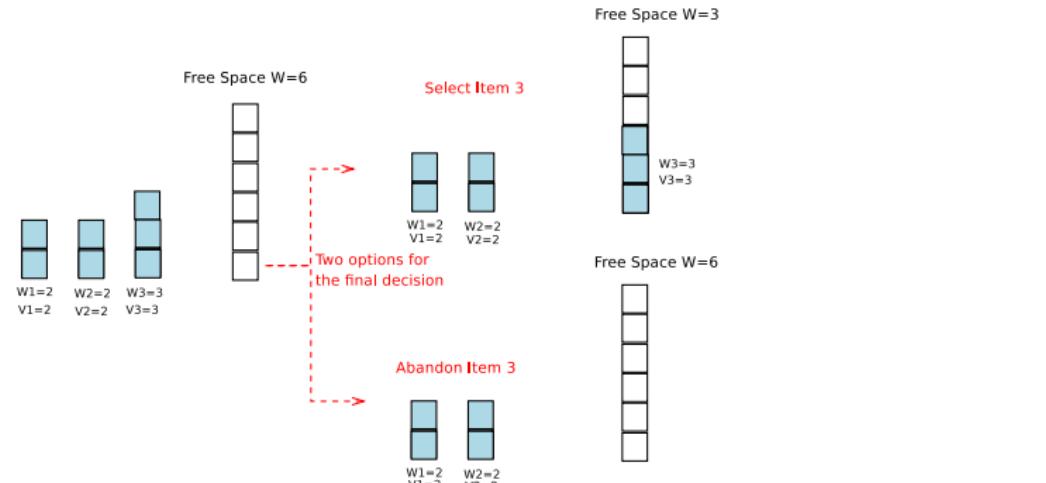
- It is not easy to solve the problem with n items. Let's whether it is possible to reduce into smaller sub-problems.
- Solution: a subset of items. Describe the solving process as a process of multiple-stage **decisions**. At the i -th decision stage, we decide whether item i should be selected.
- Suppose we have already worked out the optimal solution.
- Consider the **first** decision, i.e. whether the optimal solution contains item n or not. The decision has two options:
 - ① SELECT: then it suffices to select items as "expensive" as possible from $\{1, 2, \dots, n-1\}$ with weight limit $W - w_n$.
 - ② ABANDON: Otherwise, we should select items as "expensive" as possible from $\{1, 2, \dots, n-1\}$ with weight limit W .
- In both cases, the original problem is reduced into smaller sub-problems.

隨意的寫作

$$\max \{ \$n + w_n, W_n \}$$

- Summarizing these two cases, the general form of sub-problems is: to select items as “expensive” as possible from $\{1, 2, \dots, i\}$ with weight limit w . Denote the optimal solution value as $OPT(i, w)$.

- Optimal sub-structure property: 不装 \downarrow 装 \downarrow .
 $OPT(n, W) = \max\{OPT(n - 1, W), OPT(n - 1, W - w_n) + v_n\}$
 (Enumerating two possible decisions for item n .)



KNAPSACK(n, W)

```
1: for  $w = 1$  to  $W$  do
2:    $OPT[0, w] = 0;$ 
3: end for
4: for  $i = 1$  to  $n$  do
5:   for  $w = 1$  to  $W$  do
6:      $OPT[i, w] = \max\{OPT[i-1, w], v_i + OPT[i-1, w-w_i]\};$ 
7:   end for
8: end for
```

Example I

1: {2kg, 2\$}

2: {2kg, 2\$}

3: {3kg, 3\$}

w=6 kg

Initially all $OPT[0,w] = 0$

前*i*個物

包有W丁空

	W = 0	1	2	3	4	5	6
i=3							
i=2							
i=1							
i=0	0	0	0	0	0	0	0

添值多少.

Example II

$$\begin{aligned} \text{OPT}[1,2] &= \max\{ \\ &\quad \text{OPT}[0,2] (=0), \\ &\quad \text{OPT}[0,0] + V_1 (=0+2) \} \\ &= 2 \end{aligned}$$

	W = 0	1	2	3	4	5	6
i=3							
i=2							
i=1	0	0	2	2	2	2	2
i=0	0	0	0	0	0	0	0

Example III

$$\begin{aligned} \text{OPT}[2,4] &= \max\{ \\ &\text{OPT}[1,4] (=2), \\ &\text{OPT}[1,2] + V_2 (=2+2) \} \\ &= 4 \end{aligned}$$

	W = 0	1	2	3	4	5	6
i=3							
i=2	0	0	2	2	4	4	4
i=1	0	0	2	2	2	2	2
i=0	0	0	0	0	0	0	0

Example IV

$$\begin{aligned} \text{OPT}[3,3] &= \max\{ \\ &\text{OPT}[2,3] (=2), \\ &\text{OPT}[2,0] + V_3 (=0+3) \} \\ &= 3 \end{aligned}$$

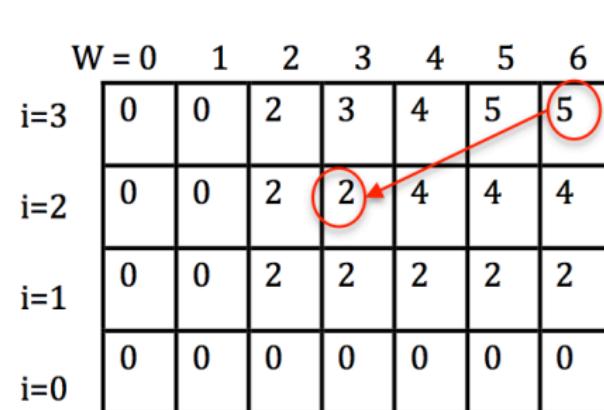
	W = 0	1	2	3	4	5	6
i=3	0	0	2	3	4	5	5
i=2	0	0	2	2	4	4	4
i=1	0	0	2	2	2	2	2
i=0	0	0	0	0	0	0	0

Backtracking: step 1

Backtracking

$$\begin{aligned} \text{OPT}[3,6] = \max\{ & \\ & \text{OPT}[2,6] (=4), \\ & \text{OPT}[2,3] + V_3 (=2+3) \} \\ &= 5 \end{aligned}$$

Decision: Select item 3



	W = 0	1	2	3	4	5	6
i=3	0	0	2	3	4	5	5
i=2	0	0	2	2	4	4	4
i=1	0	0	2	2	2	2	2
i=0	0	0	0	0	0	0	0

回溯

Backtracking: step 2

Backtracking

$$\begin{aligned} \text{OPT}[3,6] &= \max\{ \\ &\text{OPT}[2,6] (=4), \\ &\text{OPT}[2,3] + V_3 (=2+3) \} \\ &= 5 \end{aligned}$$

Decision: Select item 3

$$\begin{aligned} \text{OPT}[2,3] &= \max\{ \\ &\text{OPT}[1,3] (=2), \\ &\text{OPT}[1,1] + V_2 (=0+2) \} \\ &= 2 \end{aligned}$$

Decision: Select item 2

	W = 0	1	2	3	4	5	6
i=3	0	0	2	3	4	5	5
i=2	0	0	2	2	4	4	4
i=1	0	0	2	2	2	2	2
i=0	0	0	0	0	0	0	0

($n+1$ 行, $(W+1)$ 列), 每个单元需 2 时间.

- Time complexity: $O(nW)$. (Hint: for each entry in the matrix, only a comparison is needed; we have $O(nW)$ entries in the matrix.)

W 值多大呢?

- Notes:

① This algorithm is inefficient when W is large, say $W = 1M$.

② Remember that a polynomial time algorithm costs time polynomial in the **input length**. However, this algorithm costs time $mW = m2^{\log W} = m2^{\text{input length}}$. Exponential!

W 长度的指数.

③ Pseudo-polynomial time algorithm: polynomial in the **value** of W rather than the **length** of W ($\log W$).

④ We will revisit this algorithm in approximation algorithm design.

伪多项式时间算法.

To select $S \subseteq \{1, 2, \dots, n\}$, max OPT. $w \rightarrow$ 慢.

STP 2^n 神選擇

- Here the items were considered in an order. Why?
- Let's consider two general forms of sub-problems:
 - ① Selecting items as “expensive” as possible from a subset s with weight limit w : the number of sub-problems is **exponential**.
 - ② Selecting items as “expensive” as possible from $\{1, 2, \dots, i\}$ with weight limit w : the number of sub-problems is **$O(n)$** .
- In fact, the first one is a **recursion over sets**, while the second one is a **recursion over sequences**.

- *Cryptosystems based on the knapsack problem were among the first public key systems to be invented, and for a while were considered to be among the most promising. However, essentially all of the knapsack cryptosystems that have been proposed so far have been broken. These notes outline the basic constructions of these cryptosystems and attacks that have been developed on them.*

See *The Rise and Fall of Knapsack Cryptosystems* for details.

VERTEX COVER: recursion over **trees**

- Practical problem:

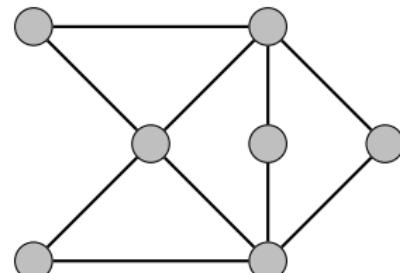
Given n sites connected with paths, how many guards (or cameras) should be deployed on sites to surveille all the paths?

Formalized Definition:

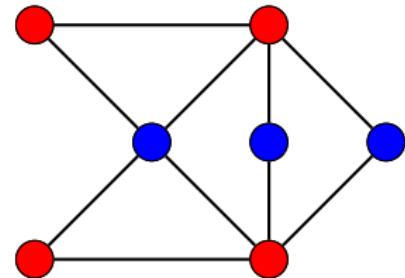
Input: Given a graph $G = \langle V, E \rangle$

Output: the minimum of nodes $S \subseteq V$, such that each edge has at least one of its endpoints in S

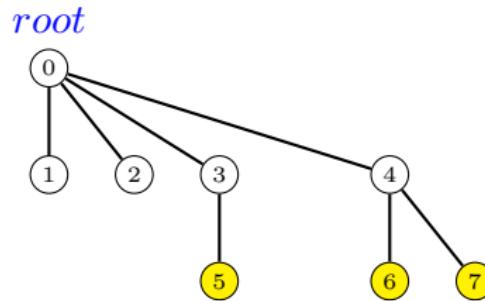
- An example: how many nodes are needed to cover all edges in the following graph?



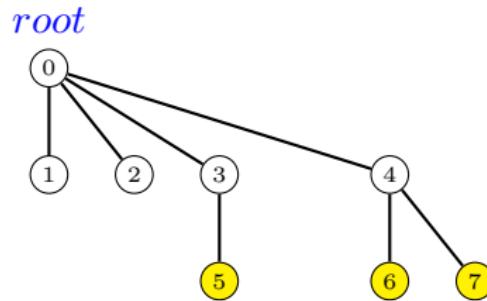
NP-hard



- The nodes in red form a vertex cover.
- VERTEX COVER is a hard problem for general graph.
- However, it is easy to find the minimum vertex cover for **trees**.



- It is not easy to solve the problem with n nodes. Let's whether it is possible to reduce into smaller sub-problems.
- Solution: selection a subset of nodes. Describe the solving process as a process of multiple-stage **decisions**. At each decision stage, we decide whether a node should be selected.
- Suppose we have already worked out the optimal solution.
- Consider the **first** decision, i.e. whether the optimal solution contains the **root** node or not. The decision has two options:
 - ① **SELECT**: it suffices to consider the sub-trees;
 - ② **ABANDON**: we should select all the children nodes, and then consider all grand-children.



- In both cases, the original problem is reduced into smaller sub-problems.
- General form of sub-problems: find minimum vertex cover on a tree rooted at node v . Let's denote the optimal solution as $OPT(v)$.
- Thus we have the following recursion:

Root 裝 → $OPT(root) = \begin{cases} 1 + \sum_c OPT(c) & c \text{ is a children of root} \\ k + \sum_g OPT(c) & c \text{ is a grand-children of root, } k = \#children \end{cases}$

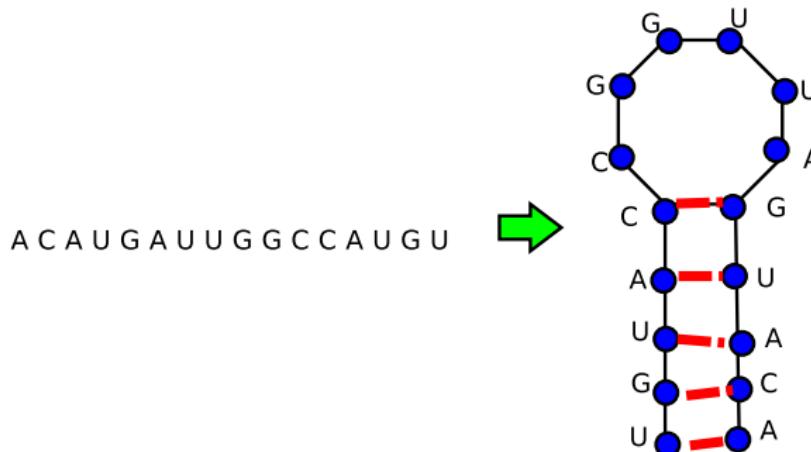
Root 不裝 →

- Time complexity: $O(n)$

每个儿子节点都装.

RNA SECONDARY STRUCTURE PREDICTION: recursion over
trees 跳过

- RNA is a sequence of nucleic acids. It will automatically form structures in water through the formation of bonds $A - U$ and $C - G$.
- The native structure is the conformation with the lowest energy. Here, we simply use the number of base pairs as the energy function.



INPUT:

A sequence in alphabet $\Sigma = \{A, U, C, G\}$;

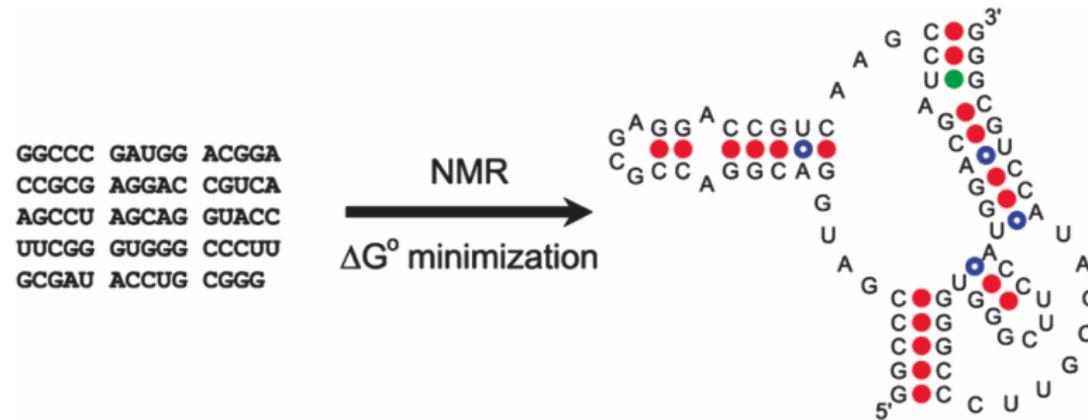
OUTPUT:

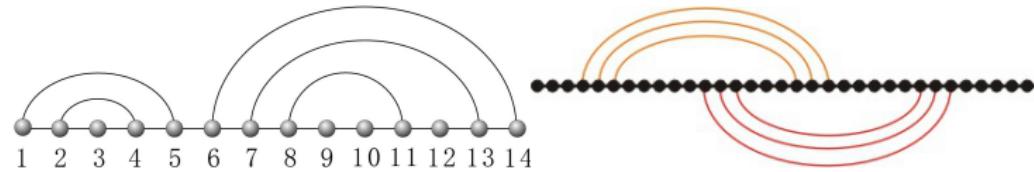
A pairing scheme with the maximum pairing number

Requirements of base pairs:

- ① Watson-Crick pair: A pairs with U , and C pairs with G ;
- ② There is no base occurring in more than 1 base pairs;
- ③ No cross-over (nesting): there is no crossover under the assumption of free pseudo-knots.
- ④ And two bases i, j ($|i - j| \leq 4$) cannot form a base pair.

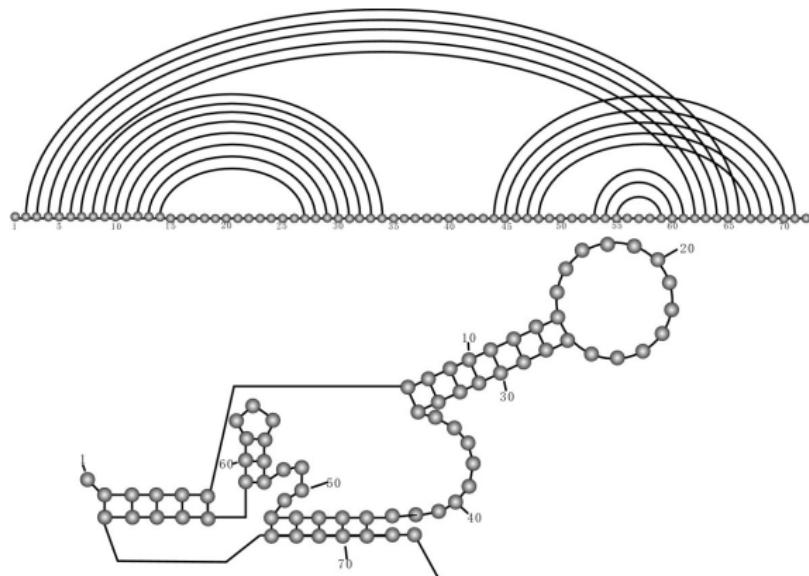
An example





Left: nesting of base pairs (no cross-over); Right: pseudo-knots (cross-over);

Feymann graph



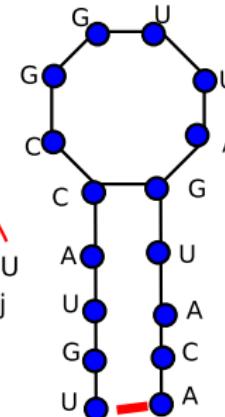
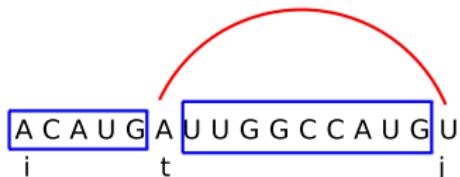
Feymann graph: an intuitive representation form of RNA secondary structure, i.e. two bases are connected by an edge if they form a Watson-Crick pair.

- Solution: a set of **nested** base pairs. Describe the solving process as a process of multiple-stage **decisions**. At the i -th decision stage, we determine whether base i forms pair or not.
- Suppose we have already worked out the optimal solution.
- Consider the *first* decision made for base n . There are two options:
 - ① Base n pairs with a base i : we should calculate optimal pairs for regions $i + 1 \dots n - 1$ and $1..i - 1$.
Note: these two sub-problems are independent due to the “nested” property.
 - ② Base n doesn’t form a pair: we should calculate optimal pairs for regions $1 \dots n - 1$.
- Thus we can design the general form of sub-problems as: to calculate the optimal pairs for region $i \dots j$. (Denote the optimal solution value as: $OPT(i, j)$.)

Key observation II

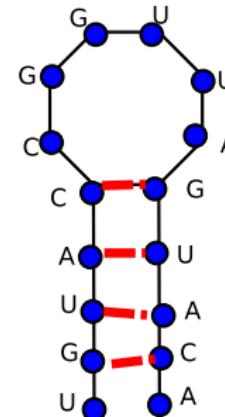
Option 1:

j and t form a pair, $t \leq j-1$



Option 2:

j forms no pair



- Optimal substructures property:

RNA2D(n)

```
1: Initialize all  $OPT[i, j]$  with 0;  
2: for  $i = 1$  to  $n$  do  
3:   for  $j = i + 5$  to  $n$  do  
4:      $OPT[i, j] = \max\{OPT[i, j - 1], \max_t\{1 + OPT[i, t - 1] + OPT[t + 1, j - 1]\}\};$   
5:     /*  $t$  and  $j$  can form Watson-Crick base pair. */  
6:   end for  
7: end for
```

INPUT:

ACCGGUAGU

i=3	0	0	0	
i=2	0	0		
i=1	0			
i=0				

j = 6 7 8 9

INPUT:

A C C G G U A G U

i=3	0	0	0	1
i=2	0	0	1	
i=1	0	0		
i=0	1			
	j = 6	7	8	9

INPUT:

A C C G G U A G U

i=3	0	0	0	1
i=2	0	0	1	1
i=1	0	0	1	
i=0	1	1		

j = 6 7 8 9

INPUT:

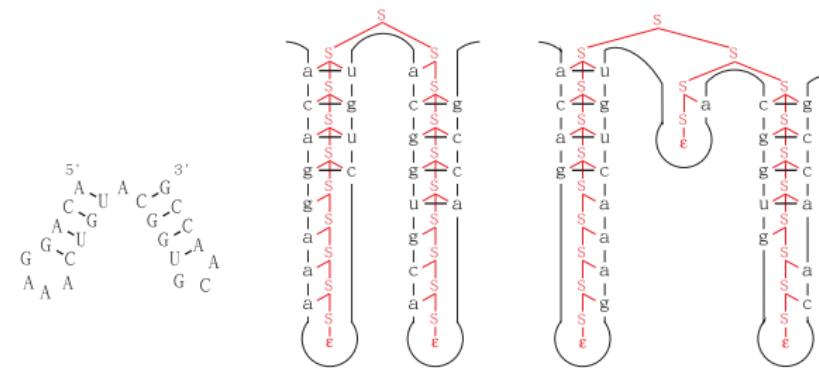
A C C G G U A G U

i=3	0	0	0	1
i=2	0	0	1	1
i=1	0	0	1	1
i=0	1	1	1	2

j = 6 7 8 9

Time complexity: $O(n^3)$.

Extension: RNA is a good example of SCFG.



(see extra slides)

对齐问题.

SEQUENCE ALIGNMENT problem: recursion over **sequence pairs**

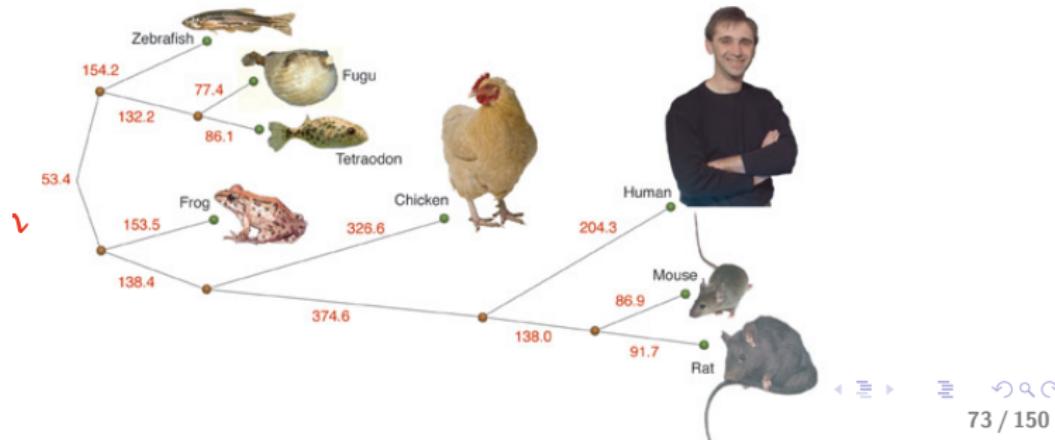
Practical problem: genome similarity

- To identify homology genes of two species, say *Human* and *Mouse*. E.g., Human and Mouse NHPPEs (in KRAS genes) show a high sequence homology (Ref: Cogoi, S., et al. NAR, 2006).

- GGGCGGTGTGGAA-GAGGGAAG-AGGGGGAG
- |||| || | ||||| ||||||| | |||| |
- GGGAGG-GAGGGAAGGAGGGAGGGAGGGAG--

在不同的地方插空格. <字符串不同>
对比核上字符串相似度.

- Having calculating the similarity of genomes of various species, a reasonable phylogeny tree can be estimated (See <https://www.llnl.gov/str/June05/Ovcharenko.html>)



- When you type in “OCURRANCE”, spell tools might guess what you really want to type through the following alignment, i.e. “OCURRANCE” is very similar to “OCCURRENCE” except for INS/DEL/MUTATION operations.
 - O-CURRANCE
 - OCCURRENCE
- But the following instance is a bit difficult:
 - abbbaa-bbbbbaab
 - ababaaaabbba-b

INPUT:

Two sequences S and T , $|S| = m$, and $|T| = n$;

OUTPUT:

To identify an alignment of S and T that maximizes a scoring function.

Note: for the sake of simplicity, the following indexing schema is used: $S = S_1S_2\dots S_m$.

- An example of alignment:
 - O-CURRANCE
 - I | | | | | | |
 - OCCURRENCE
- Basic idea:
 - ① Alignment is usually used to describe the generating process of an erroneous word from the correct word.
 - ② Make the two sequences to have the same length through adding space ' ', i.e. changing S to S' through adding spaces at some positions, and changing T to T' through adding spaces at some positions, too. The only requirement is: $|S'| = |T'|$. There are three cases:
 - ① $T'[i] = ' - '$: $S'[i]$ is simply an INSERTION.
 - ② $S'[i] = ' - '$: $S'[i]$ is simply a DELETION of $T'[i]$.
 - ③ Otherwise, $S'[i]$ is a copy of $T'[i]$ (with possible MUTATION).
 - ③ Thus, an alignment clearly illustrates how to change T into S with a series of INS/DEL/MUTATION operations.

How to measure an alignment in the sense of sequence similarity?

The similarity is defined as the sum of score of aligned letter pairs,
i.e.

$$d(S, T) = \sum_{i=1}^{|S'|} \delta(S'[i], T'[i])$$

目标：最大化分數。

The simplest $\delta(a, b)$ is:

- ① Match: +1 , e.g. $\delta('C', 'C') = 1$.
- ② Mismatch: -1, e.g. $\delta('E', 'A') = -1$.
- ③ Ins/Del: -3, e.g. $\delta('C', '-') = -3$.

打分

3

³Ideally, the score function is designed such that $d(S, T)$ is proportional to $\log \Pr[S \text{ is generated from } T]$. See extra slides for the statistical model for sequence alignment, and better similarity definition, say BLOSUM62, PAM250 substitution matrix, etc.

Alignment is useful

- Observation 1: Using alignment, we can determine the most likely source of “OCURRANCE”.

① $T = \text{“OCCUPATION”}$:

S' :	O	C	_U	RRA	___	NCE					
T' :	O	C	C	U	P	A	T	I	O	N	__

$$d(S', T') = 1 + 1 - 3 + 1 - 3 - 3 - 1 + 1 - 3 - 3 - 3 + 1 - 3 - 3 = -28.$$

② $T = \text{“OCCURRENCE”}$:

S' :	O	C	U	R	R	A	N	C	E	
T' :	O	C	C	U	R	R	E	N	C	E

$$d(S', T') = 1 - 3 + 1 + 1 + 1 + 1 - 1 + 1 + 1 + 1 = 4.$$

- Conjecture: it is more likely that “ocurrance” comes from “occurrence” relative to “occupation”.

- Observation 2: In addition, we can also determine the most likely operations changing “occurrence” into “ocurrance”.

① Alignment 1:

S' :	O_CURRANCE
T' :	OCCURRENCE

“ How? ”
“ Does it? ”

$$d(S', T') = 1 - 3 + 1 + 1 + 1 + 1 - 1 + 1 + 1 + 1 = 4.$$

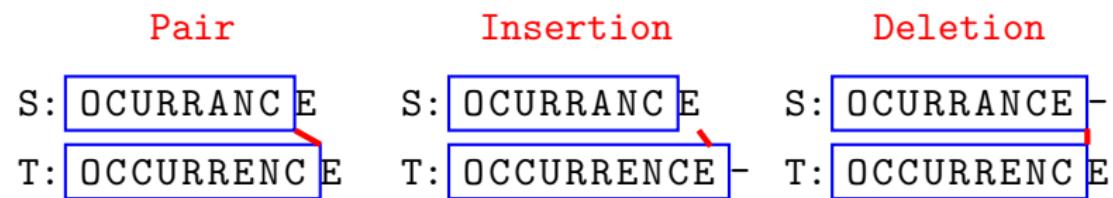
② Alignment 2:

S' :	O_CURR_ANCE
T' :	OCCURRE-NCE

$$d(S', T') = 1 - 3 + 1 + 1 + 1 - 3 - 3 + 1 + 1 + 1 = -1.$$

- Conjecture: the first alignment might describes the real generating process of “ocurrance” from “occurrence”.

- It is not easy to consider long sequences directly. Let's consider whether it is possible to reduce into smaller subproblem.
- Solution: alignment of two sequences. Describe the solving process as a process of multiple-stage **decisions**. At each decision stage, we decide how to align two letters, i.e., how to generate $S[i]$ from $T[j]$.



- Suppose we have already worked out the optimal solution. Consider the **first** decision made for $S[m]$. There are three cases:
 - $S[m]$ pairs with $T[n]$, i.e. $S[m]$ comes from $T[n]$. Then it suffices to align $S[1..m - 1]$ and $T[1..n - 1]$;
 - $S[m]$ pairs with a space ' ', i.e. $S[m]$ is an **INSERTION**. Then we need to align $S[1..m - 1]$ and $T[1..n]$;
 - $T[n]$ pairs with a space ' ', i.e. $S[m]$ is a **DELETION** of a letter in T . Then we need to align $S[1..m]$ and $T[1..n - 1]$.
- In the three cases, the original problem can be reduced into smaller sub-problems.

Pair	Insertion	Deletion
S: OCURRANCE T: OCCURRENCE	S: OCURRANCE T: OCCURRENCE	S: OCURRANCE T: OCCURRENCE
	 ↓	 ↓

子问题

- Thus, we can design the general form of sub-problems as: **前缀对**
alignment a **prefix** of S (denoted as $S[1..i]$) and **prefix** of T (denoted as $T[1..j]$). Denote the optimal solution value as $OPT(i, j)$.
- Optimal substructure property:

$$OPT(i, j) = \max \begin{cases} \delta(S_i, T_j) + OPT(i - 1, j - 1) & i \neq j \\ \delta(' ', T_j) + OPT(i, j - 1) & S_i \text{ 为空} \\ \delta(S_i, ' ') + OPT(i - 1, j) & T_j \text{ 为空} \end{cases}$$

i 和 j 不配
S_i 空
T_j 空

NEEDLEMAN_WUNSCH(S, T)

```
1: for  $i = 0$  to  $m$ ; do
2:    $OPT[i, 0] = -3 * i;$ 
3: end for
4: for  $j = 0$  to  $n$ ; do
5:    $OPT[0, j] = -3 * j;$ 
6: end for
7: for  $i = 1$  to  $m$  do
8:   for  $j = 1$  to  $n$  do
9:      $OPT[i, j] = \max\{OPT[i - 1, j - 1] + \delta(S_i, T_j), OPT[i - 1, j] - 3, OPT[i, j - 1] - 3\};$ 
10:  end for
11: end for
12: return  $OPT[m, n]$  ;
```

Note: the first row is introduced to describe the alignment of prefixes $T[1..i]$ with an empty sequence ϵ , so does the first column.

The first row/column of the alignment score matrix

S:	'	O	C	U	R	R	A	N	C	E
T:	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O										
C	-6									
C	-9									
U	-12									
R	-15									
R	-18									
E	-21									
N	-24									
C	-27									
E	-30									

$$\begin{aligned}
 & S("OC", "O") \\
 \Rightarrow & \left\{ \begin{array}{l} \textcircled{1} S("C", "O") + S("O", "") \\ -1 \\ \textcircled{2} S("C", "") + S("O", "O") \\ \textcircled{3} \end{array} \right.
 \end{aligned}$$

Score:
Alignment:

$d("OCU", "") = -9$
 $S' = OCU$
 $T' = ---$

Score:
Alignment:

$d("", "OC") = -6$
 $S' = --$
 $T' = OC$

Why should we introduce the first row/column?

S:	'	O	C	U	R	R	A	N	C	E	
T:	'	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23	
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19	
U	-9	-5	-1	1	-2	-5	-8	-11	-12	-15	
R	-12	-8	-4	0	0	-3	-6	-9	-12	13	
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11	
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9	
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5	
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4	
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0	
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4	

Score: $d("OC", "O") = \max \begin{cases} d("OC", "") & -3 \\ d("O", "") & -1 \\ d("O", "O") & -3 \end{cases} \quad (-9) \quad (-4) \quad (-2)$

Alignment: $S' = OC$
 $T' = O-$

General cases

S:	'	O	C	U	R	R	A	N	C	E	
T:	'	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23	
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19	
U	-9	-5	-1	1	-2	-5	-8	-11	-12	-15	
R	-12	-8	-4	0	0	-3	-6	-9	-12	13	
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11	
E	-18	-14	-10	-6	-2	2	-	-3	-6	-9	
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5	
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4	
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0	
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4	

Score: $d("OCUR", "OC") = \max \begin{cases} d("OCUR", "O") & -3 \quad (= -11) \\ d("OCU", "O") & -1 \quad (= -6) \\ d("OCU", "OC") & -3 \quad (= -4) \end{cases}$

Alignment: S' = OCUR
T' = OC--

The final entry

S: ''	O	C	U	R	R	A	N	C	E	
T: '	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19
U	-9	-5	-1	1	-2	-5	-8	-11	-12	-15
R	-12	-8	-4	0	0	-3	-6	-9	-12	13
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4

Score: $d("OCURRANCE", "OCCURRENC") = \max \left\{ \begin{array}{ll} d("OCURRANCE", "OCCURRENC") & -3 \quad (-3) \\ d("OCURRANC", "OCCURRENC") & +1 \quad (=4) \\ d("OCURRANC", "OCCURRENCE") & -3 \quad (=3) \end{array} \right.$

Alignment: S= OCURRANCE
T= OCCURRENCE

Question: how to find the alignment with the highest score?

Find the optimal alignment via backtracking

S:	'	O	C	U	R	R	A	N	C	E
T:	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19
C	-9	-5	-1	1	-2	-5	-8	-11	-12	-15
U	-12	-8	-4	0	0	-3	-6	-9	-12	13
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4

Optimal Alignment:
S' = O-CURRANCE
T' = OCCURRENCE

⑥ 目标函数很难设计

⑦ 次优解更易找

→ ⑧ 找一些次优解求

解，然后对其进行聚

类求解

次优解

- It should be noted that in practice, **sub-optimal alignments (as an ensemble)** are more robust than the optimal alignment due to inaccuracy in the scoring model.
- Please refer to Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids for details.

节省空间方案.

Space efficient algorithm: reducing the space requirement from $O(mn)$ to $\underline{O(m + n)}$ (D. S. Hirschberg, 1975)

Technique 1: two arrays are enough if only score is needed

- Key observation 1: it is easy to calculate **the final score** $OPT(S, T)$ **only!**, i.e. **the alignment information are not recorded.**

S: ''	0	C	U	R	R	A	N	C	E	
T: ''	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19
C	-9	-5	-1	1	-2	-5	-8	-11	-12	-15
U	-12	-8	-4	0	0	-3	-6	-9	-12	13
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4

Technique 1: two arrays are enough if only score is needed

- Why? Only column $j - 1$ is needed to calculate column i . Thus, we use two arrays $score[1..m]$ and $newscore[1..m]$ instead of the matrix $OPT[1..m, 1..n]$.

S:	'	O	C	U	R	R	A	N	C	E	
T:	'	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23	
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19	
C	-9	-5	-1	1	-2	-5	-8	-11	-12	-15	
U	-12	-8	-4	0	0	-3	-6	-9	-12	13	
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11	
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9	
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5	
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4	
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0	
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4	

只读两列，
前面部分被丢弃

Technique 1: two arrays are enough if only score is needed

- Why? Only column $j - 1$ is needed to calculate column i . Thus, we use two arrays $score[1..m]$ and $newscore[1..m]$ instead of the matrix $OPT[1..m, 1..n]$.

S:	'	O	C	U	R	R	A	N	C	E	
T:	'	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	-3	1	-2	-5	-8	-11	-14	-17	-20	-23	
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19	
C	-9	-5	-1	1	-2	-5	-8	-11	-12	-15	
U	-12	-8	-4	0	0	-3	-6	-9	-12	13	
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11	
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9	
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5	
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4	
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0	
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4	

Technique 1: two arrays are enough if only score is needed

- Why? Only column $j - 1$ is needed to calculate column i . Thus, we use two arrays $score[1..m]$ and $newscore[1..m]$ instead of the matrix $OPT[1..m, 1..n]$.

S:	'	O	C	U	R	R	A	N	C	E	
T:	'	0	-3	-6	-9	-12	-15	-18	-21	-24	-27
O	0	-3	1	-2	-5	-8	-11	-14	-17	-20	-23
C	-6	-2	2	-1	-4	-7	-10	-13	-16	-19	
C	-9	-5	-1	1	-2	-5	-8	-11	-12	-15	
U	-12	-8	-4	0	0	-3	-6	-9	-12	13	
R	-15	-11	-7	-3	1	1	-2	-5	-8	-11	
R	-18	-14	-10	-6	-2	2	-	-3	-6	-9	
E	-21	-17	-13	-9	-5	-1	1	-1	-4	-5	
N	-24	-20	-16	-12	-8	-4	-2	2	-1	-4	
C	-27	-23	-19	-15	-11	-7	-5	-1	3	0	
E	-30	-26	-22	-18	-14	-10	-8	-4	0	4	

PREFIX-SPACE-EFFICIENT-ALIGNMENT(S, T, SCORE)

```
1: for  $i = 0$  to  $m$  do
2:    $score[i] = -3 * i;$ 
3: end for
4: for  $i = 1$  to  $m$  do
5:    $newscore[0] = 0;$ 
6:   for  $j = 1$  to  $n$  do
7:      $newscore[j] = \max\{score[j - 1] + \delta(S_i, T_j), score[j] -$ 
      $3, newscore[j - 1] - 3\};$ 
8:   end for
9:   for  $j = 1$  to  $n$  do
10:     $score[j] = newscore[j];$ 
11:   end for
12: end for
13: return  $score[n]$  ;
```

Technique 2: aligning suffixes instead of prefixes

- Key observation: Similarly, we can align **suffixes** of S and T instead of **prefixes** and obtain the same score and alignment.

4	0	-4	-10	-12	-16	-18	-22	-26	-30
5	3	-1	-7	-9	-13	-15	-19	-23	-27
3	6	2	-4	-6	-10	-12	-16	-20	-24
-1	2	5	-1	-3	-7	-9	-13	-17	-21
-5	-2	1	4	0	-4	-6	-10	-14	-18
-9	-6	-3	0	3	-1	-3	-7	-11	-15
-13	-10	-7	-4	-1	2	0	-4	-8	-12
-15	-12	-9	-6	-3	0	3	-1	-5	-9
-19	-16	-13	-10	-7	-4	-1	2	-2	-6
-23	-20	-17	-14	-11	-8	-5	-2	1	-3
-27	-24	-21	-18	-15	-12	-9	-6	-3	0

O C U R R A N C E ‘’ S

0 1 2 3 4 5 6 7 8 9

- ① However, **only the recent two columns** of the matrix were kept, the optimal alignment cannot be restored via **backtracking**.
- ② A clever idea: Suppose we have already obtained the optimal alignment. Consider the position **where $S_{[\frac{m}{2}]}$ is aligned to** (denoted as q). We have

$$OPT(S, T) = OPT(S[1..\frac{m}{2}], T[1..q]) + OPT(S[\frac{m}{2}+1..m], T[q+1..n])$$

- ③ Notes:

- Things will be easy as soon as q was determined. The equality holds due to the definition of $d(S, T)$.
- $\frac{m}{2}$ is chosen for the sake of time-complexity analysis.

$$\frac{m}{2}$$

S: OCURRANCE

T: OCCURRENCE

$$1 \leq q \leq n$$

LINEAR-SPACE-ALIGNMENT(S, T)

- 1: Allocate two arrays f and b ; each array has a size of m .
- 2: PREFIX-SPACE-EFFICIENT-ALIGNMENT($S[1.. \frac{m}{2}], T, f$);
- 3: SUFFIX-SPACE-EFFICIENT-ALIGNMENT($S[\frac{m}{2} + 1..m], T, b$);
- 4: Let $q^* = \text{argmax}_q(f[q] + b[q])$;
- 5: Free arrays f and b ;
- 6: Record $\langle \frac{m}{2}, q^* \rangle$ in array A ;
- 7: LINEAR-SPACE-ALIGNMENT($S[1.. \frac{m}{2}], T[1..q^*]$);
- 8: LINEAR-SPACE-ALIGNMENT($S[\frac{m}{2} + 1..m], T[q^* + 1..n]$);
- 9: **return** A ;

$$T(m,n) = T(\frac{m}{2}, q^*) + T(\frac{m}{2}, n - q^*)$$

- Key observation: at each iteration step, only $2n$ space is needed.
- How to determine q ? **Identifying the largest entry in $f[q] + b[q]$.**

Step 1: Determine the optimal aligned position of $S_{[\frac{m}{2}]}$

S: ""		O	C	U	R		
T:	""	0	-3	-6	-9	-12	-24
O		-3	1	-2	-5	-8	-17
C		-6	-2	2	-1	-4	-10
C		-9	-5	-1	1	-2	-5
U		-12	-8	-4	0	0	0
R		-15	-11	-7	-3	1	4
R		-18	-14	-10	-6	-2	-3
E		-21	-17	-13	-9	-5	-8
N		-24	-20	-16	-12	-8	-15
C		-27	-23	-19	-15	-11	-22
E		-30	-26	-22	-18	-14	-29
		<i>f</i>		<i>f+b</i>			
		<i>b</i>					
		R	A	N	C	E	""
							S

The **value** of the largest item: Recall that 4 is actually the **optimal score** of S and T .

Step 2: Recursively solve sub-problems

The **position** of the largest item: Generate two independent sub-problems.

The total space requirement: $O(m + n)$.

- PREFIX_SPACE_EFFICIENT_ALIGNMENT($S[1..\frac{m}{2}], T, f$) needs only $O(n)$ space;
- SUFFIX_SPACE_EFFICIENT_ALIGNMENT($S[\frac{m}{2} + 1..m], T, b$) needs only $O(n)$ space;
- Line 4 (Record $< \frac{n}{2}, q^* >$ in array A) needs only $O(m)$ space;

Theorem

Algorithm LINEAR-SPACE-ALIGNMENT(S, T) still takes $O(mn)$ time.

Proof.

- The algorithm implies the following recursion:

$$T(m, n) = cmn + T\left(q, \frac{n}{2}\right) + T\left(m - q, \frac{n}{2}\right);$$

- Difficulty: we have no idea of q before algorithm ends; thus, the master theorem cannot apply directly. **Guess and substitution!!!**
- Guess: $T(m', n') \leq km'n'$ follows for any $m' < m$ and $n' < n$.
- Substitution:

$$T(m, n) = cmn + T\left(q, \frac{n}{2}\right) + T\left(m - q, \frac{n}{2}\right) \quad (7)$$

$$\leq cmn + kq\frac{n}{2} + k(m - q)\frac{n}{2} \quad (8)$$

$$= cmn + kq\frac{n}{2} + km\frac{n}{2} - kq\frac{n}{2} \quad (9)$$

$$\leq \left(c + \frac{k}{2}\right)mn \quad (10)$$

$$= kmn \quad (\text{set } k = 2c) \quad (11)$$



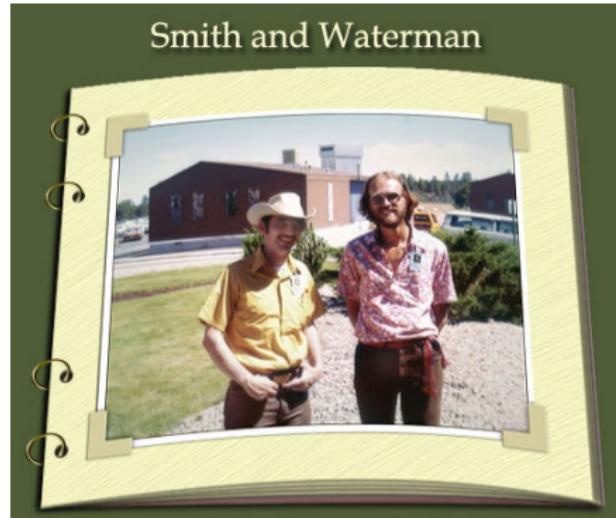
Extended Reading 1: From global alignment to local alignment

From Global alignment to Local alignment: Smith-Waterman algorithm

- Global alignment: to identify similarity between **two whole sequences**;
- Local alignment: It is often that we wish to find **similar SEGMENTS (sub-sequences)**.

Smith-Waterman algorithm [1981]

- Needleman-Wunsch **global alignment** algorithm was developed by biologists in 1970s, about twenty years later than Bellman-Ford algorithm was developed.
- Then Smith-Waterman **local alignment** algorithm was proposed.



Please refer to Smith and Waterman1981 for details.

Extended Reading 2: How to derive a reasonable scoring schema?

$$\text{Opt}(i, j) = \max \left\{ \begin{array}{l} \delta(s_i, T_j) + \text{Opt}(i-1, j-1) \\ \delta(' - ', T_j) + \text{Opt}(i, j-1) \\ \delta(s_i, ' - ') + \text{Opt}(i-1, j) \end{array} \right.$$

PAM250: one of the most popular substitution matrices in Bioinformatics

分子进化的距离：

A_i : OCRR
O - RR

B : THE
TEH

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	B	Z	X	*
A	2	-2	0	0	-2	0	0	1	-1	-1	-2	-1	-1	-3	1	1	1	-6	-3	0	0	0	0	-8
R	-2	6	0	-1	-4	1	-1	-3	2	-2	-3	3	0	-4	0	0	-1	2	-4	-2	-1	0	-1	-8
N	0	0	2	2	-4	1	1	0	2	-2	-3	1	-2	-3	0	1	0	-4	-2	-2	2	1	0	-8
D	0	-1	2	4	-5	2	3	1	1	-2	-4	0	-3	-6	-1	0	0	-7	-4	-2	3	3	-1	-8
C	-2	-4	-4	-5	12	-5	-5	-3	-3	-2	-6	-5	-5	-4	-3	0	-2	-8	0	-2	-4	-5	-3	-8
Q	0	1	1	2	-5	4	2	-1	3	-2	-2	1	-1	-5	0	-1	-1	-5	-4	-2	1	3	-1	-8
E	0	-1	1	3	-5	2	4	0	1	-2	-3	0	-2	-5	-1	0	0	-7	-4	-2	3	3	-1	-8
G	1	-3	0	1	-3	-1	0	5	-2	-3	-4	-2	-3	-5	0	1	0	-7	-5	-1	0	0	-1	-8
H	-1	2	2	1	3	3	1	-2	6	-2	-2	0	-2	-2	0	-1	-1	-3	0	-2	1	2	-1	-8
I	-1	-2	-2	-2	-2	-2	-3	-2	5	2	-2	2	1	-2	-1	0	-5	-1	4	-2	-2	-1	-8	
L	-2	-3	-3	-4	-6	-2	-3	-4	-2	2	6	-3	4	2	-3	-3	-2	-2	-1	2	-3	-3	-1	-8
K	-1	3	1	0	-5	1	0	-2	0	-2	-3	5	0	-5	-1	0	0	-3	-4	-2	1	0	-1	-8
M	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6	0	-2	-2	-1	-4	-2	2	-2	-2	-1	-8
F	-3	-4	-3	-6	-4	-5	-5	-5	-2	1	2	-5	0	9	-5	-3	-3	0	7	-1	-4	-5	-2	-8
P	1	0	0	-1	-3	0	-1	0	0	-2	-3	-1	-2	-5	6	1	0	-6	-5	-1	-1	0	-1	-8
S	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	2	1	-2	-3	-1	0	0	0	-8
T	1	-1	0	0	-2	-1	0	0	-1	0	-2	0	-1	-3	0	1	3	-5	-3	0	0	-1	0	-8
W	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	0	-6	-2	-5	17	0	-6	-5	-6	-4	-8
Y	-3	-4	-2	-4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10	-2	-3	-4	-2	-8
V	0	-2	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	-1	0	-6	-2	4	-2	-2	-1	-8
B	0	-1	2	3	-4	1	3	0	1	-2	-3	1	-2	-4	-1	0	0	-5	-3	-2	3	2	-1	-8
Z	0	0	1	3	-5	3	3	0	2	-2	-3	0	-2	-5	0	0	-1	-6	-4	-2	2	3	-1	-8
X	0	-1	0	-1	-3	-1	-1	-1	-1	-1	-1	-1	-2	-1	0	0	-4	-2	-1	-1	-1	-1	-8	
*	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	

$$d(0,0) = \log \frac{P_{00}}{q_{00} q_{00}}$$

$$d(0,-) = ?$$

max likelihood.

$\max L(A_i, \theta)$

Please refer to "PAM matrix for Blast algorithm" (by C. Alexander, 2002) for the details to calculate PAM matrix.

Extended Reading 3: How to measure the significance of an alignment?

- When two random sequences of length m and n are compared, the probability of finding a pair of segments with a score greater than or equal to S is $1 - e^{-y}$, where $y = Kmne^{-\lambda S}$.

Please refer to Altschul1990 for details.

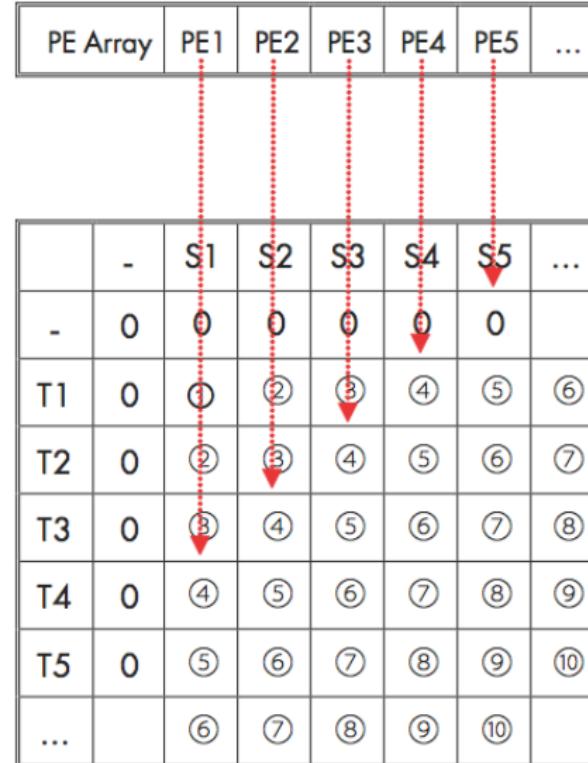
Extended Reading 4: An FPGA implementation of
Smith-Waterman algorithm

The potential parallelity of SmithWaterman algorithm

	-	S1	S2	S3	S4	S5	...
-	0	0	0	0	0	0	
T1	0	①	②	③	④	⑤	⑥
T2	0	②	③	④	⑤	⑥	⑦
T3	0	③	④	⑤	⑥	⑦	⑧
T4	0	④	⑤	⑥	⑦	⑧	⑨
T5	0	⑤	⑥	⑦	⑧	⑨	⑩
...		⑥	⑦	⑧	⑨	⑩	

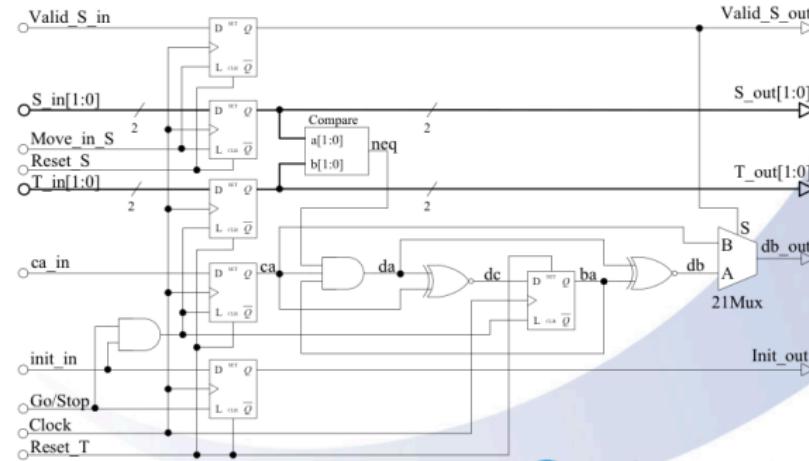
For example, in the first cycle, only one element marked as (1) could be calculated. In the second cycle, two elements marked as (2) could be calculated. In the third cycle, three elements marked as (3) could be calculated, etc., and this feature implies that the algorithm has a very good potential parallelity.

Mapping Smithg-Waterman algorithm on PE



See *Implementation of the Smith-Waterman Algorithm on a Reconfigurable Supercomputing Platform* for details.

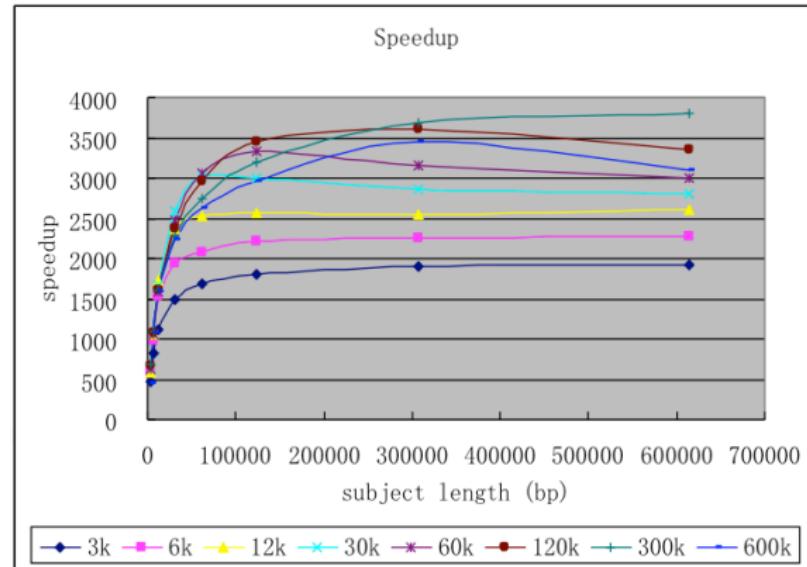
PE design of a card for Dawning 4000L



Smith-Waterman card for Dawning 4000L



Performance of Dawning 4000L



4

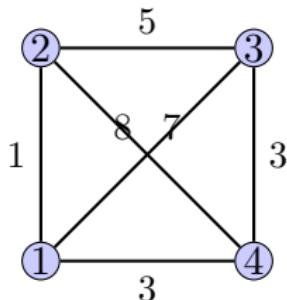
⁴Some pictures were excerpted from *Introduction to algorithms*

BELLMAN-HELD-KARP for TSP problem: recursion over **graphs**

TRAVELLING SALESMAN PROBLEM

INPUT: a list of n cities (denoted as V), and the distances between each pair of cities d_{ij} ($1 \leq i, j \leq n$);

OUTPUT: the shortest tour that visits each city exactly once and returns to the origin city

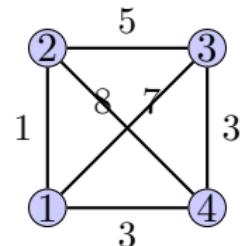


#Tours: 6

- Tour 1: $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 1$ (12)
- Tour 2: $1 \rightarrow 2 \rightarrow 4 \rightarrow 3 \rightarrow 1$ (21)
- Tour 3: $1 \rightarrow 3 \rightarrow 2 \rightarrow 4 \rightarrow 1$ (23)
-

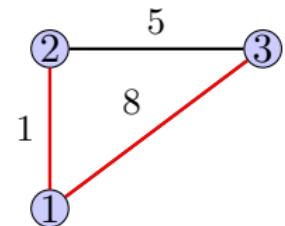
Definition

$D(S, e) =$ the minimum distance, starting from city 1, visiting all cities in S , and finishing at city $e \in S$.



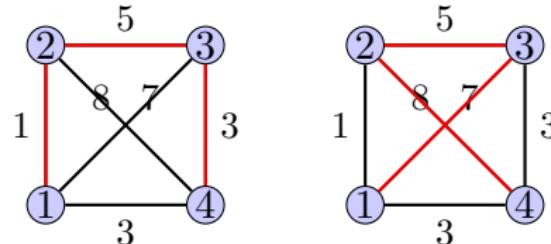
- It suffices to calculate $D(S, e)$ for any $S \in \{1, 2, \dots, n\}$ and city e since:
 - There are 3 cases of the city from which we return to 1.
 - Thus, the shortest tour can be calculated as:
$$\min\{D(\{1, 2, 3, 4\}, 2) + d_{2,1}, D(\{1, 2, 3, 4\}, 3) + d_{3,1}, D(\{1, 2, 3, 4\}, 4) + d_{4,1}\}$$

- It is trivial to calculate $D(S, e)$ when S consists of only 1 cities.



- $D(\{2\}, 2) = d_{12}$;
- $D(\{3\}, 3) = d_{13}$;
- But how to solve a large problem, say $D(\{2, 3, 4\}, 4)$?

Divide a larger problem into smaller problems



- $D(\{2, 3, 4\}, 4) = \min\{D(\{2, 3\}, 3) + d_{34}, D(\{2, 3\}, 2) + d_{24}\};$
- Optimal substructure property:

$$D(S, e) = \begin{cases} d_{1e} & \text{if } S = \{e\} \\ \min_{m \in S - \{e\}} (D(S - \{e\}, m) + d_{me}) & \text{otherwise} \end{cases}$$

```
function  $D(S, e)$ 
1: if  $S = \{e\}$  then
2:   return  $d_{1e}$ ;
3: end if
4:  $d = \infty$ ;
5: for all city  $m \in S$ , and  $m \neq e$  do
6:   if  $D(S - \{e\}, m) + d_{me} < d$  then
7:      $d = D(S - \{e\}, m) + d_{me}$ ;
8:   end if
9: end for
10: return  $d$ ;
```

- Space complexity: $\sum_{k=2}^{n-1} k \binom{n-1}{k} + n - 1 = (n-1)2^{n-2}$
- Time complexity: $\sum_{k=2}^{n-1} k(k-1) \binom{n-1}{k} + n - 1 = O(2^n n^2)$.

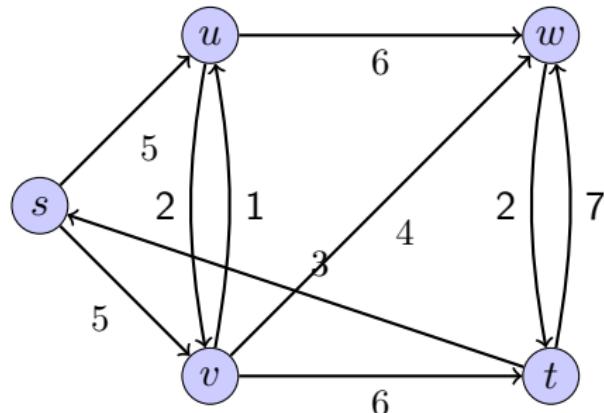
SINGLESOURCESHORTESTPATH problem: recursion over **graphs**

INPUT:

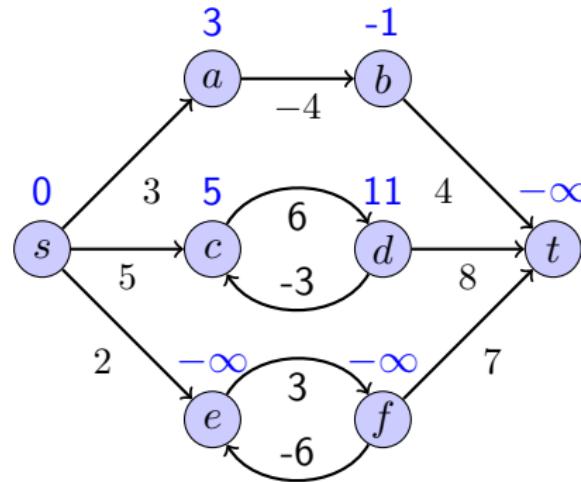
A directed graph $G = \langle V, E \rangle$. Each edge $e = \langle i, j \rangle$ has a weight or distance $d(i, j)$. Two special nodes: source s , and destination t ;

OUTPUT:

A shortest path from s to t ; that is, the sum weight of the edges is minimized.



SHORTESTPATH problem: cycles



- Here $d(i, j)$ might be negative; however, there should be **no negative cycle**, i.e. the sum weight of edges in any cycle should be greater than 0.
- In fact, a negative cycle means an $-\infty$ shortest-path weight. Since e and f form a negative-weight cycle reachable from s , they have shortest-path weight of $-\infty$ from s .

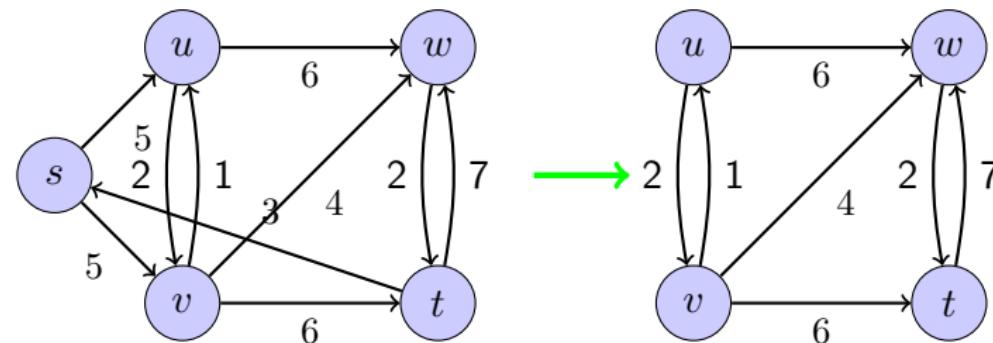
Trial 1: describing the sub-problem as finding the shortest path in
a graph

- Solution: a path. Describe the solving process as series of decisions. At each decision stage, we need to **determine an edge** to the subsequent node.
- Suppose we have already worked out the optimal solution O .
- Consider **the first decision** in O . The options are:
 - All edges starting from s : Suppose we use an edge $\langle s, v \rangle$. Then it suffices to calculate the shortest path in graph $G' = \langle V', E' \rangle$, where node s and related edges are removed.
 - Thus, the original problem can be reduced into smaller sub-problems.
 - General form of sub-problem: to find the shortest path from node v to t in graph G .

多与决策

- General form of sub-problem: to find the shortest path from node v to t in graph G . Denote the optimal solution value as $OPT(G, v)$.
 - Optimal substructure:

$$OPT(G, s) = \min_{v: \langle s, v \rangle \in E} \{OPT(G', v) + d(s, v)\}$$



- **Infeasible!** The number of sub-problems is **exponential**.

Trial 2: another problem form with a new variable

Trial 2: simplifying sub-problem form via limiting path length

$\text{opt}(v, t, k)$: 框图

k 条边的最短路径.



$$\min \begin{cases} \text{opt}(v, t, k) = \min(d(v, w) + \text{opt}(w, t, k-1)), \forall w \\ \text{OPT}(v, t, k-1) \\ \text{OPT}(v, t, k-2) \\ \vdots \end{cases}$$

- Solution: the shortest path from node s to t is a path with **at most n nodes** (Why? no negative cycle \Rightarrow removing cycles in a path can shorten the path). Describe the solving process as a process of multiple-stage **decisions**; at each decision stage, we decide **the subsequent node** from current node.

- Suppose we have already worked out the optimal solution O .
- Consider the first decision in O . The feasible options are:
 - All adjacent nodes of s : Suppose we choose an edge $< s, v >$ to node v . Then the left-over is to find the shortest path from v to t via at most $n - 2$ edges.
- Thus the general form of subproblem can be designed as: to find the shortest path from node v to t with **at most k edges** ($k \leq n - 1$). Denote the optimal solution value as $OPT(v, t, k)$.
- Optimal substructure: 最多 k 条边

$$OPT[v, t, k] = \min \left\{ OPT[v, t, k-1], \min_{<v, w> \in E} \{ OPT[w, t, k-1] + d(v, w) \} \right\}$$

从 v 到 t 经过最多
 k 条边.

$$\begin{array}{l} V^{\max} \\ T^{\max} \\ K^{\max} \end{array} \Rightarrow \text{问题规模是 } n^3$$

Bellman-Ford algorithm [1956, 1958]

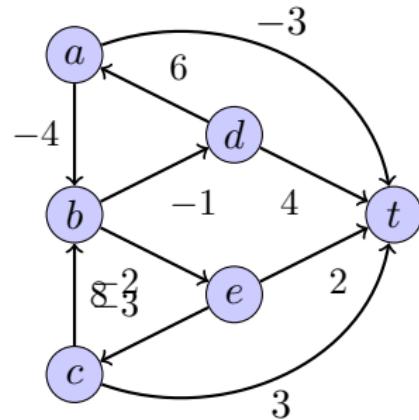
```
BELLMAN-FORD( $G, s, t$ )
1: for any node  $v \in V$  do
2:    $OPT[v, t, 0] = \infty$ ;
3: end for
4: for  $k = 0$  to  $n - 1$  do
5:    $OPT[t, t, k] = 0$ ;
6: end for
7: for  $k = 1$  to  $n - 1$  do
8:   for all node  $v$  (in an arbitrary order) do
9:      $OPT[v, t, k] =$ 
         $\min \left\{ OPT[v, t, k - 1], \min_{v, w \in E} \{ OPT[w, t, k - 1] + d(v, w) \} \right\}$ 
10:  end for
11: end for
12: return  $OPT[s, t, n - 1]$ ;
```

Note that the algorithm actually finds the shortest path from every possible source to t (or from s to every possible destination) by constructing a **shortest path tree**.



See "Richard Bellman on the birth of dynamic programming" (S. Dreyfus, 2002) and "On the routing problem" (R. Bellman, 1958) for details.

An example

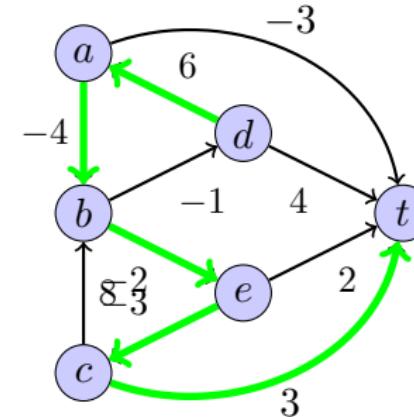


Source node	$k = 0$	$k = 1$	$k = 2$	$k = 3$	$k = 4$	$k = 5$
t	0	0	0	0	0	0
a	-	-3	-3	-4	-6	-6
b	-	-	0	-2	-2	-2
c	-	3	3	3	3	3
d	-	4	3	3	2	0
e	-	2	0	0	0	0

最短路径树

Shortest path tree

杂居：①多少个问题
②每个问题花多少时间。



Source node	$k = 0$	$k = 1$	$k = 2$	$k = 3$	$k = 4$	$k = 5$
t	0	0	0	0	0	0
a	-	-3	-3	-4	-6	-6
b	-	-	0	-2	-2	-2
c	-	3	3	3	3	3
d	-	4	3	3	2	0
e	-	2	0	0	0	0

Note: the shortest paths from all nodes to t form a *shortest path tree*.

- ① Cursory analysis: $O(n^3)$. (There are n^2 subproblems, and for each subproblem, we need at most $O(n)$ operations in line 7.)
- ② Better analysis: $O(mn)$. (Efficient for sparse graph, i.e. $m << n^2$.)
 - For each node v , line 7 need $O(d_v)$ operations, where d_v denotes the degree of node v ;
 - Thus **the inner for loop** (lines 6-8) needs $\sum_v d_v = O(m)$ operations;
 - Thus **the outer for loop** (lines 5-9) needs $O(nm)$ operations.

Extension: detecting negative cycle

Theorem

If t is reachable from node v , and v is contained in a negative cycle, then we have: $\lim_{k \rightarrow \infty} OPT(v, t, k) = -\infty$.

负圈检测

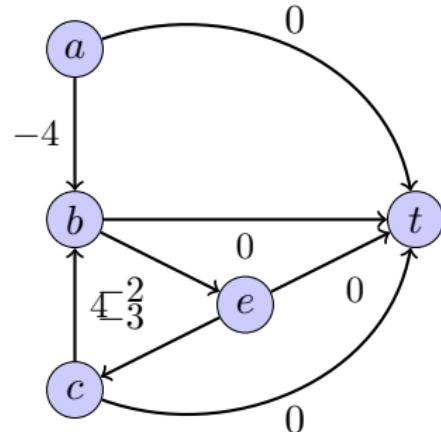
和每个边加一个正值
不等式.

例: $t \Rightarrow +100$.

$b \rightarrow t: 00$ | 0

$b \rightarrow e \rightarrow t: 198$ | 2

增量不同.



Intuition: a traveling of the negative cycle leads to a shorter length. Say,

$$\text{length}(b \rightarrow t) = 0$$

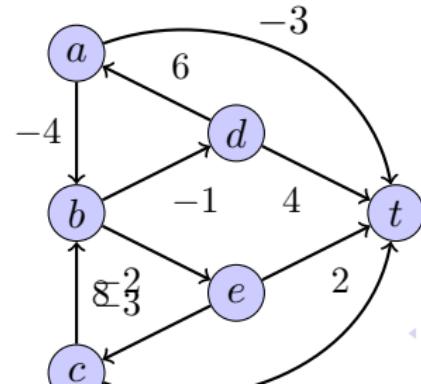
$$\text{length}(b \rightarrow e \rightarrow c \rightarrow b \rightarrow t) = -1$$

$$\text{length}(b \rightarrow e \rightarrow c \rightarrow b \rightarrow e \rightarrow c \rightarrow b \rightarrow t) = -2$$

.....

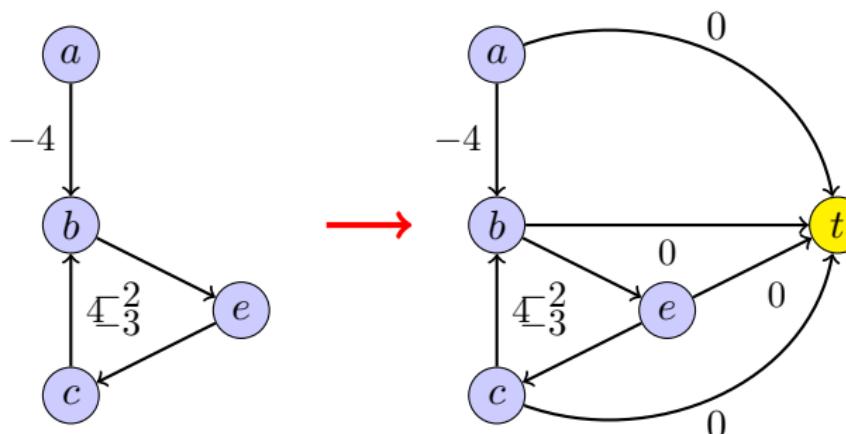
Corollary

If there is no negative cycle in G , then for all node v , and $k \geq n$, $OPT(v, t, k) = OPT(v, t, n)$.

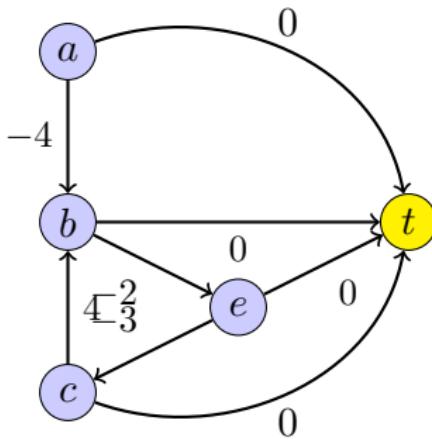


Detecting negative cycle via adding edges and a node t

- Expanding G to G' to guarantee that t is reachable from the negative cycle:
 - ① Adding a new node t ;
 - ② For each node v , adding a new edge $\langle v, t \rangle$ with $d(v, t) = 0$;
- Property: G has a negative cycle C , say, $b \rightarrow e \rightarrow c \rightarrow b \Rightarrow t$ is reachable from a node in C . Thus, the above theorem applies.



An example of negative cycle

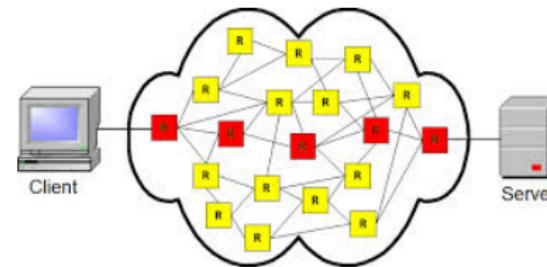


Source node	k=0	k=1	k=2	k=3	k=4	k=5	k=6	k=7	k=8	...
t	0	0	0	0	0	0	0	0	0	...
a	-	0	-4	-6	-9	-9	-11	-11	-12	...
b	-	0	-2	-5	-5	-7	-7	-8	-8	...
c	-	0	0	0	-2	-3	-3	-3	-4	...
e	-	0	-3	-3	-5	-5	-6	-6	-6	...

Application of Bellman-Ford algorithm: Internet router protocol

Problem statement:

- Each node denotes a route, and the weight denotes the **transmission delay** of the link from router i to j .
- The objective to design a protocol to determine the quickest route when router s wants to send a package to t .



Internet router protocol: Dijkstra's algo vs. Bellman-Ford algo

- Choice: Dijkstra algorithm.
- However, the algorithm needs **global knowledge**, i.e. the knowledge of the whole graph, which is (almost) impossible to obtain.
- In contrast, the Bellman-Ford algorithm **needs only local information**, i.e. the information of its **neighbourhood** rather than **the whole network**.

ASYNCHRONOUSSHORTESTPATH(G, s, t)

- 1: Initially, set $OPT[t, t] = 0$, and $OPT[v, t] = \infty$;
- 2: Label node t as "active";
- 3: **while** exists an active node **do**
- 4: arbitrarily select an active node w ;
- 5: remove w 's active label;
- 6: **for all** edges $< v, w >$ (in an arbitrary order) **do**
- 7: $OPT[v, t] = \min \begin{cases} OPT[v, t] \\ OPT[w, t] + d(v, w) \end{cases}$
- 8: **if** $OPT[v, t]$ was updated **then**
- 9: label v as "active";
- 10: **end if**
- 11: **end for**
- 12: **end while**

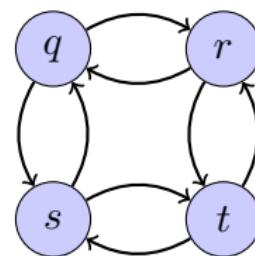
A related problem: LONGESTPATH problem

INPUT:

A directed graph $G = \langle V, E \rangle$. Each edge (u, v) has a distance $d(u, v)$. Two nodes s and t ;

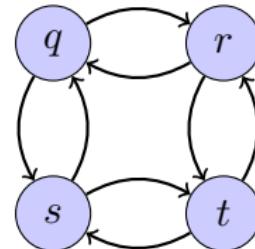
OUTPUT:

The longest simple path from s to t ;



Hardness: LONGESTPATH problem is NP-hard. (Hint: it is obvious that LONGESTPATH problem contains *Hamilton* path as its special case.)

- Divide: Wrong! The subproblems are not **independent**.
- Consider dividing problem to find a path from q to t into two subproblems: to find a path from q to r , and to find a path from r to t .



- Suppose we have already solved the sub-problems. Let's try to combine the solutions to the two sub-problems:

- ① $P(q, r) = q \rightarrow s \rightarrow t \rightarrow r$
 - ② $P(r, t) = r \rightarrow q \rightarrow s \rightarrow t$,
- 子问题不独立.

We will obtain a path $q \rightarrow s \rightarrow t \rightarrow r \rightarrow q \rightarrow s \rightarrow t$, which is not simple.

- In other words, the use of s in the first subproblem prevents us from using s in the second subproblem. However, we cannot obtain the optimal solution to the second subproblem without using s .

- In contrast, the SHORTESTPATH does not have this difficulty.
- Why? The solutions to the subproblems share no node. *无公共节点,*
Suppose the shortest paths $P(q, r)$ and $P(r, t)$ share a node $w(w \neq r)$. Then there will be a cycle
 $w \rightarrow \dots \rightarrow r \rightarrow \dots \rightarrow w$. Removing this cycle leads to a shorter path (no negative cycle). A contradiction.
- This means that the two subproblems are independent: the solution of one subproblem does not affect the solution to another subproblem.

A greedy algorithm exists when posing a stricter limit, i.e., all edges have a positive weight.

We will talk about this in next lectures.