

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 31**

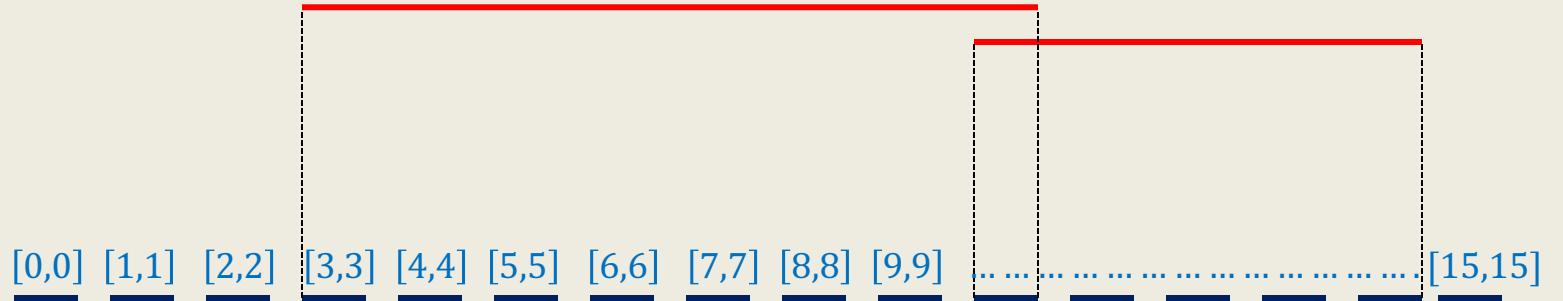
- **Magical applications of Binary trees -II**

# **RECAP OF LAST LECTURE**

# Intervals

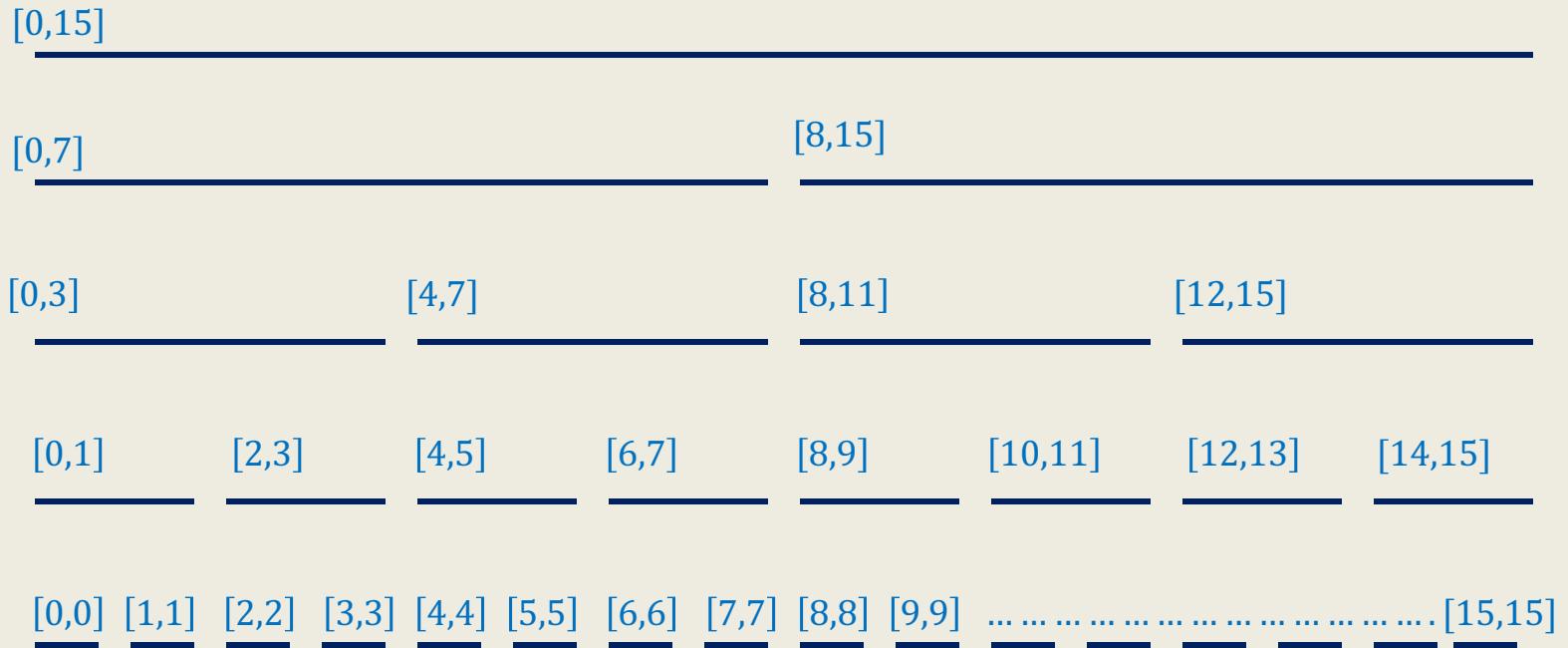
$$\mathbf{S} = \{[\mathbf{i}, \mathbf{j}], 0 \leq i \leq j < n\}$$

**Question:** Can we have a small set  $\mathbf{X} \subset \mathbf{S}$  of **intervals** s.t.  
every interval in  $\mathbf{S}$  can be expressed as a union of a few intervals from  $\mathbf{X}$  ?

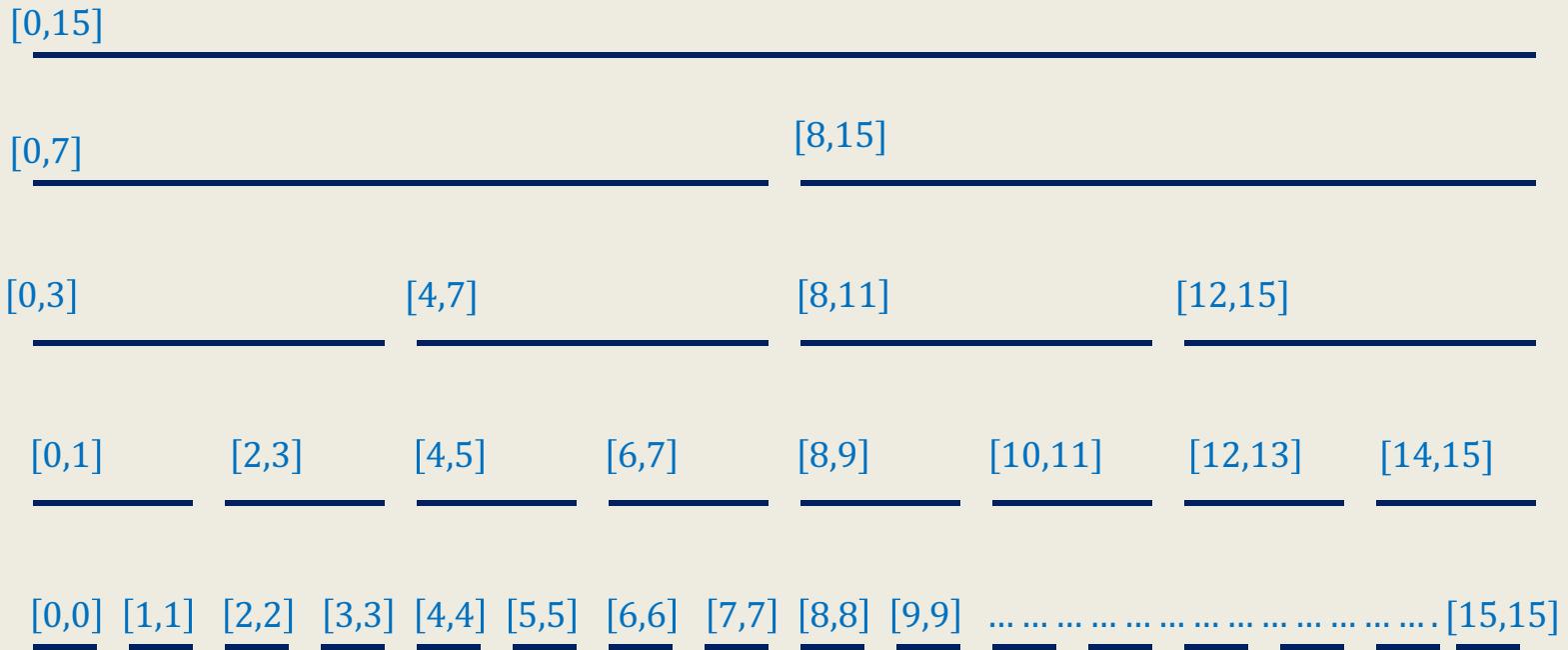


**Answer:** yes😊

# Hierarchy of intervals

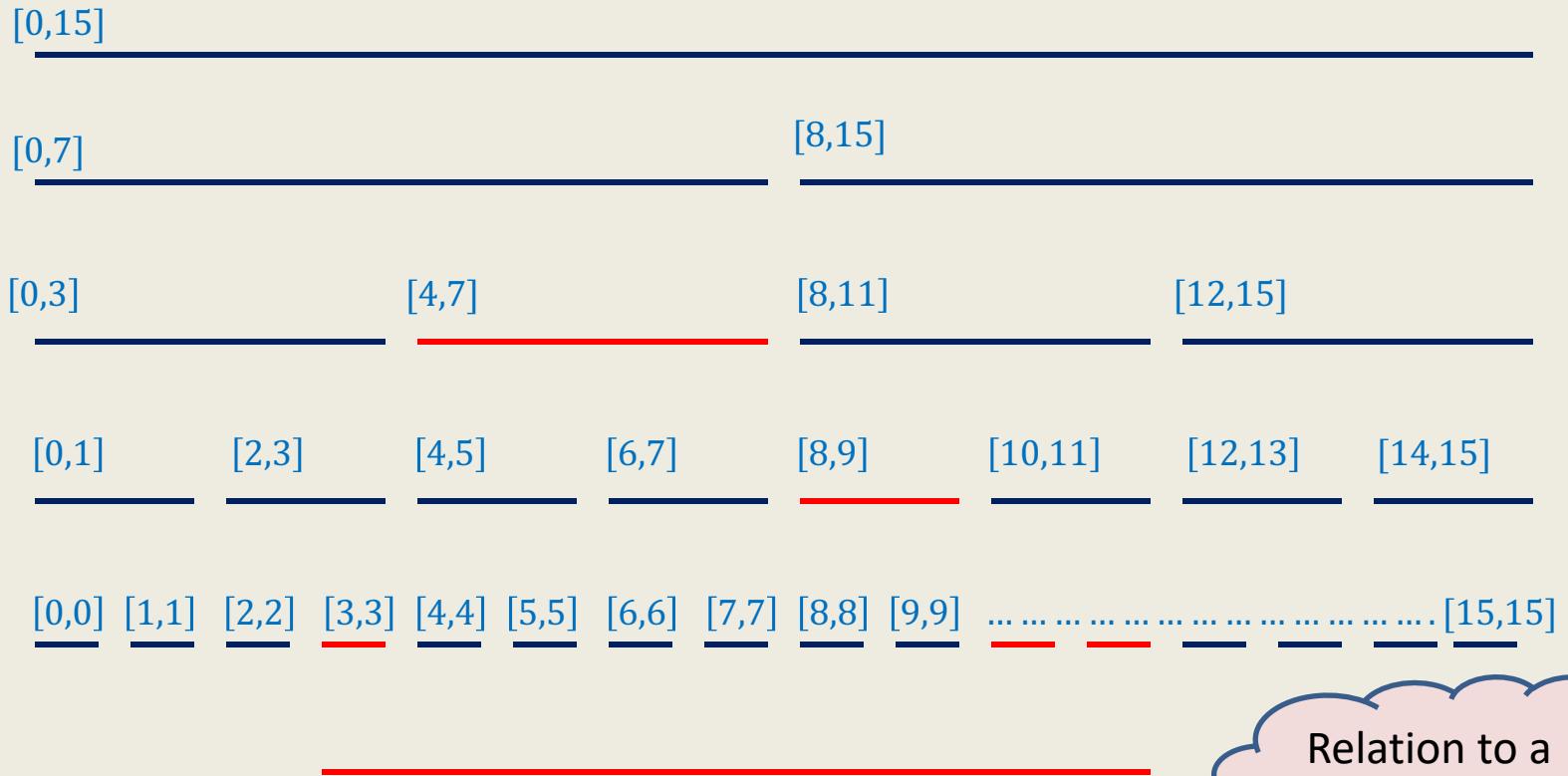


# Hierarchy of intervals



**Observation:** There are  $2n$  intervals such that  
any interval  $[i, j]$  can be expressed as **union** of  $O(\log n)$  basic intervals ☺

# Hierarchy of intervals



**Observation:** There are  $2n$  intervals such that

any interval  $[i, j]$  can be expressed as union of  $O(\log n)$  basic intervals 😊

Relation to a  
sequence ?

# Which data structure emerges ?

[0,15]

[0,7]

[8,15]

[0,3]

[4,7]

[8,11]

[12,15]

[0,1]

[2,3]

[4,5]

[6,7]

[8,9]

[10,11]

[12,13]

[14,15]

[0,0]

[1,1]

[2,2]

[3,3]

[4,4]

[5,5]

[6,6]

[7,7]

[8,8]

[9,9]

.....

.....

.....

.....

[15,15]

Sequence

$x_0$

$x_1$

$x_2$

$x_3$

$x_4$

$x_5$

$x_6$

$x_7$

$x_8$

$x_9$

$x_{10}$

$x_{11}$

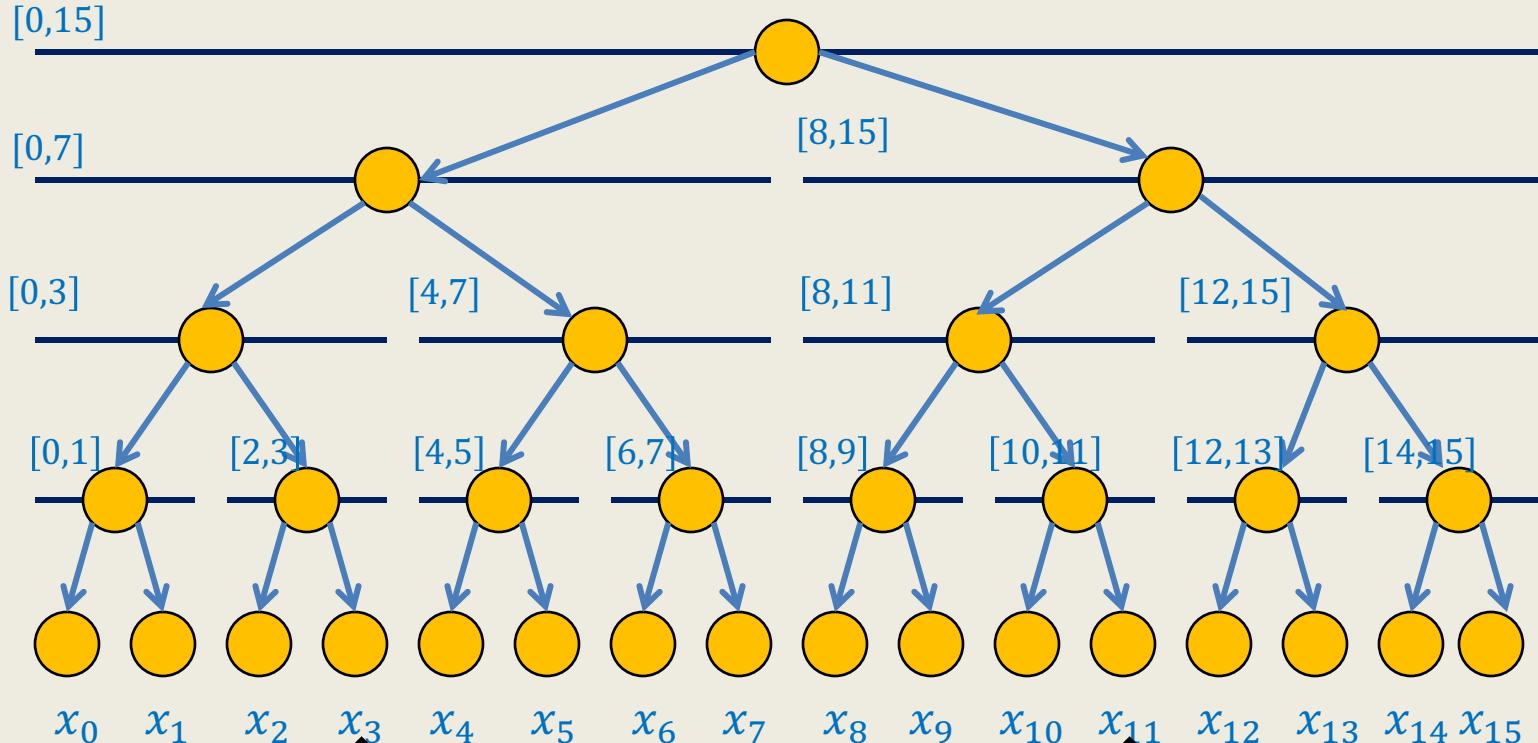
$x_{12}$

$x_{13}$

$x_{14}$

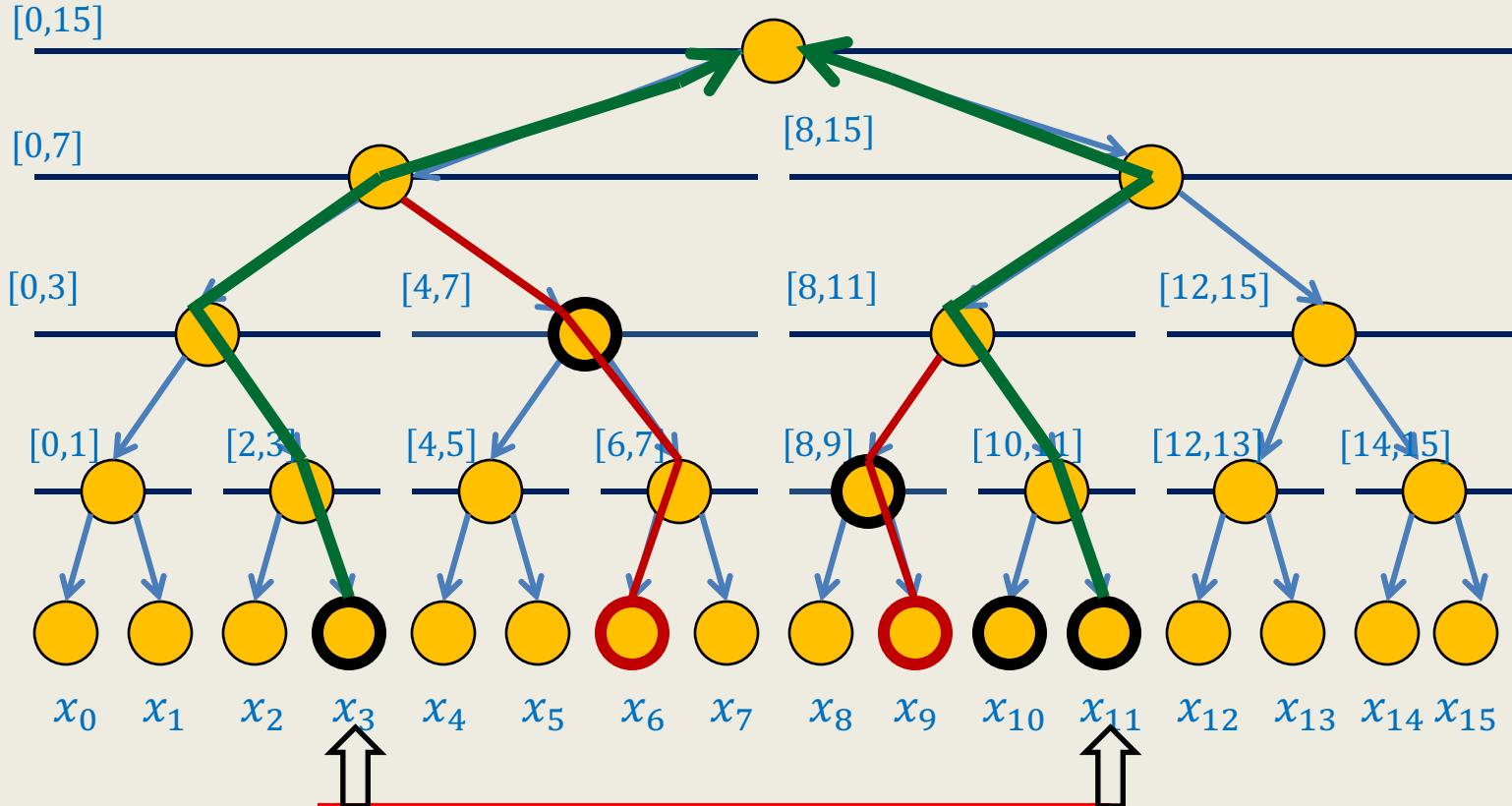
$x_{15}$

## A Binary tree



What value should you keep  
in internal nodes ?

How to perform Operation on an interval ?



**How to perform Operation on an interval ?**

# **Problem 2**

**Dynamic Range-minima**

# Dynamic Range Minima Problem

Given an initial sequence  $S = \langle x_0, \dots, x_{n-1} \rangle$  of numbers, maintain a compact data structure to perform the following operations efficiently for any  $0 \leq i < j < n$ .

- **ReportMin( $i, j$ ):**  
Report the minimum element from  $\{x_k \mid \text{for each } i \leq k \leq j\}$
  - **Update( $i, a$ ):**  
 $a$  becomes the new value of  $x_i$ .
- 

## Example:

Let the initial sequence be  $S = \langle 14, 12, 3, 49, 4, 21, 322, -40 \rangle$

**ReportMin(1, 5)** returns **3**

**ReportMin(0, 3)** returns **3**

**Update(2, 19)** update  $S$  to  $\langle 14, 12, 19, 49, 4, 21, 322, -40 \rangle$

**ReportMin(1, 5)** returns **4**

**ReportMin(0, 3)** returns **12**

# Dynamic Range Minima Problem

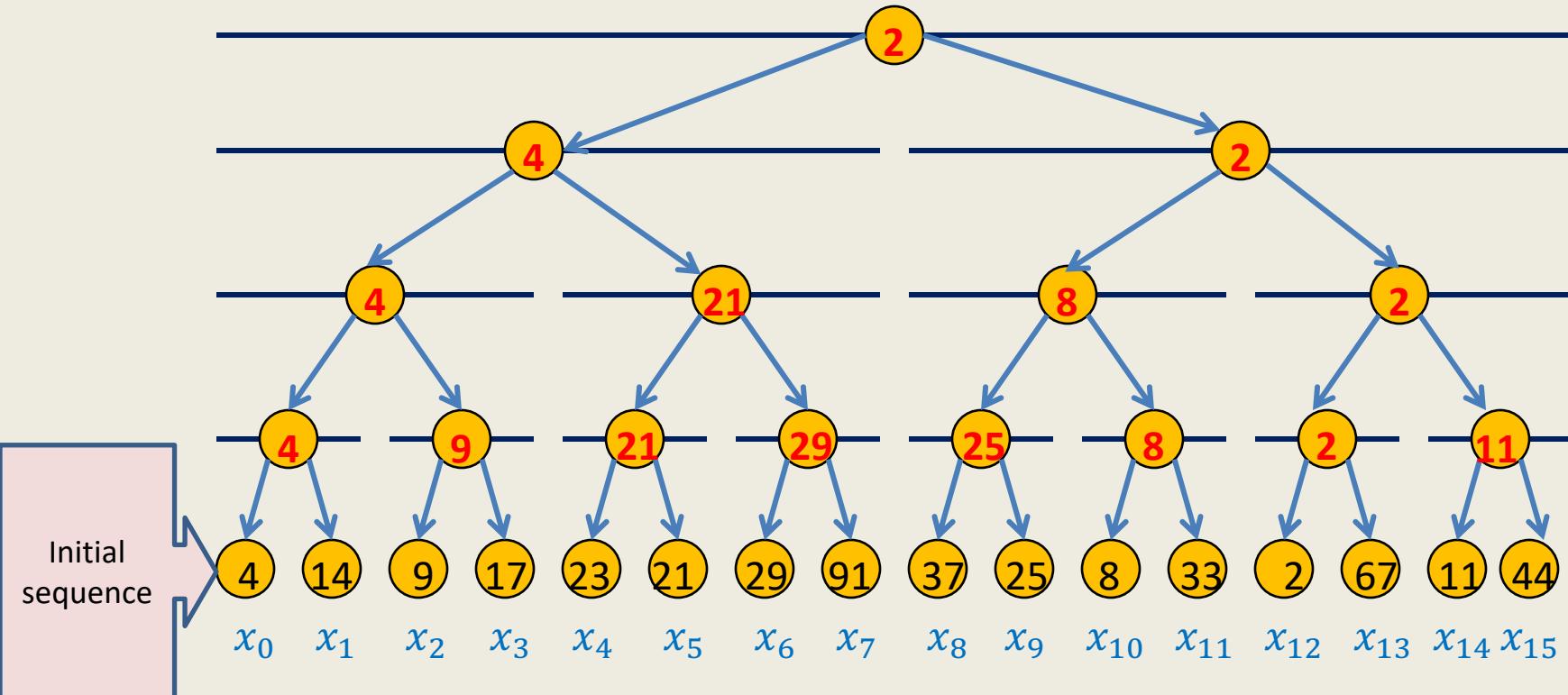
Given an initial sequence  $S = \langle x_0, \dots, x_{n-1} \rangle$  of numbers, maintain a compact data structure to perform the following operations efficiently for any  $0 \leq i < j < n$ .

- **ReportMin( $i, j$ ):**  
Report the minimum element from  $\{x_k \mid \text{for each } i \leq k \leq j\}$
  - **Update( $i, a$ ):**  
 $a$  becomes the new value of  $x_i$ .
- 

## AIM:

- $O(n)$  size data structure.
- ReportMin( $i, j$ ) in  $O(\log n)$  time.
- Update( $i, a$ ) in  $O(\log n)$  time.

# Data structure for dynamic range minima

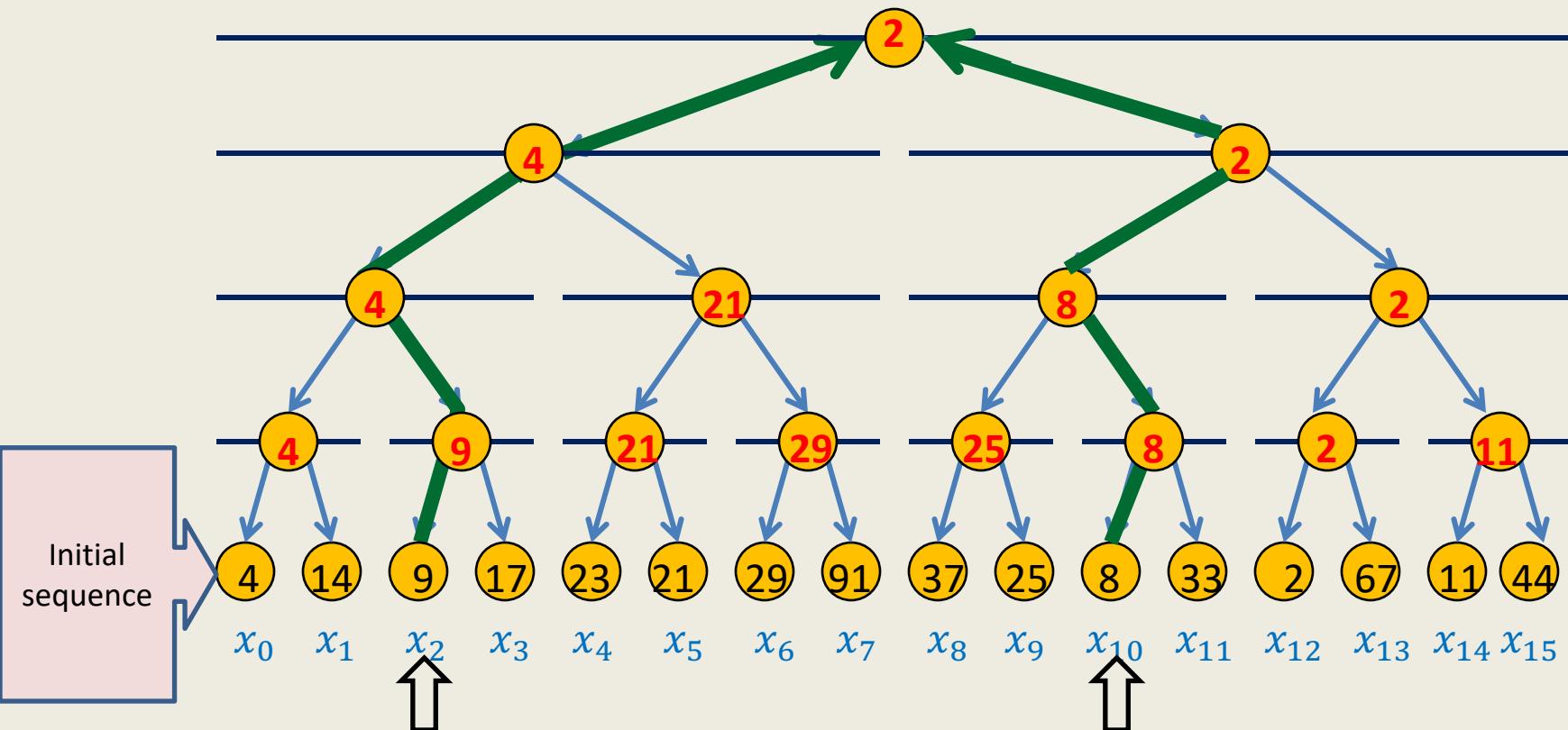


**Question:** What should be stored in an internal node  $v$  ?

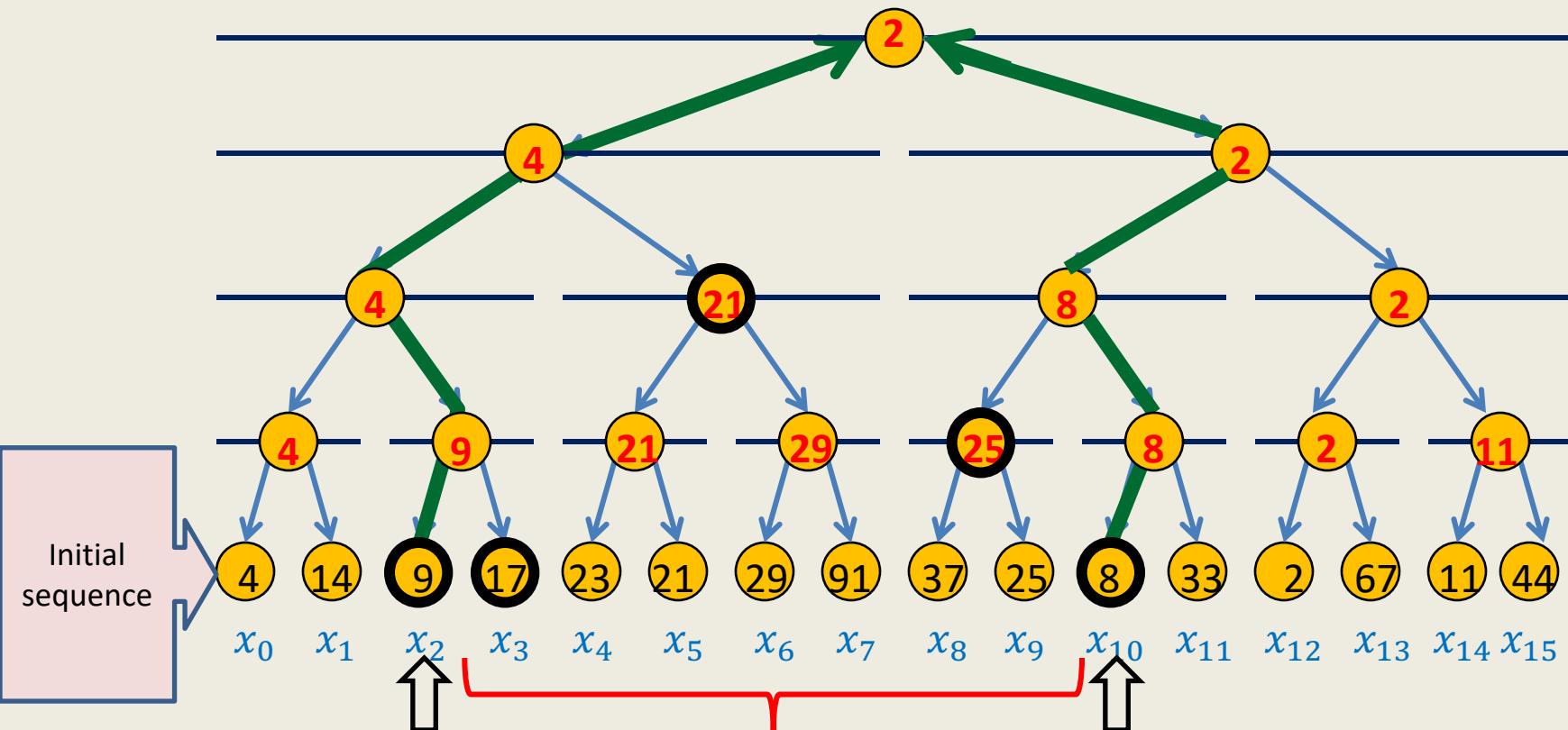
**Answer:**

minimum value

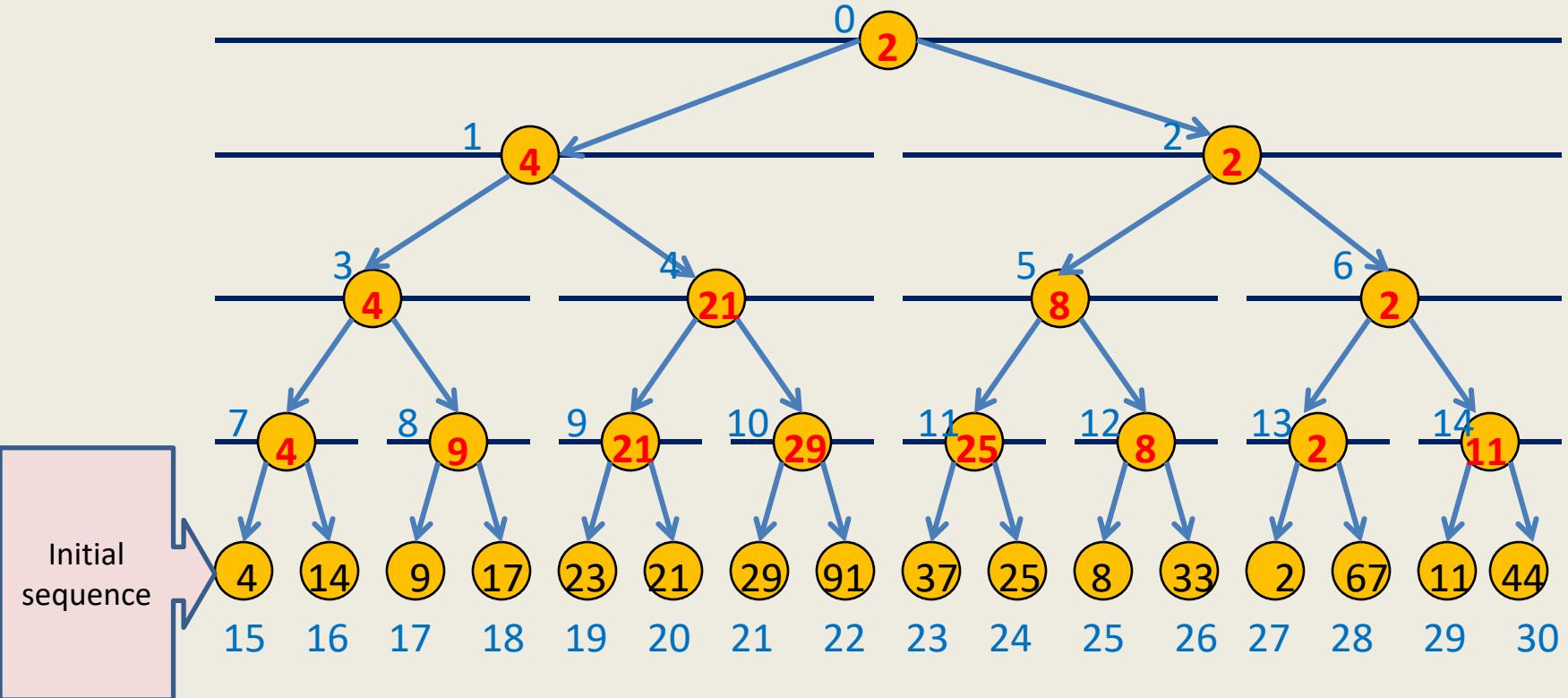
# Data structure for dynamic range minima



# Data structure for dynamic range minima



# Data structure for dynamic range minima



**Data structure:** An array  $A$  of size  $2n-1$ .

Copy the sequence  $S = \langle x_0, \dots, x_{n-1} \rangle$  into  $A[n-1] \dots A[2n-2]$

Leaf node corresponding to  $x_i = A[(n-1) + i]$

How to check if a node is left child or right child of its parent ?

(if index of the node is odd, then the node is left child, else the node is right child)

# Update(*i, a*)

Update(*i, a*)

*i*  $\leftarrow$  (*n - 1*) + *i* ;

A[*i*]  $\leftarrow$  *a* ;

*i*  $\leftarrow$  [(*i - 1*)/2] ;

While( ?? )

{

## Homework

At the end of the lecture slides, an **incorrect** pseudocode is given for Update(*i, a*).  
It is followed by a **correct** one.

Ponder over it ☺

}

# Report-Min(*i,j*)

## Report-Min(*i,j*)

```
i < (n - 1) + i ;  
j < (n - 1) + j ;  
min <- A(i) ;  
If (j > i)  
{     If (A(j) < min)   min <- A(j) ;  
    While(  $\lfloor (\iota - 1)/2 \rfloor \neq \lfloor (j - 1)/2 \rfloor$  )  
    {  
        If(   i%2=1 and A(i + 1) < min   )   min <- ;  
        If(   j%2=0 and A(j - 1) < min   )   min <- ;  
        i < ;  
        j < ;  
    }  
}  
return min ;
```

## **Another interesting problem on sequences**

# Practice Problem

Given an initial sequence  $S = \langle x_0, \dots, x_{n-1} \rangle$  of  $n$  numbers,  
maintain a compact data structure to perform the following operations efficiently :

- **Report\_min( $i, j$ ):**

Report the minimum element from  $\{x_i, \dots, x_j\}$ .

- **Multi-Increment( $i, j, \Delta$ ):**

Add  $\Delta$  to each  $x_k$  for each  $i \leq k \leq j$

---

## Example:

Let the initial sequence be  $S = \langle 14, 12, 3, 12, 111, 51, 321, -40 \rangle$

- **Report\_min(1, 4):**

returns 3

- **Multi-Increment(2, 6, 10):**

$S$  becomes  $\langle 14, 12, 13, 22, 121, 61, 331, -40 \rangle$

- **Report\_min(1, 4):**

returns 12

# An **challenging problem** on sequences

**For summer vacation  
(not for the exam)**

# \* Problem

Given an initial sequence  $S = \langle x_0, \dots, x_{n-1} \rangle$  of  $n$  numbers,  
maintain a compact data structure to perform the following operations efficiently :

- **Report\_min( $i, j$ ):**

Report the minimum element from  $\{x_i, \dots, x_j\}$ .

- **Multi-Increment( $i, j, \Delta$ ):**

Add  $\Delta$  to each  $x_k$  for each  $i \leq k \leq j$

- **Rotate( $i, j$ ):**

$x_i \leftrightarrow x_j, x_{i+1} \leftrightarrow x_{j-1}, \dots$

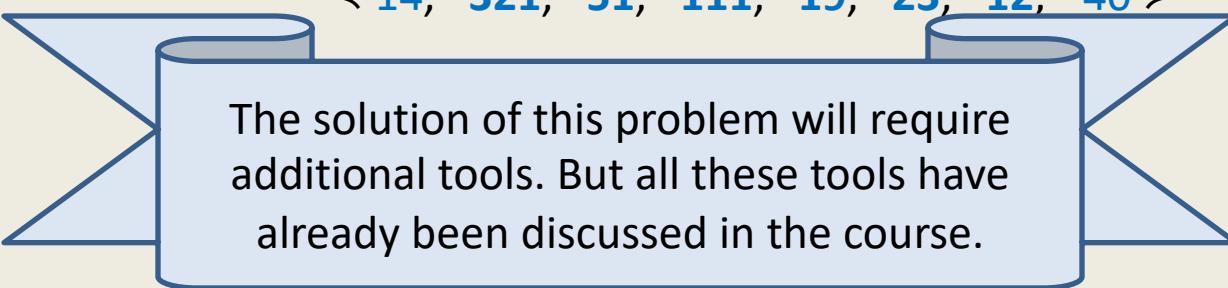
---

## Example:

Let the initial sequence be  $S = \langle 14, 12, 23, 19, 111, 51, 321, -40 \rangle$

After **Rotate(1,6)**,  $S$  becomes

$\langle 14, 321, 51, 111, 19, 23, 12, -40 \rangle$



The solution of this problem will require additional tools. But all these tools have already been discussed in the course.

# **Problem 4**

**A data structure for sets**

# Sets under operations

**Given:** a collection of  $n$  singleton sets  $\{0\}, \{1\}, \{2\}, \dots \{n - 1\}$

**Aim:** a compact data structure to perform

- **Union( $i, j$ ):**  
Unite the two sets containing  $i$  and  $j$ .
  - **Same\_sets( $i, j$ ):**  
Determine if  $i$  and  $j$  belong to the same set.
- 

## Trivial Solution 1

Keep an array **Label[]** such that

**Label[ $i$ ]=Label[ $j$ ]** if and only if  $i$  and  $j$  belong to the same set.

- **Same\_sets( $i, j$ ):**  $O(1)$  time
- check if **Label[ $i$ ]=Label[ $j$ ]** ?
- **Union( $i, j$ ):**  $O(n)$  time
- For each  $0 \leq k < n$
- if (**Label[ $k$ ]=Label[ $i$ ]**)      **Label[ $k$ ] ← Label[ $j$ ])**

# Sets under operations

**Given:** a collection of  $n$  singleton sets  $\{0\}, \{1\}, \{2\}, \dots \{n - 1\}$

**Aim:** a compact data structure to perform

- **Union( $i, j$ ):**  
Unite the two sets containing  $i$  and  $j$ .
  - **Same\_sets( $i, j$ ):**  
Determine if  $i$  and  $j$  belong to the same set.
- 

## Trivial Solution 2

Treat the problem as a graph problem: ??

Connected component

- $V = \{0, \dots, n - 1\}$ ,  $E =$  empty set initially.
- A set  $\Leftrightarrow$
- Keep array **Label[]** such that **Label[ $i$ ] = Label[ $j$ ]** iff  $i$  and  $j$  belong to the same component.



**Union( $i, j$ ) :**

**O( $n$ ) time**

add an edge  $(i, j)$  and  
**recompute** connected components using **BFS/DFS**.

# Sets under operations

**Given:** a collection of  $n$  singleton sets  $\{0\}, \{1\}, \{2\}, \dots \{n - 1\}$

**Aim:** a compact data structure to perform

- **Union( $i, j$ ):**  
Unite the two sets containing  $i$  and  $j$ .
  - **Same\_sets( $i, j$ ):**  
Determine if  $i$  and  $j$  belong to the same set.
- 

**Efficient solution:**

- A data structure which supports each operation in  $O(\log n)$  time.
- **An additional heuristic**  
→ time complexity of an operation :

We shall discuss it in the next  
lecture.

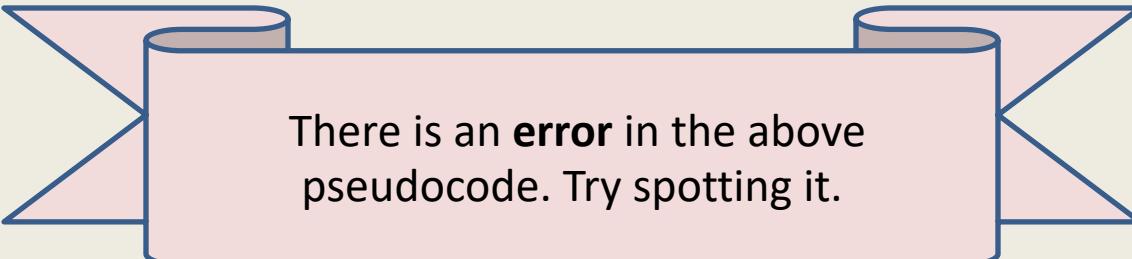
# Homework

For **Dynamic Range-minima** problem

# Update(*i, a*)

Update(*i, a*)

```
i < (n - 1) + i ;  
A[i] <= a ;  
i <= [(i - 1)/2] ;  
While(           i ≥ 0           )  
{  
    If(a < A[i])  A[i] <= a;  
    i <= [(i - 1)/2] ;  
}
```

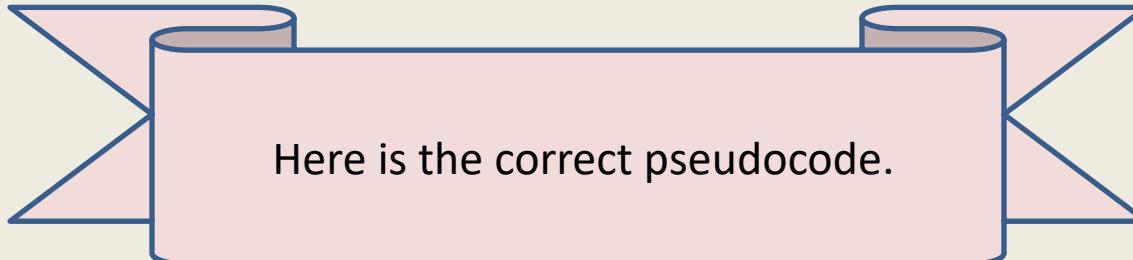


There is an **error** in the above pseudocode. Try spotting it.

# Update(*i, a*)

Update(*i, a*)

```
i ← (n - 1) + i ;  
A[i] ← a ;  
i ← [(i - 1)/2] ;  
While( i ≥ 0 )  
{  
    If( A[2i + 1] < A[2i + 2])  
        A[i] ← A[2i + 1]  
    else  
        A[i] ← A[2i + 2] ;  
    i ← [(i - 1)/2] ;  
}
```



Here is the correct pseudocode.

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 32**

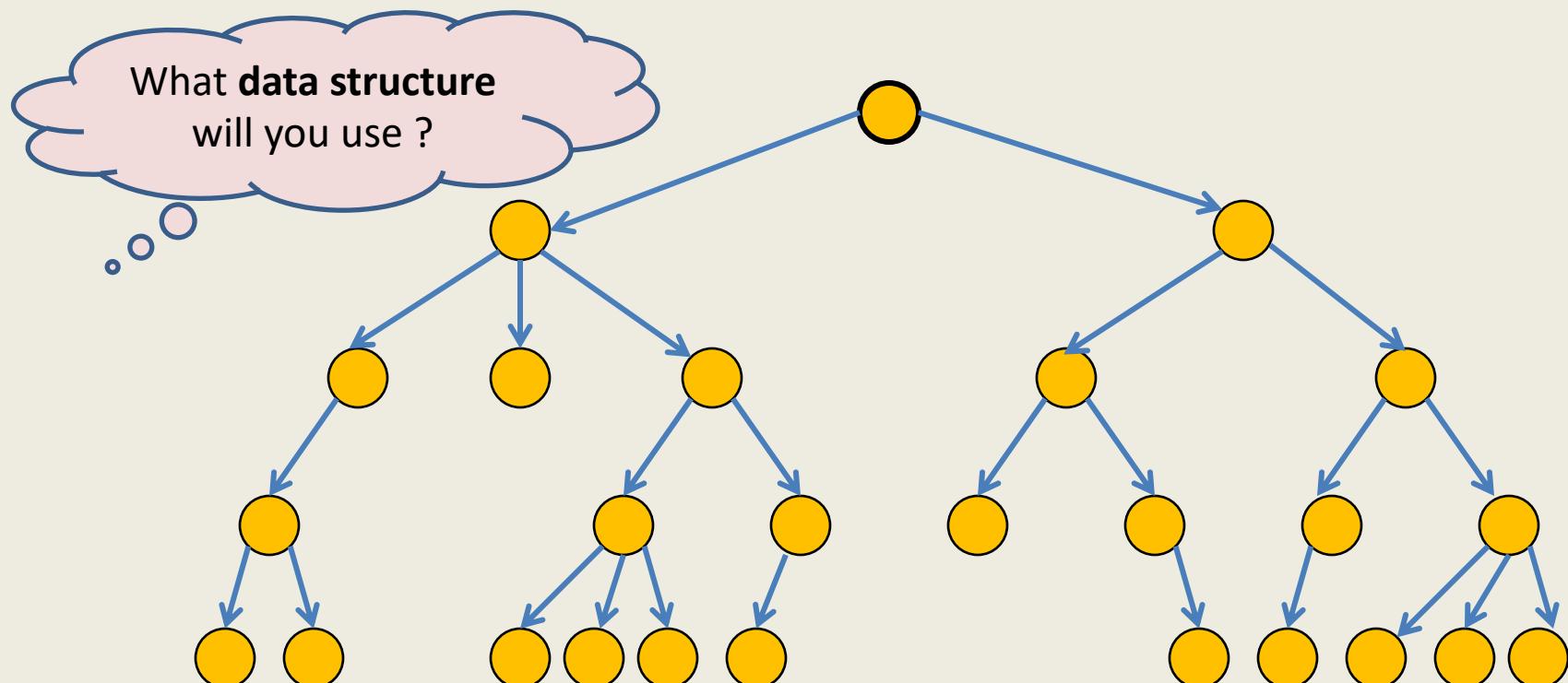
- **Magical application of binary trees – III**

**Data structure for sets**

# **Rooted tree**

**Revisiting and extending**

# A typical rooted tree we studied



**Definition we gave:**

Every vertex, except **root**, has exactly one incoming edge and has a path **from** the root.

**Examples:**

Binary search trees,

DFS tree,

BFS tree.

# A typical rooted tree we studied

**Question:** what data structure can be used for representing a rooted tree ?

**Answer:**

**Data structure 1:**

- Each node stores a list of its children.
- To access the tree, we keep a pointer to the root node.

(there is no way to access any node (other than root) directly in this data structure)

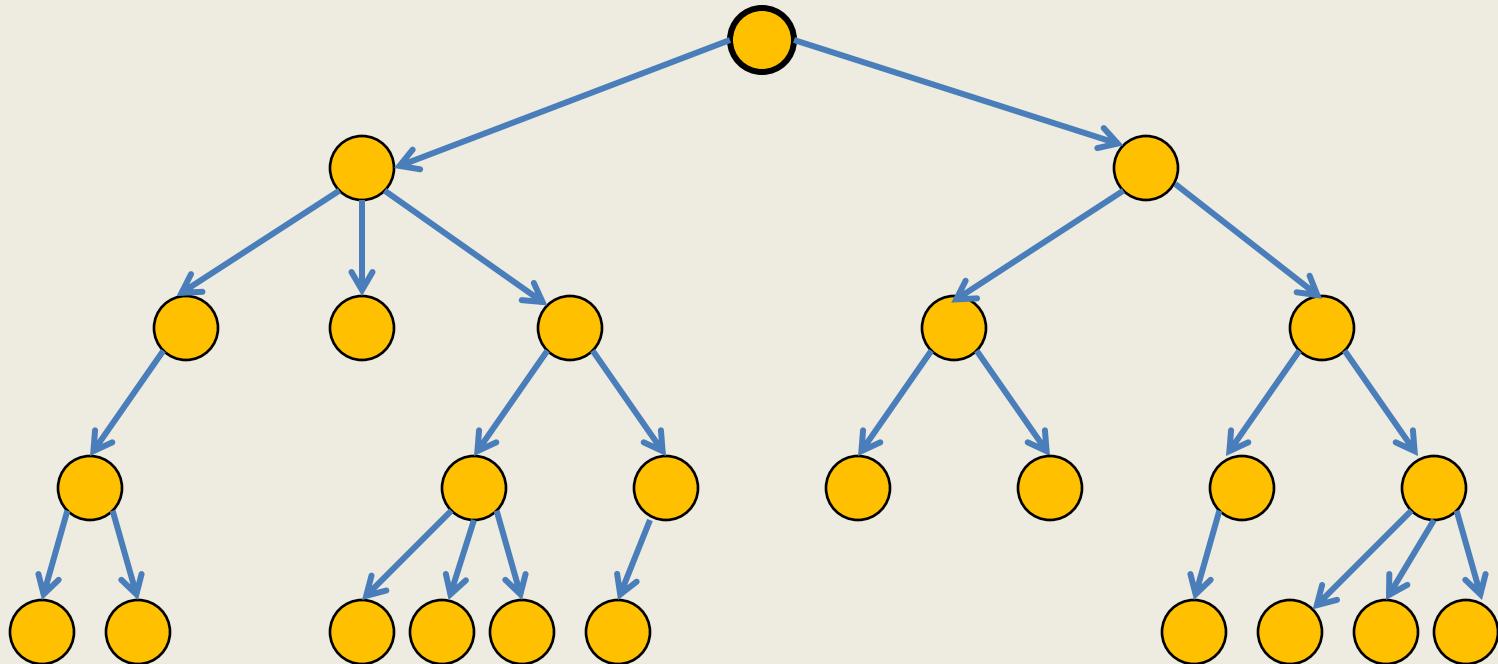
**Data structure 2:** (If nodes are labeled in a contiguous range [0..n-1])

rooted tree becomes an instance of a **directed graph**.

So we may use **adjacency list** representation.

**Advantage:** We can access each node directly.

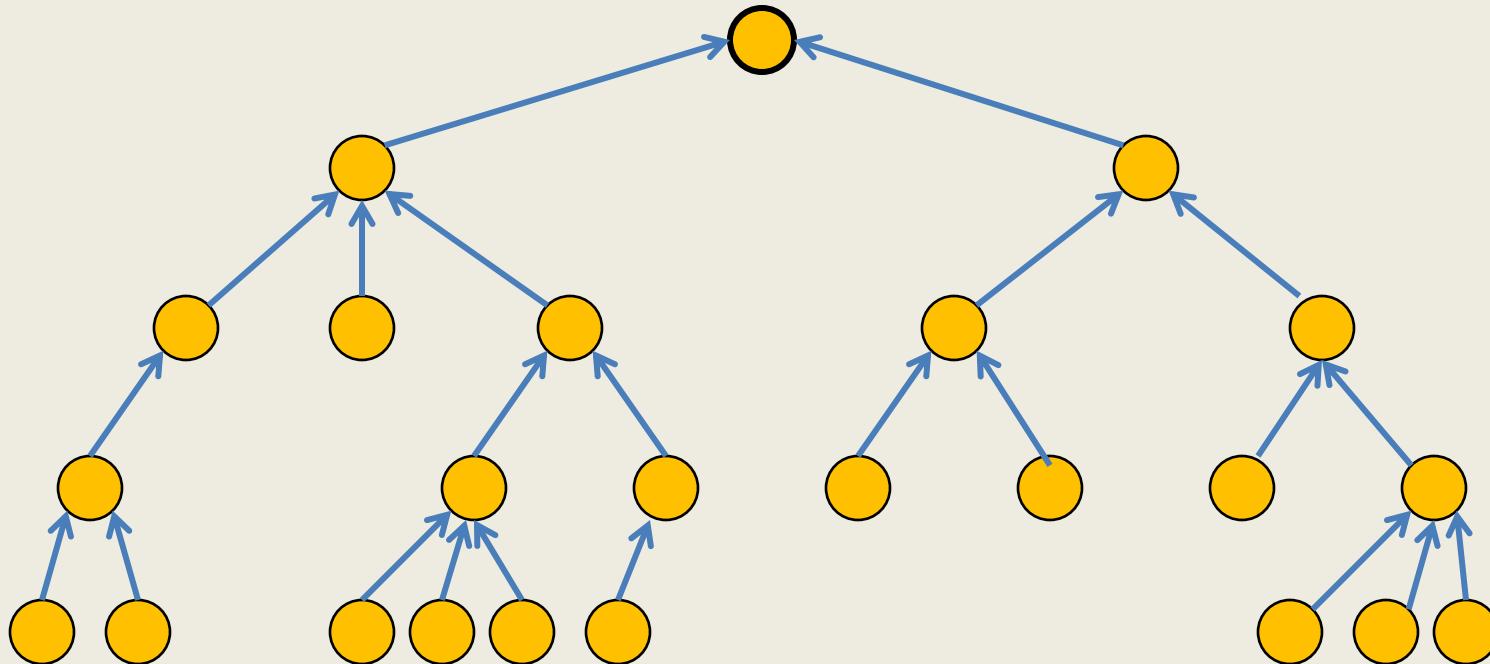
# Extending the definition of rooted tree



## Extended Definition:

Type 1: Every vertex, except **root**, has exactly one incoming edge and has a path **from** the root.

# Extending the definition of rooted tree



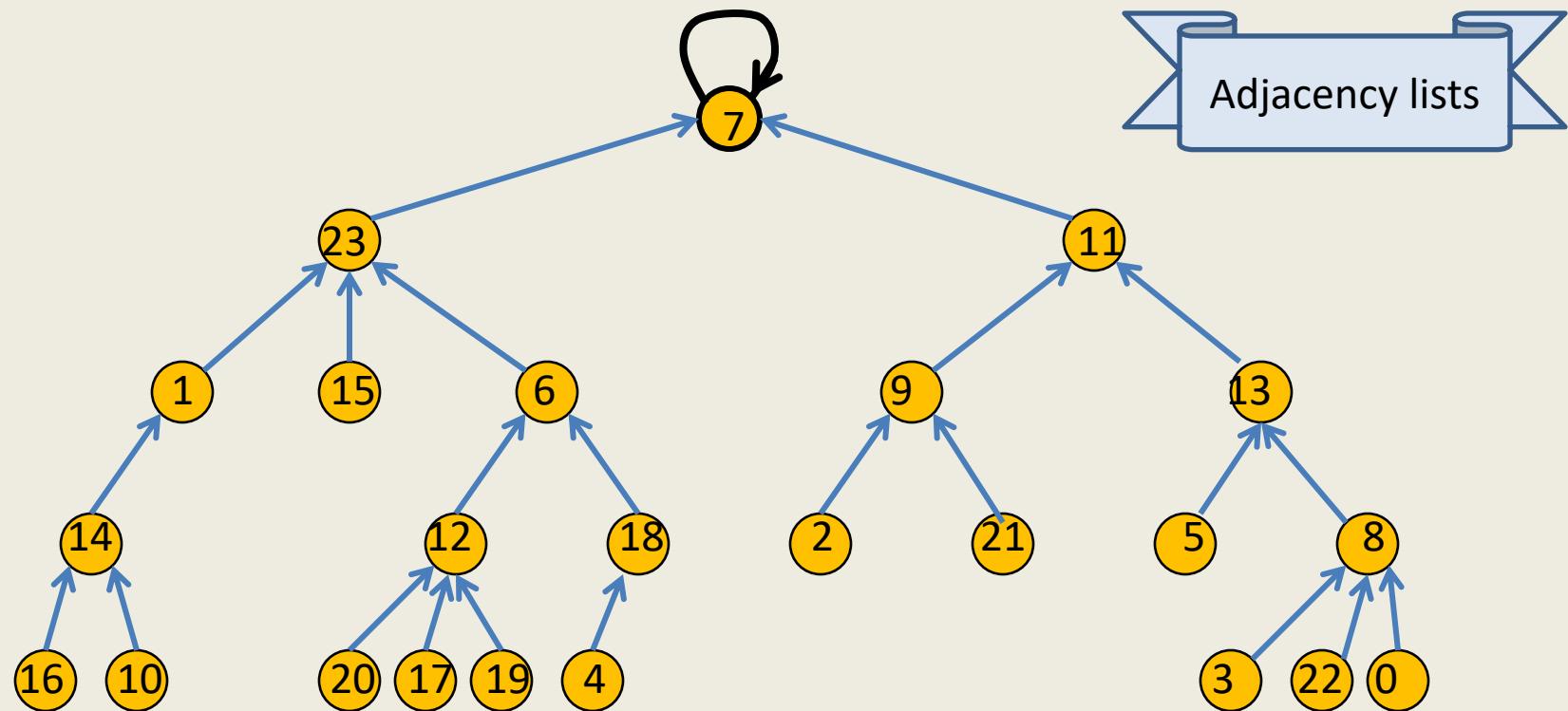
## Extended Definition:

Type 1: Every vertex, except **root**, has exactly one incoming edge and has a path **from** the root.

OR

Type 2: Every vertex, except root, has exactly one outgoing edge and has a path **to** the root.

# Data structure for rooted tree of type 2



If nodes are labeled in a contiguous range [0.. $n - 1$ ],  
there is even simpler and more compact data structure

Guess ??

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

# **Application** of rooted tree of **type 2**

Maintaining sets

# Sets under two operations

**Given:** a collection of  $n$  singleton sets  $\{0\}, \{1\}, \{2\}, \dots \{n - 1\}$

**Aim:** a compact data structure to perform

- **Union( $i, j$ ):**  
Unite the two sets containing  $i$  and  $j$ .
  - **Same\_sets( $i, j$ ):**  
Determine if  $i$  and  $j$  belong to the same set.
- 

## Trivial Solution

Treat the problem as a graph problem: **Connected component**

- $V = \{0, \dots, n - 1\}$ ,  $E =$  empty set initially.
- A set  $\Leftrightarrow$
- Keep array **Label[]**



**Union( $i, j$ ):**

if (**Same\_sets( $i, j$ )** = **false**)  
add an edge ( $i, j$ )

**O( $n$ ) time**

# Sets under two operations

**Given:** a collection of  $n$  singleton sets  $\{0\}, \{1\}, \{2\}, \dots \{n - 1\}$

**Aim:** a compact data structure to perform

- **Union( $i, j$ ):**  
Unite the two sets containing  $i$  and  $j$ .
  - **Same\_sets( $i, j$ ):**  
Determine if  $i$  and  $j$  belong to the same set.
- 

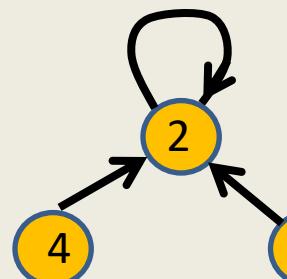
## Efficient solution:

- A data structure which supports each operation in  $O(\log n)$  time.
- **An additional heuristic**  
→ time complexity of an operation : practically  $O(1)$ .

# Data structure for sets

Maintain each set as a rooted tree .

Union( $i, j$ ) ?



{2,4,6}

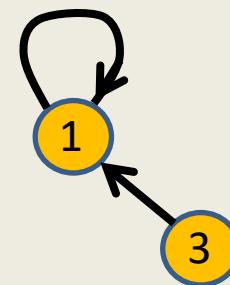


{0}



{5}

SameSet( $i, j$ ) ?



{1,3}

# Data structure for sets

Maintain each set as a rooted tree .

**Question:** How to perform operation  $\text{Same\_sets}(i, j)$  ?

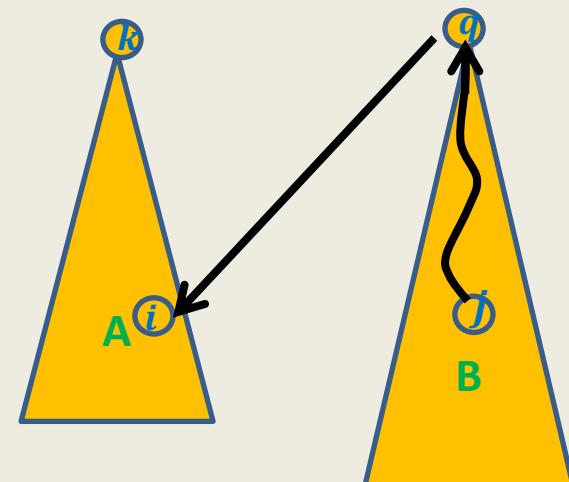
**Answer:** Determine if  $i$  and  $j$  belong to the same tree.

→ find root of  $i$

**Question:** How to perform  $\text{Union}(i, j)$  ?

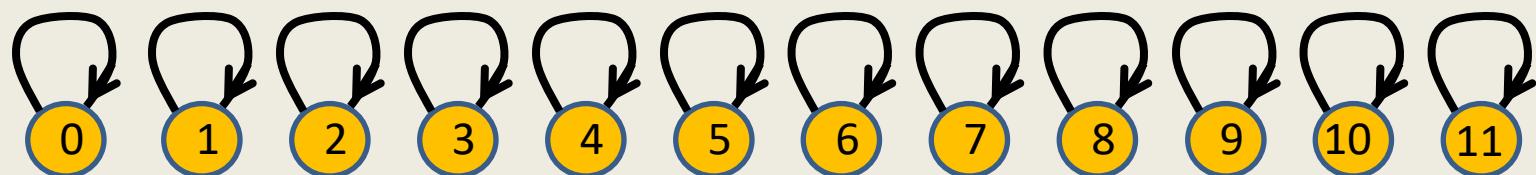
**Answer:**

- find root of  $j$ ; let it be  $q$ .
- $\text{Parent}(q) \leftarrow i$ .



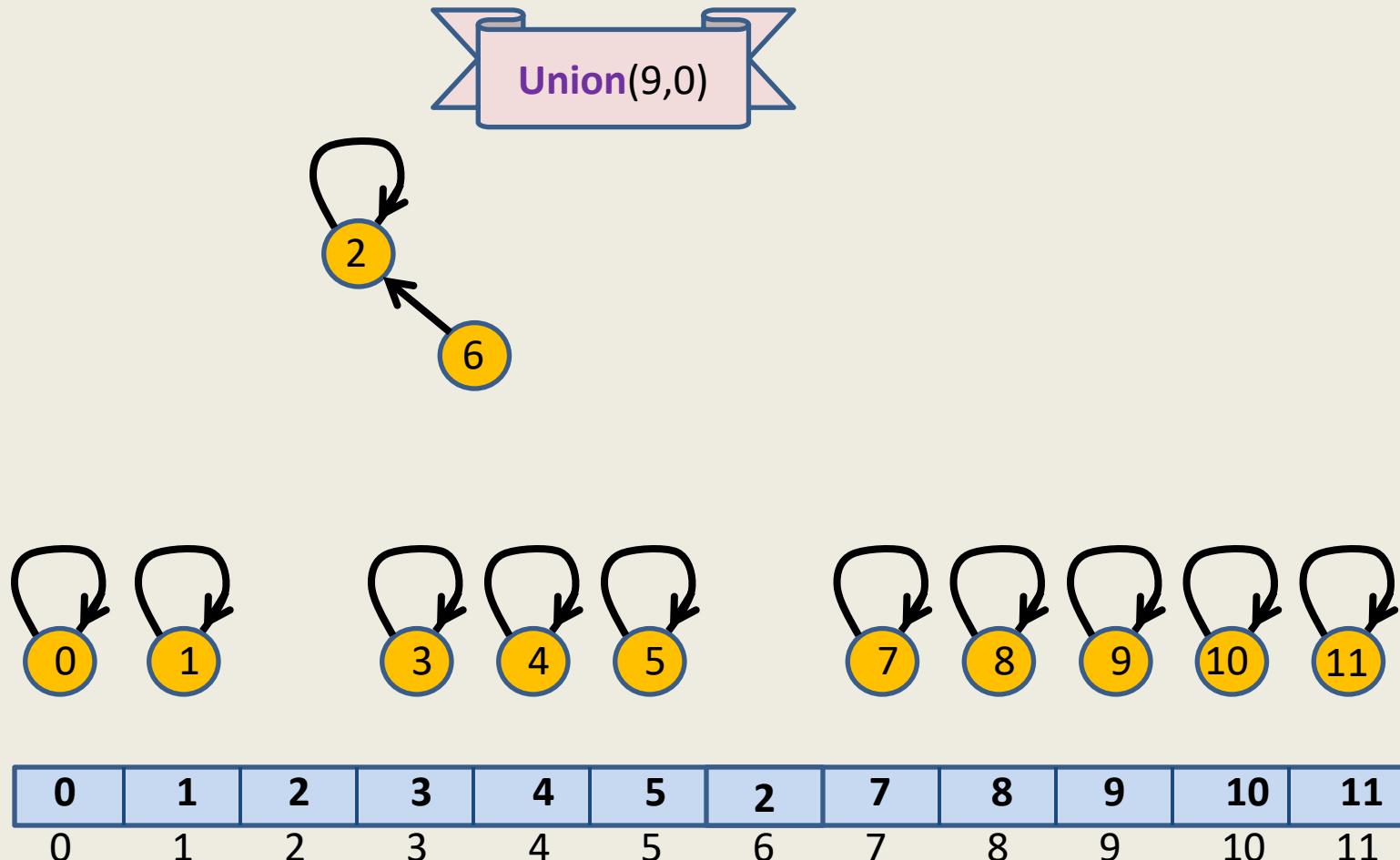
# A rooted tree as a data structure for sets

Union(2,6)

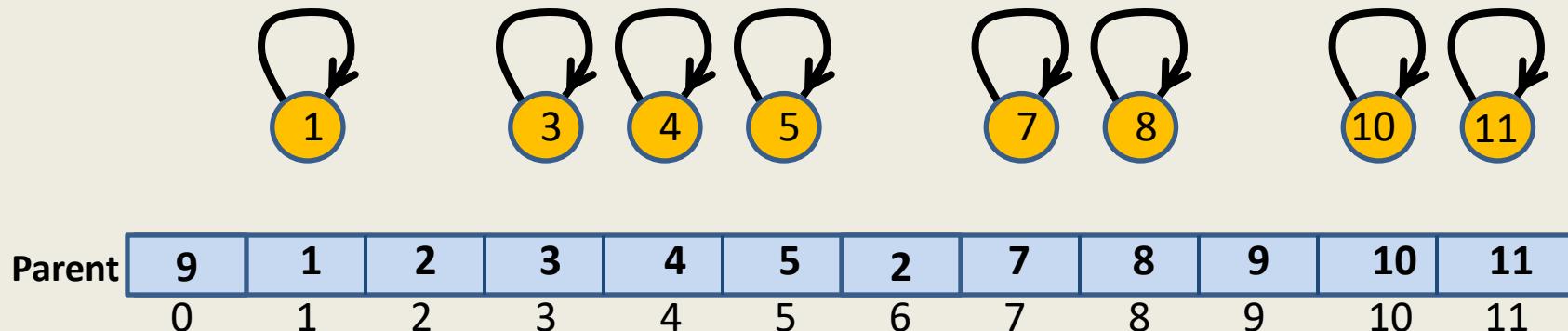
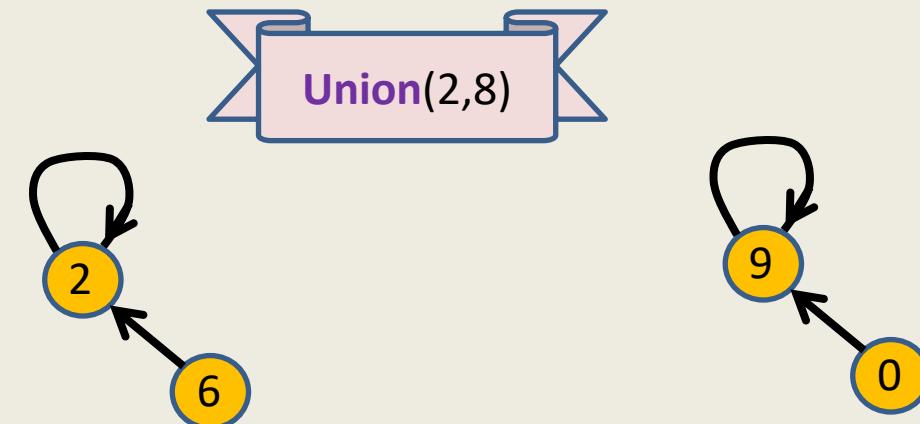


Parent	0	1	2	3	4	5	6	7	8	9	10	11
	0	1	2	3	4	5	6	7	8	9	10	11

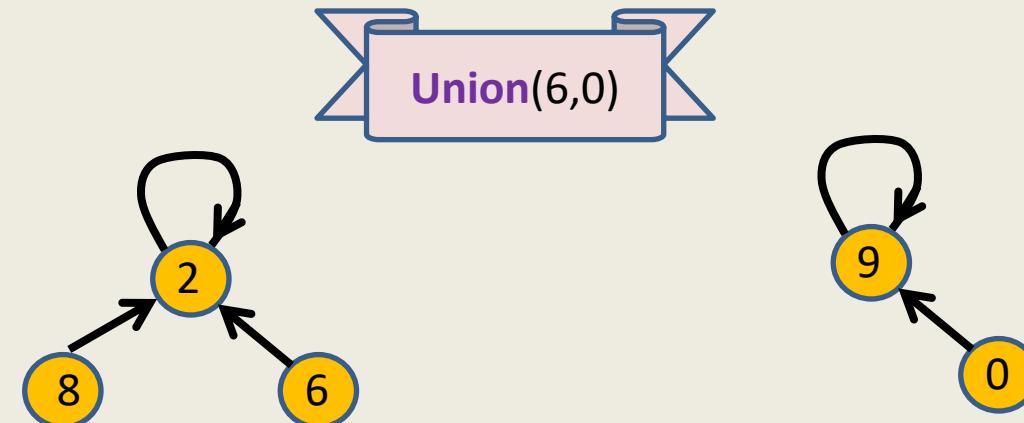
# A rooted tree as a data structure for sets



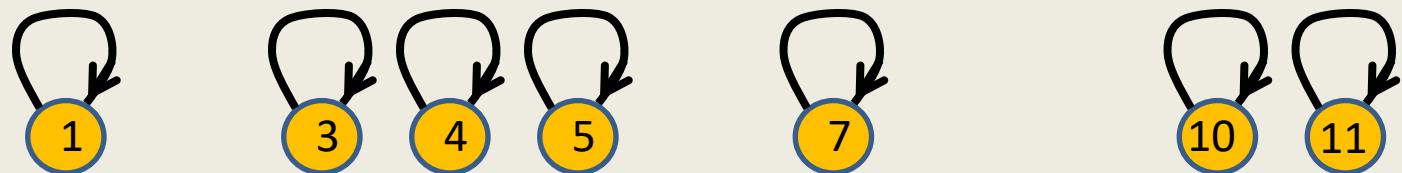
# A rooted tree as a data structure for sets



# A rooted tree as a data structure for sets

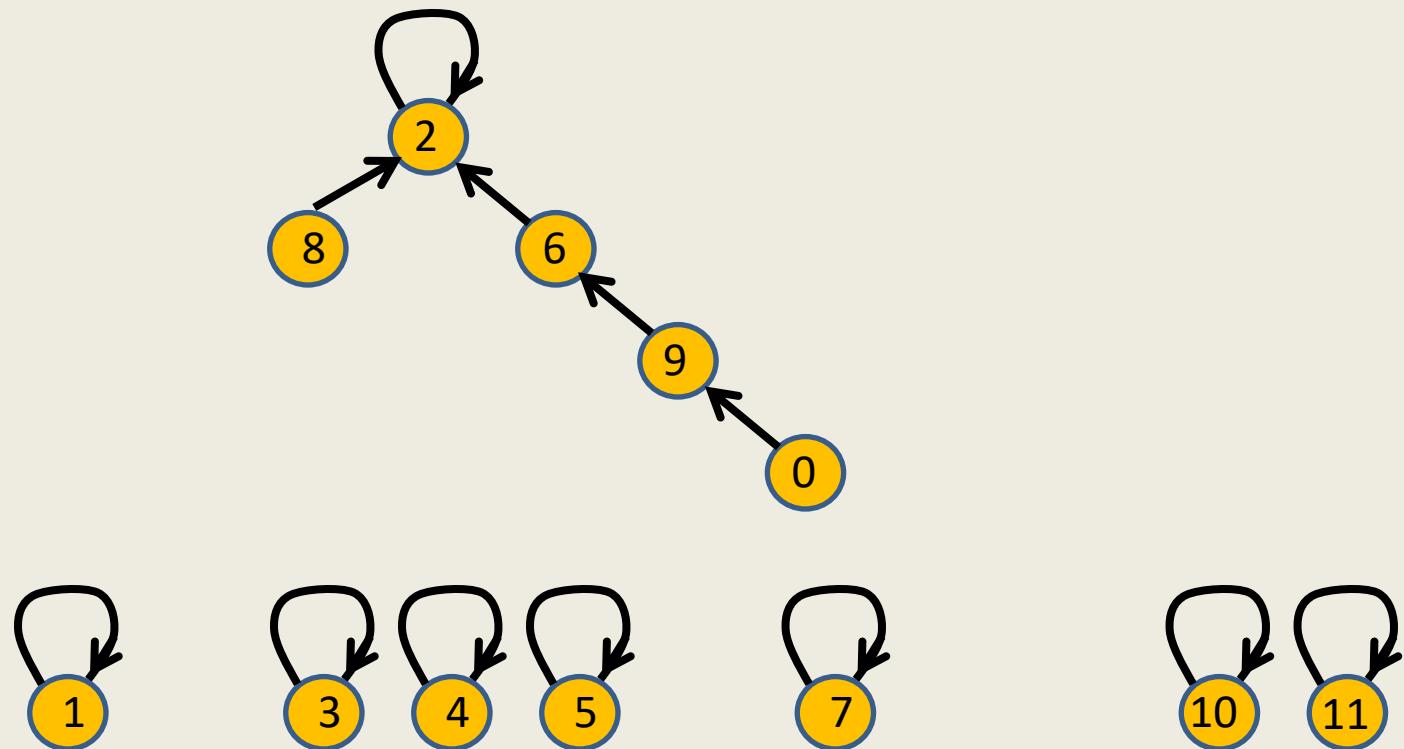


Union(6,0)



Parent	9	1	2	3	4	5	2	7	2	9	10	11
	0	1	2	3	4	5	6	7	8	9	10	11

# A rooted tree as a data structure for sets



Parent	9	1	2	3	4	5	2	7	2	6	10	11
	0	1	2	3	4	5	6	7	8	9	10	11

# Pseudocode for Union and SameSet()

**Find(*i*)** // subroutine for finding the root of the tree containing *i*

```
If (Parent(i) = i)    return i ;  
else return Find(Parent(i));
```

---

**SameSet(*i, j*)**

```
k  $\leftarrow$  Find(i);  
l  $\leftarrow$  Find(j);  
If (k = l)    return true else return false
```

**Union(*i, j*)**

```
k  $\leftarrow$  Find(j);  
Parent(k)  $\leftarrow$  i;
```

**Observation:** Time complexity of **Union(*i, j*)** as well as **Same\_sets(*i, j*)** is governed by the time complexity of **Find(*i*)** and **Find(*j*)**.

**Question:** What is time complexity of **Find(*i*)** ?

**Answer:** **depth** of the node *i* in the tree containing *i*.

# Time complexity of $\text{Find}(i)$

$\text{Union}(0,1)$

$\text{Union}(1,2)$

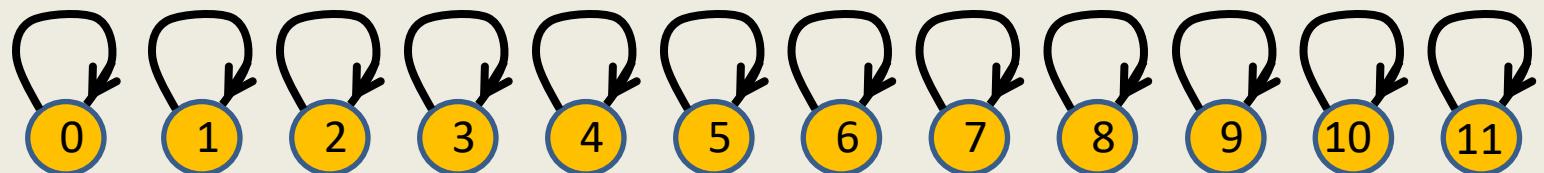
$\text{Union}(2,3)$

...

$\text{Union}(9,10)$

$\text{Union}(10,11)$

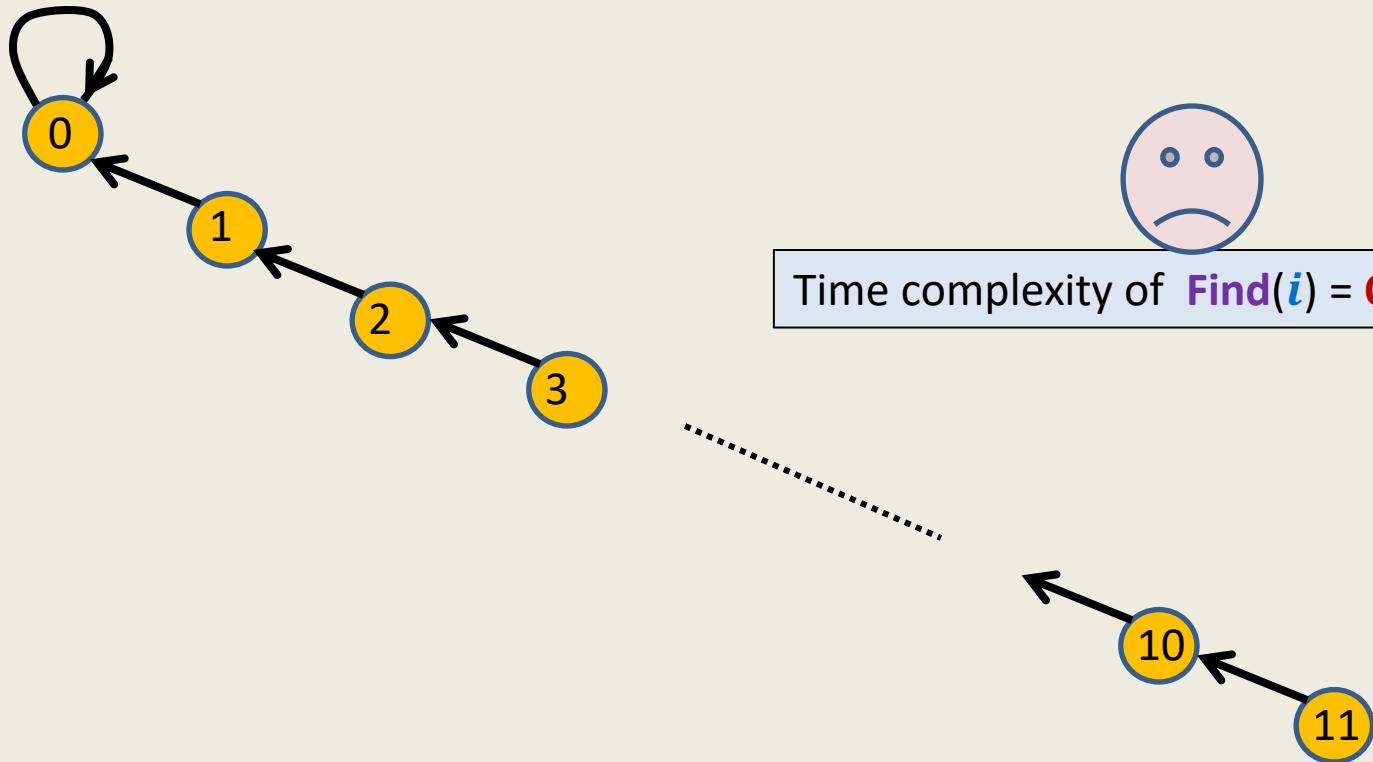
What will be the rooted tree structures  
after these union operations ?



Parent	0	1	2	3	4	5	6	7	8	9	10	11
	0	1	2	3	4	5	6	7	8	9	10	11

# Time complexity of $\text{Find}(i)$

$\text{Union}(0,1)$   
 $\text{Union}(1,2)$   
 $\text{Union}(2,3)$   
...  
 $\text{Union}(9,10)$   
 $\text{Union}(10,11)$

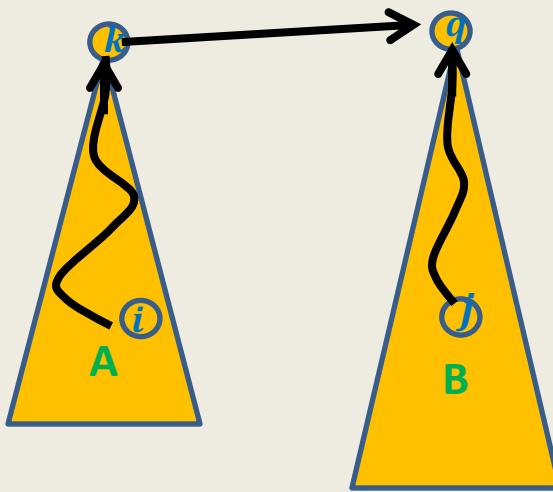


Parent	0	0	1	2	3	4	5	6	7	8	9	10
	0	1	2	3	4	5	6	7	8	9	10	11

# Improving the time complexity of Find(*i*)

**Heuristic 1: Union by size**

# Improving the Time complexity



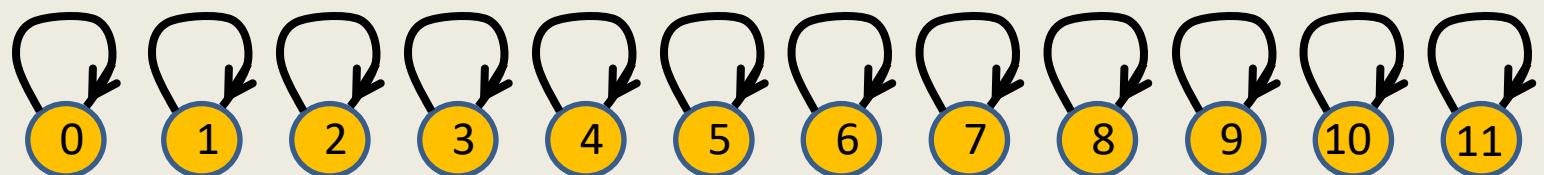
**Key idea:** Change the **union(*i,j*)** .

While doing **union(*i,j*)**, hook the **smaller size tree to the root of the bigger size tree**.

For this purpose, keep an array **size[0,..,n-1]**

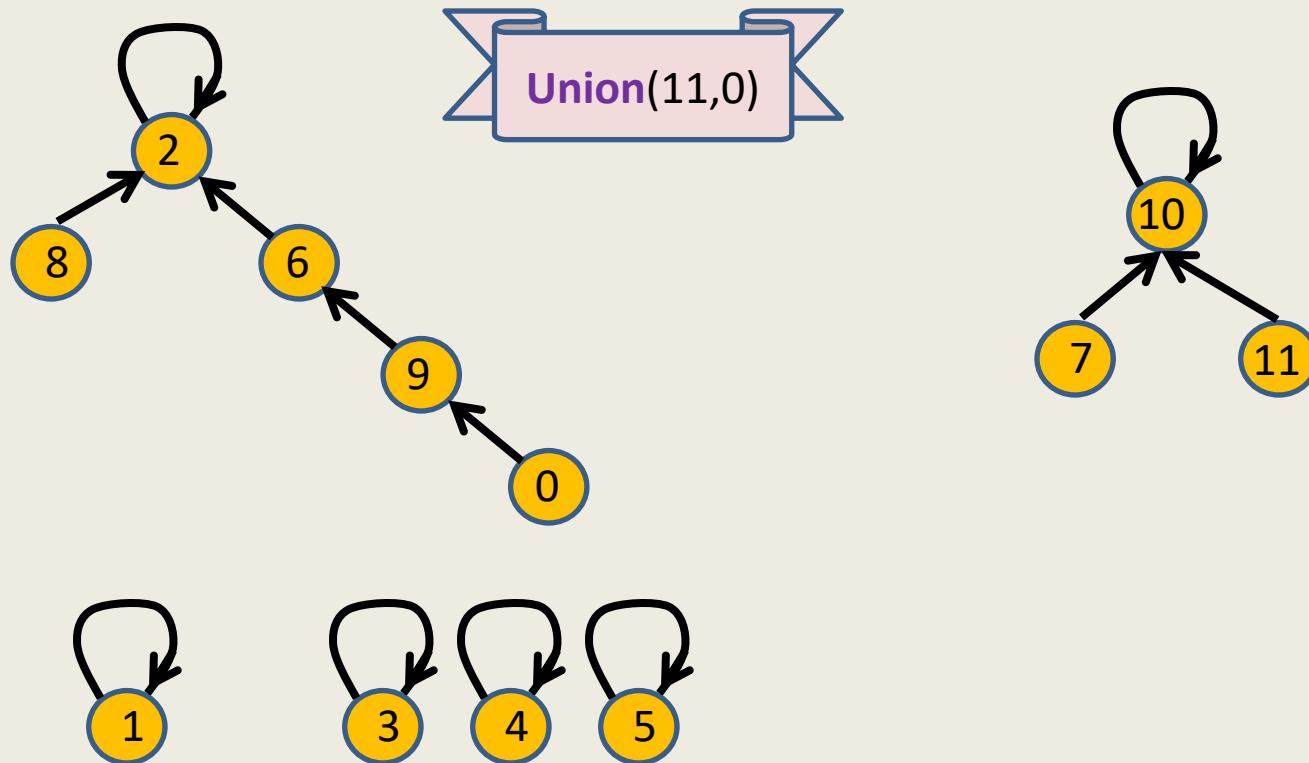
**size[*i*]** = number of nodes in the tree containing *i*  
(if *i* is a **root** and zero otherwise)

# Efficient data structure for sets



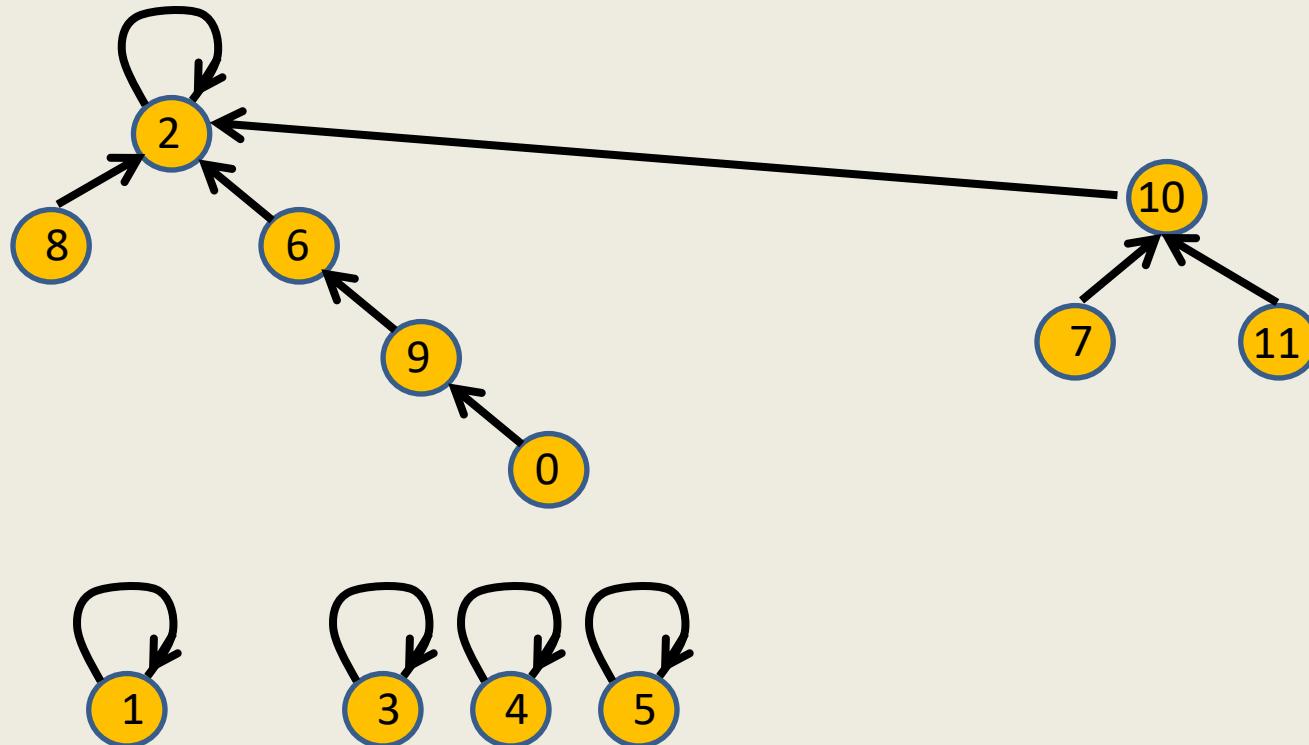
Parent	0	1	2	3	4	5	6	7	8	9	10	11
size	1	1	1	1	1	1	1	1	1	1	1	1

# Efficient data structure for sets



Parent	9	1	2	3	4	5	2	10	2	6	10	10
size	0	1	2	3	4	5	6	7	8	9	10	11
size	0	1	5	1	1	1	0	0	0	0	3	0

# Efficient data structure for sets



Parent	9	1	2	3	4	5	2	10	2	6	10	10
size	0	1	5	1	1	1	0	0	0	3	0	
size	0	1	5	1	1	1	0	0	0	3	0	

# Pseudocode for modified Union

**Union( $i, j$ )**

$k \leftarrow \text{Find}(i);$

$l \leftarrow \text{Find}(j);$

**If**(size( $k$ ) < size( $l$ ) )

$l \leftarrow \text{Parent}(k);$

        size( $l$ )  $\leftarrow$  size( $k$ ) + size( $l$ );

        size( $k$ )  $\leftarrow$  0;

**Else**

$k \leftarrow \text{Parent}(l);$

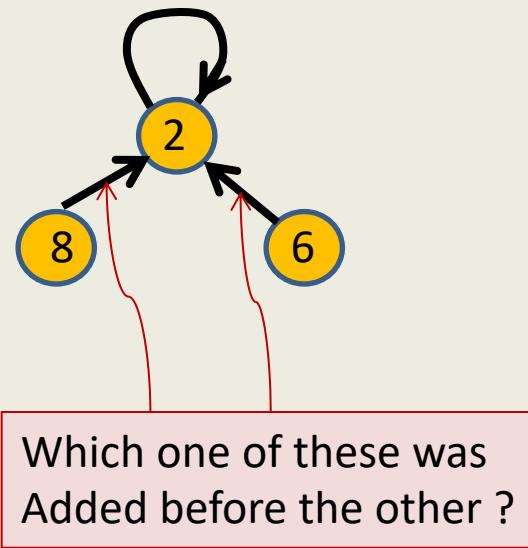
        size( $k$ )  $\leftarrow$  size( $k$ ) + size( $l$ );

        size( $l$ )  $\leftarrow$  0;

**Question:** How to show that **Find( $i$ )** for any  $i$  will now take  $O(\log n)$  time only ?

**Answer:** It suffices if we can show that **Depth( $i$ )** is  $O(\log n)$ .

# Can we infer history of a tree?

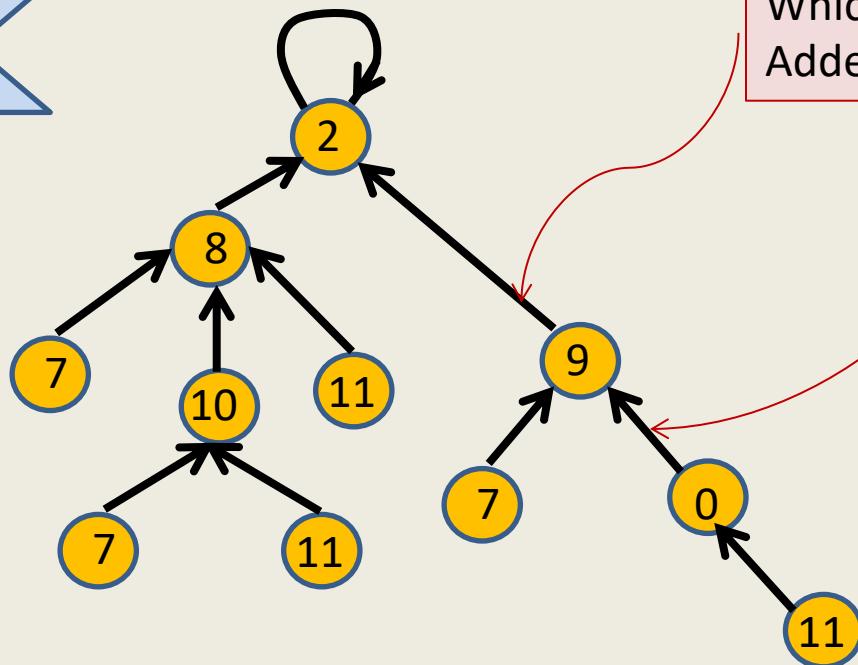


Answer: Can not be inferred with any certainty 😔.

# Can we infer history of a tree?

During union, we join **roots** of two trees.

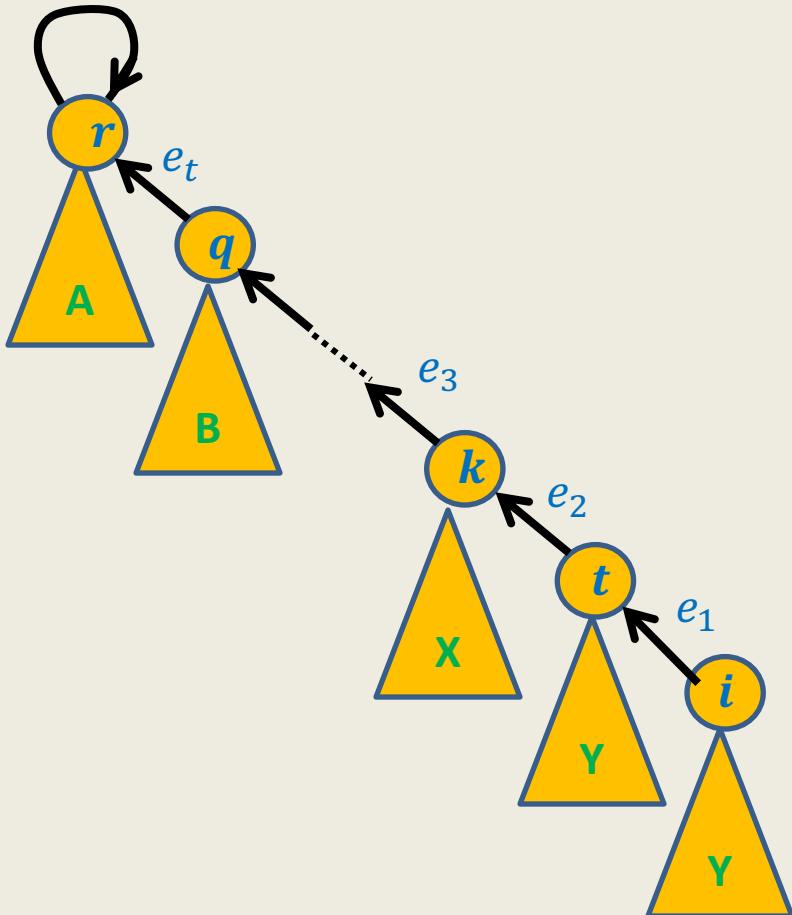
Which one of these was  
Added before the other ?



→  $(0 \rightarrow 9)$  was added **before**  $(9 \rightarrow 2)$ .

**Theorem:** The edges on a **path** from node  $v$  to root were inserted in the order they appear on the **path**.

# How to show that depth of any element = $O(\log n)$ ?



Let  $e_1, e_2, \dots, e_t$  be the edges on the path from  $i$  to the **root**.

Let us visit the history.  
(how this tree came into being ? ).

Edges  $e_1, e_2, \dots, e_t$  would have been added in the order:

$e_1$   
 $e_2$   
...  
 $e_t$

# How to show that depth of any element = $O(\log n)$ ?

Let  $e_1, e_2, \dots, e_t$  be the edges on the path from  $i$  to the root.

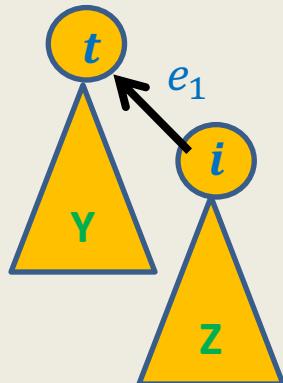
Consider the moment just before edge  $e_1$  is inserted.

Let no. of elements in subtree  $T(i)$  at that moment be  $n_i$ .

We added edge  $i \rightarrow t$  (and not  $t \rightarrow i$ ).

→ no. of elements in  $T(t) \geq n_i$ .

→ After the edge  $i \rightarrow t$  is inserted,  
no. of element in  $T(t) \geq 2n_i$

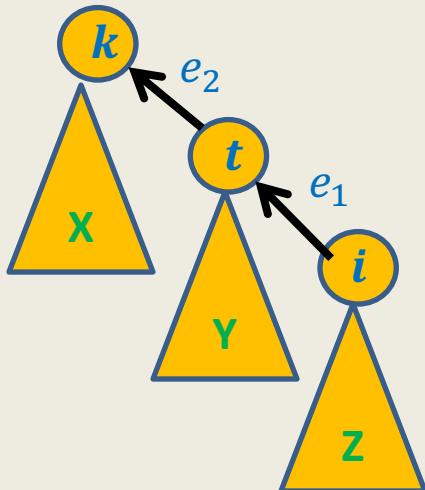


# How to show that depth of any element = $O(\log n)$ ?

Let  $e_1, e_2, \dots, e_t$  be the edges on the path from  $i$  to the root.

Consider the moment just before edge  $e_2$  is inserted.

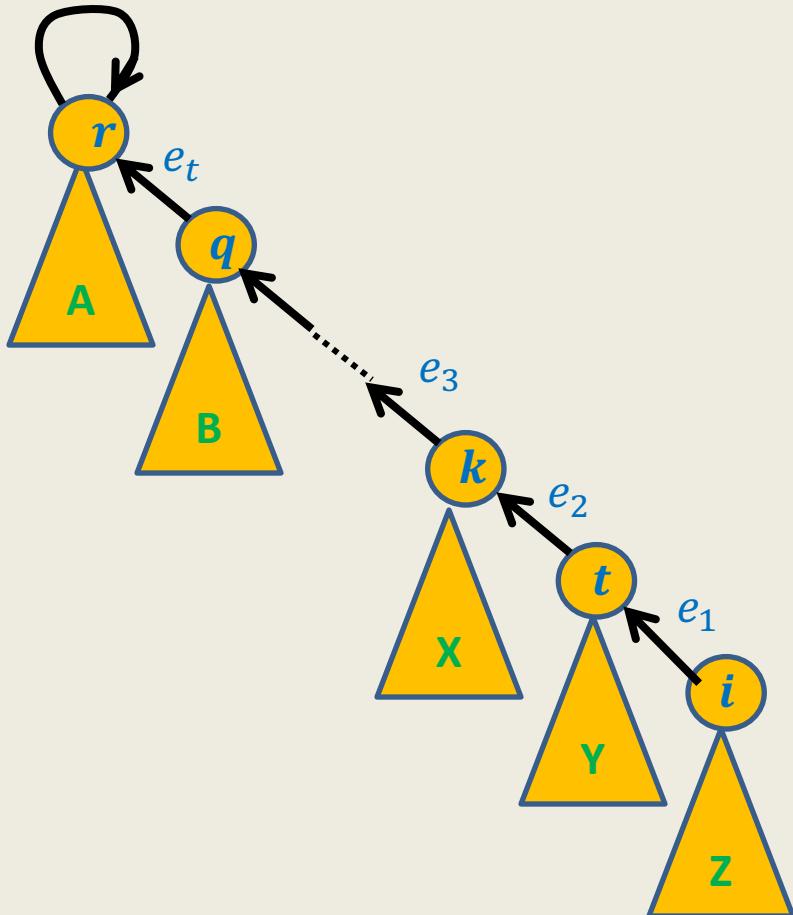
no. of element in  $T(t) \geq 2n_i$



We added edge  $t \rightarrow k$  (and not  $k \rightarrow t$ ).

- # elements in  $T(k) \geq 2n_i$ .
- After the edge  $t \rightarrow k$  is inserted,  
no. of element in  $T(k) \geq 4n_i$

# How to show that depth of any element = $O(\log n)$ ?



Let  $e_1, e_2, \dots, e_t$  be the edges on the path from  $i$  to the root.

Arguing in a similar manner for edge  $e_3, \dots, e_t \rightarrow$

# elements in  $T(r)$  after insertion of  $e_t \geq 2^t n_i$

Obviously  $2^t n_i \leq n$

→

Theorem:  $t \leq \log_2 n$

**Theorem:** Given a collection of  $n$  singleton sets followed by a sequence of **union** and **find** operations, there is a data structure based that achieves  $O(\log n)$  time per operation.

**Question:** Can we achieve even better bounds ?

**Answer:** Yes.

# A new heuristic for better time complexity

**Heuristic 2: Path compression**

# This is how this heuristic got invented

- The time complexity of a **Find(*i*)** operation is proportional to the depth of the node *i* in its rooted tree.
- If the elements are stored closer to the root, faster will the **Find()** be and hence faster will be the overall algorithm.

The algorithm for **Union** and **Find** was used in some application of **data-bases**.

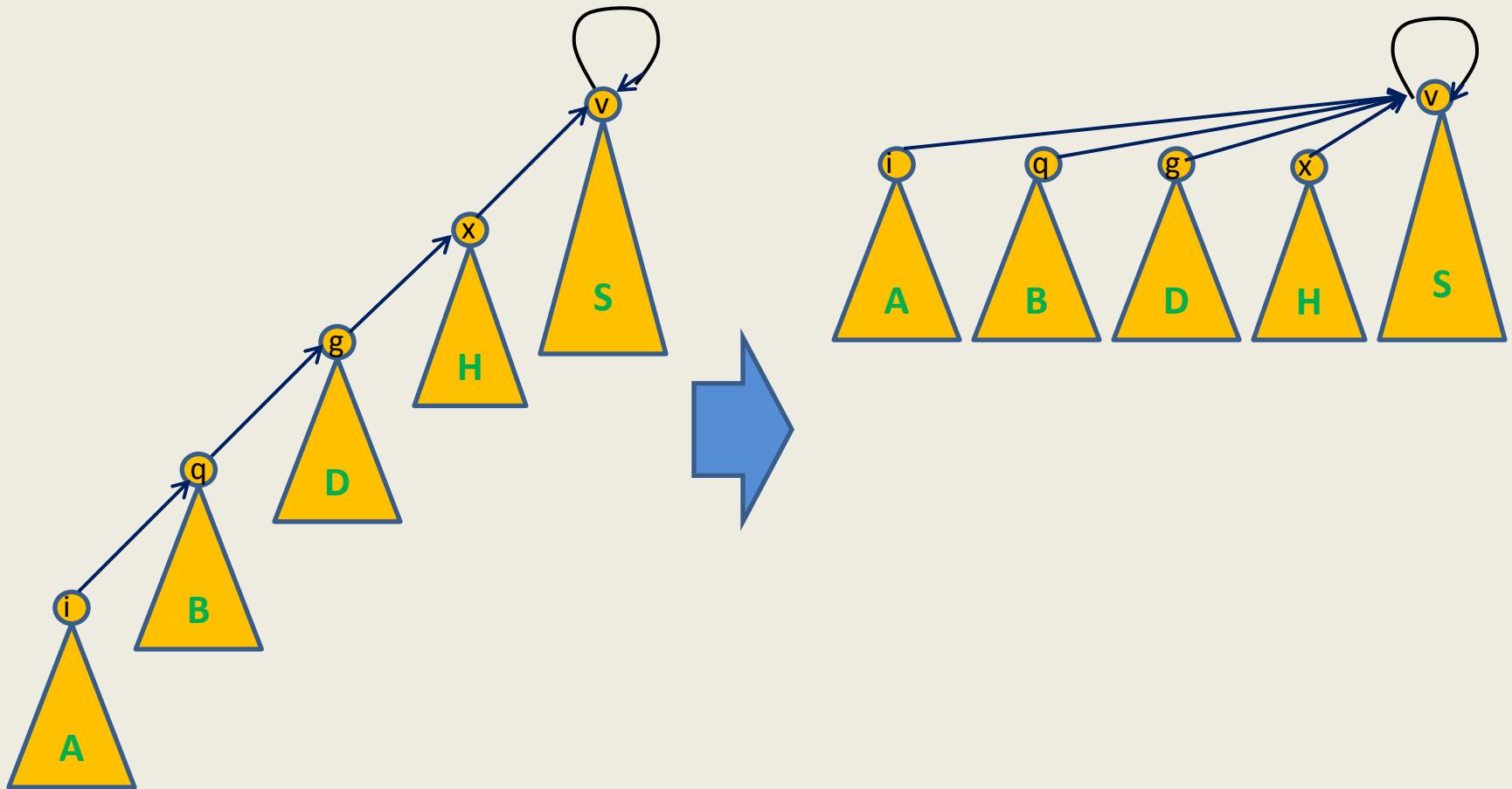
A clever programmer did the following modification to the code of **Find(*i*)**.

While executing **Find(*i*)**, we traverse the path from node *i* to the root. Let  $v_1, v_2, \dots, v_t$ , be the nodes traversed with  $v_t$  being the root node. At the end of **Find(*i*)**, if we update parent of each  $v_k$ ,  $1 \leq k < t$ , to  $v_t$ , we achieve a reduction in depth of many nodes. This modification increases the time complexity of **Find(*i*)** by at most a constant factor. But this little modification increased the overall speed of the application very significantly.

The heuristic is called **path compression**. It is shown pictorially on the following slide.

It remained a mystery for many years to provide a theoretical explanation for its practical success.

# Path compression during Find(i)



# Pseudocode for the modified Find

**Find(*i*)**

If (**Parent(*i*)** = *i*)    return *i* ;

else

*j*  $\leftarrow$  **Find(Parent(*i*));**

**Parent(*i*)**  $\leftarrow$  *j*;

return *j*

# Data Structures and Algorithms

(ESO207)

## Lecture 33

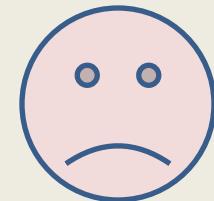
- Algorithm for  $i$ th order statistic of a set  $S$ .

# Problem definition

Given a set  $S$  of  $n$  elements  
compute  $i$ th smallest element from  $S$ .

**Applications:**

**Trivial algorithm:**



But sorting takes  $O(n \log n)$  time  
and appears to be an overkill  
for this simple problem.

**AIM:** To design an algorithm with  $O(n)$  time complexity.

**Assumption** (For the sake of **neat description** and **analysis** of algorithms of this lecture):

- All elements of  $S$  are assumed to be **distinct**.

# A motivational background

Though it was **intuitively appealing** to believe that there exists an  $O(n)$  time algorithm to compute  $i$ th smallest element, it remained a challenge for many years to design such an algorithm...

In **1972**, five well known researchers: **Blum, Floyd, Pratt, Rivest, and Tarjan** designed the  $O(n)$  time algorithm. It was designed during a **lunch break** of a conference when these five researchers sat together for the first time to solve the problem.

In this way, the problem which remained unsolved for many years got solved in less than an hour. But one should not ignore the efforts these researchers spent for years before arriving at the solution ... It was their effort whose fruit got ripened in that hour 😊.

# Notations

We shall now introduce some notations which will help in a neat description of the algorithm.

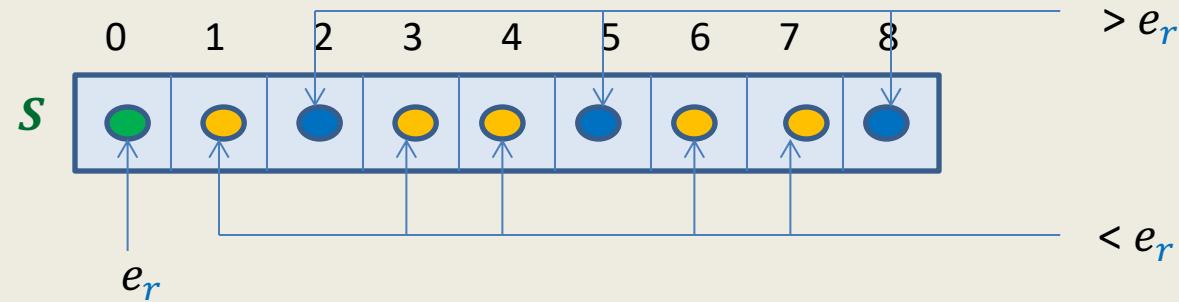
# Notations

- $S$  :  
the given set of  $n$  elements.
- $e_i$  :  
 $i$ th **smallest** element of  $S$ .
- $S_{<x}$  :  
subset of  $S$  consisting of all elements **smaller than**  $x$ .
- $S_{>x}$  :  
subset of  $S$  consisting of all elements **greater than**  $x$ .
- $\text{rank}(S,x)$  :  
 $1 +$  number of elements in  $S$  that are smaller than  $x$ .
- **Partition**( $S,x$ ):  
algorithm to partition  $S$  into  $S_{<x}$  and  $S_{>x}$ ;  
this algorithm returns  $(S_{<x}, S_{>x}, r)$  where  $r = \text{rank}(S,x)$ .

# Why should such an algorithm exist ?

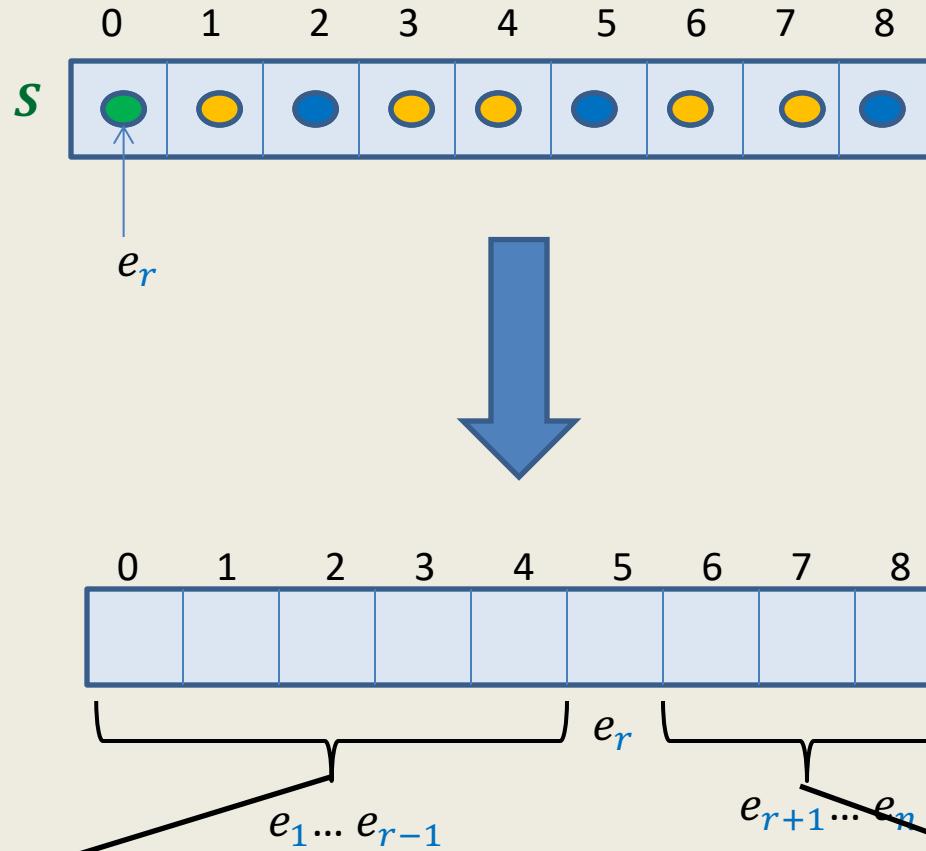
(inspiration from **QuickSort**)

# QuickSelect( $S, i$ )



What happens during  
Partition( $S, 0, 8$ )

# QuickSelect( $S, i$ )



If  $i > r$  we can discard these elements.

If  $i < r$  we can discard these elements.

# Pseudocode for QuickSelect( $S, i$ )

```
QuickSelect( $S, i$ )
{
    Pick an element  $x$  from  $S$ ;
     $(S_{<x}, S_{>x}, r) \leftarrow \text{Partition}(S, x);$ 
    If( $i = r$ ) return  $x$ ;
    Else If ( $i < r$ )
        QuickSelect( $S_{<x}, i$ )
    Else
        QuickSelect(  $S_{>x}, i - r$  );
}
```

Average case time complexity:  $O(n)$

Worst case time complexity :  $O(n^2)$

Analysis is simpler than Quick Sort.

**Towards worst case  $O(n)$  time algorithm ...**

# Key ideas

- Inspiration from some recurrences.
- Concept of approximate median
- Learning from QuickSelect( $S, i$ )

isn't it surprising that knowledge of recurrence can help in the design of an efficient algorithm? 😊

This is the usual trick:  
When a problem appears difficult, weaken the problem and try to solve it.

This is also a natural choice.  
Can we fine tune this algorithm to achieve our goal?

# Learning from recurrences

**Question:** what is the solution of recurrence  $T(n) = cn + T(9n/10)$  ?

**Answer:**  $O(n)$ .

Sketch (by gradual unfolding):

$$\begin{aligned}T(n) &= cn + c \frac{9n}{10} + c \frac{81n}{100} + \dots \\&= cn[1 + \frac{9}{10} + \frac{81}{100} + \dots] \\&= O(n)\end{aligned}$$

**Lesson 1 :**

Solution for  $T(n) = cn + T(an)$  is  $O(n)$  if  $0 < a < 1$ .

# Learning from recurrences

**Question:** what is the solution of recurrence

$$T(n) = cn + T(n/6) + T(5n/7) ?$$

**Answer:**  $O(n)$ .

**Sketch:** (by induction)

Assertion:  $T(n) \leq c_1 n$ .

**Induction step:**  $T(n) = cn + T(n/6) + T(5n/7)$

$$\leq cn + \frac{37}{42}c_1 n$$

$$\leq c_1 n \text{ if } c_1 \geq \frac{42}{5} c$$

**Lesson 2 :**

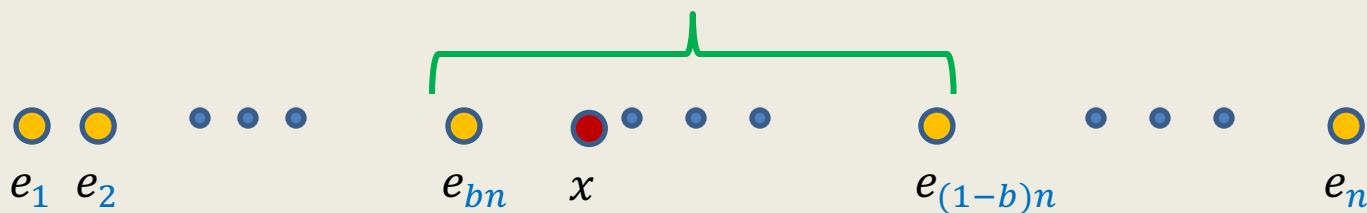
Solution for  $T(n) = cn + T(an) + T(dn)$  is  $O(n)$  if  $a + d < 1$ .

# Concept of approximate median

**Definition:** Given a constant  $0 < b \leq 1/2$ ,

an element  $x \in S$

if  $\text{rank}(x, S)$



# Learning from QuickSelect( $S, i$ )

QuickSelect( $S, i$ )

```
{   Pick an element  $x$  from  $S$ ;  
     $(S_{<x}, S_{>x}, r) \leftarrow \text{Partition}(S, x);$   
    If( $i = r$ ) return  $x$ ;  
    Else If ( $i < r$ )  
        QuickSelect( $S_{<x}, i$ )  
    Else  
        QuickSelect( $S_{>x}, i - r$ );  
}
```

Answer:  $T(n) = cn + T((1 - b)n)$   
 $= O(n)$

Lesson 1

$O(n)$

$T( (1 - b)n)$

What is time complexity  
of the algorithm

# Algorithm 2

**Select( $S, i$ )**

(A linear time algorithm)

# Overview of the algorithm

**Select( $S, i$ )**

```
{   Compute a  $b$ -approximate median, say  $x$ , of  $S$ ;      }  $O(n)$ 
     $(S_{<x}, S_{>x}, r) \leftarrow \text{Partition}(S, x);$           }  $O(n)$ 
    If( $i = r$ ) return  $x$ ;
    Else If ( $i < r$ )
        Select( $S_{<x}, i$ )
    Else
        Select( $S_{>x}, i - r$ );
}
```

**Observation:** If we can compute  $b$ -approximate median in  $O(n)$  time, we get  $O(n)$  time algo.  
But that appears too much to expect from us. Isn't it ?  
So what to do 😔 ?

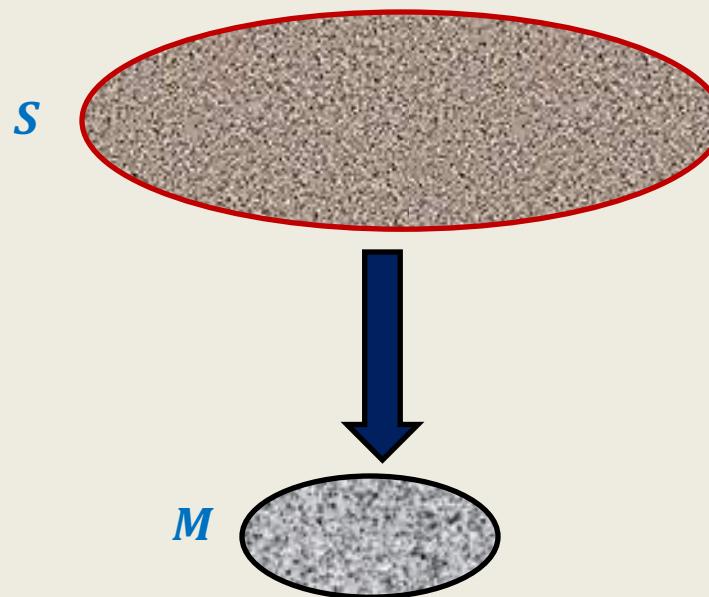
Hint: use **Lesson 2**

**Observation:** If we can compute  $b$ -approximate median  
time complexity of the algorithm will still be  $O(n)$ .

Spend some time on this observation to infer what it hints at.

**AIM:** How to compute a  $b$ -approximate median of  $S$

in  $O(n) + T(dn)$  time



**Question:** Can we form a set  $M$  of size  $dn$  such that  
exact median of  $M$  is  $b$ -approximate median of  $S$ ?

# Forming the subset $M$ with desired parameters

*This step forms the core of the algorithm and is indeed a brilliant stroke of inspiration. The student is strongly recommended to ponder over this idea from various angles.*

- Divide  $S$  into **groups** of **5** elements;
  - Compute median of each group by sorting;
  - Let  $M$  be the set of medians;
  - Compute median of  $M$ , let it be  $x$ ;
- $\left. \begin{matrix} \\ \\ \\ \end{matrix} \right\} O(n)$   
 $\left. \begin{matrix} \\ \end{matrix} \right\} T(n/5)$

**Question:** Is  $x$  an approximate median of  $S$  ?

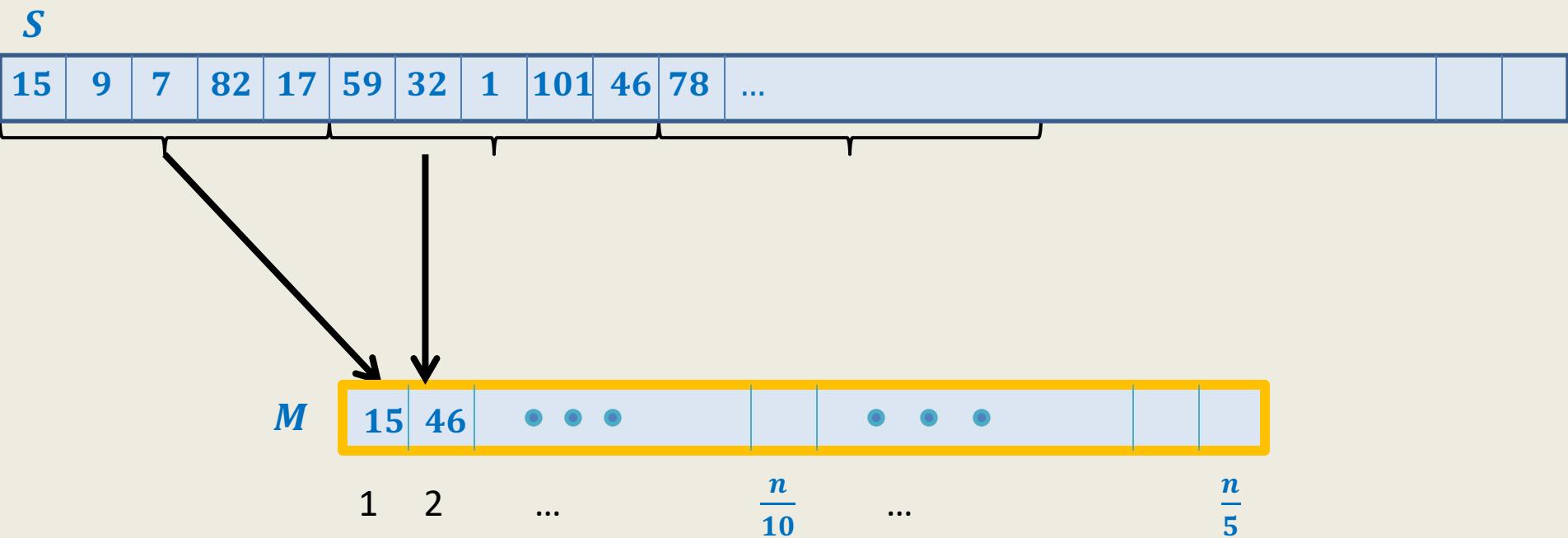
**Answer:** indeed.

The rank of  $x$  in  $M$  is  $n/10$ . Each element in  $M$  has two elements smaller than itself in its respective group. Hence there are at least  $\frac{3n}{10} - 1$  elements in  $S$  which are **smaller** than  $x$ .

In a similar way, there are at least  $\frac{3n}{10} - 1$  elements in  $S$  which are **greater** than  $x$ . Hence,  $x$  is  $\frac{3n}{10}$ -approximate median of  $S$ .

(See the animation on the following slide to get a better understanding of this explanation.)

# Forming the subset $M$



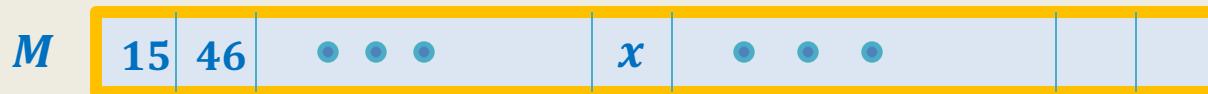
- Divide  $S$  into **groups** of 5 elements;
- Compute median of each group by sorting;

]}  $O(n)$

# Forming the subset $M$

$S$

15	9	7	82	17	59	32	1	101	46	78	...



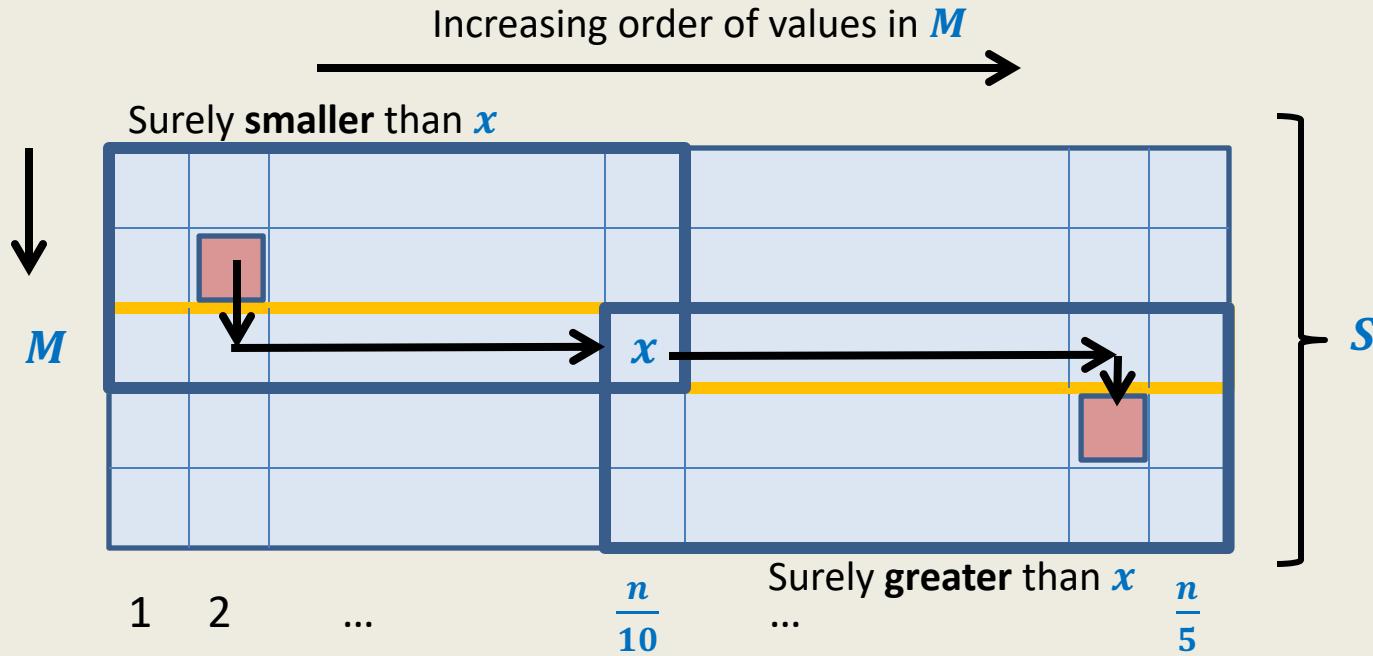
- Divide  $S$  into **groups** of 5 elements;
- Compute median of each group by sorting;
- Let  $M$  be the set of medians;
- Let  $x$  be median of  $M$ .

What can we say about rank of  $x$  in  $S$  ?

*Spend some time to answer this question before moving ahead.*

# Forming the subset $M$

Bring back the remaining 4 elements associated with each element of  $M$



→  $x$  is  $\left(\frac{3n}{10}\right)$  –approximate median of  $S$ .

Time required to form  $M$  :  $O(n)$

# Pseudocode for $\text{Select}(S, i)$

$\text{Select}(S, i)$

$M \leftarrow \emptyset;$

Divide  $S$  into **groups** of 5 elements;

Sort each group and add its **median** to  $M$ ;

$x \leftarrow \text{Select}(M, |M|/2);$

$(S_{<x}, S_{>x}, r) \leftarrow \text{Partition}(S, x);$

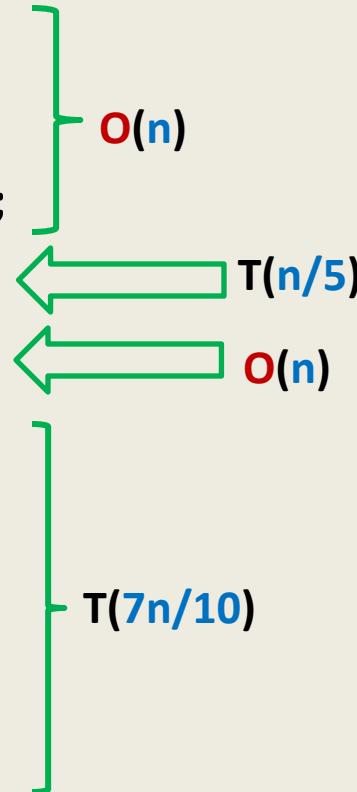
If( $i = r$ ) return  $x$ ;

Else If ( $i < r$ )

$\text{Select}(S_{<x}, i)$

Else

$\text{Select}(S_{>x}, i - r);$



# Analysis

$$\begin{aligned} T(n) &= cn + T(n/5) + T(7n/10) \\ &= O(n) \quad [\text{Learning from Recurrence of type II}] \end{aligned}$$

**Theorem:** Given any  $S$  of  $n$  elements, we can compute  $i$ th smallest element from  $S$  in  $O(n)$  worst case time.

# Exercises

(Attempting these exercises will give you a better insight into the algorithm.)

- What is magical about number **5** in the algorithm ?
- What if we divide the set **S** into groups of size **3** ?
- What if we divide the set **S** into groups of size **7** ?
- What if we divide the set **S** into groups of even size (e.g. **4** or **6**) ?

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 34**

- A new algorithm design paradigm: Greedy strategy  
part I

# Path to the solution of a problem



But there is a **systematic approach** which usually works 😊

**Today's lecture will demonstrate this approach 😊**

## **Problem : JOB Scheduling**

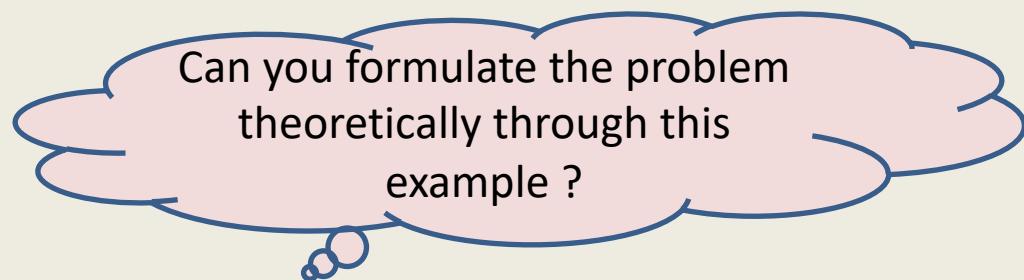
**Largest subset of non-overlapping job**

# A motivating example

## Antaragni 2021

- There are ***n*** large-scale activities to be performed in Auditorium.
- Each large scale activity has a **start time** and **finish time**.
- There is **overlap** among various activities.

**Aim:** What is the largest subset of activities that can be performed ?



## Formal Description

# A job scheduling problem

### INPUT:

- A set  $J$  of  $n$  jobs  $\{j_1, j_2, \dots, j_n\}$
- job  $j_i$  is specified by two real numbers
  - $s(i)$ : start time of job  $j_i$
  - $f(i)$ : finish time of job  $j_i$
- A single server

### Constraints:

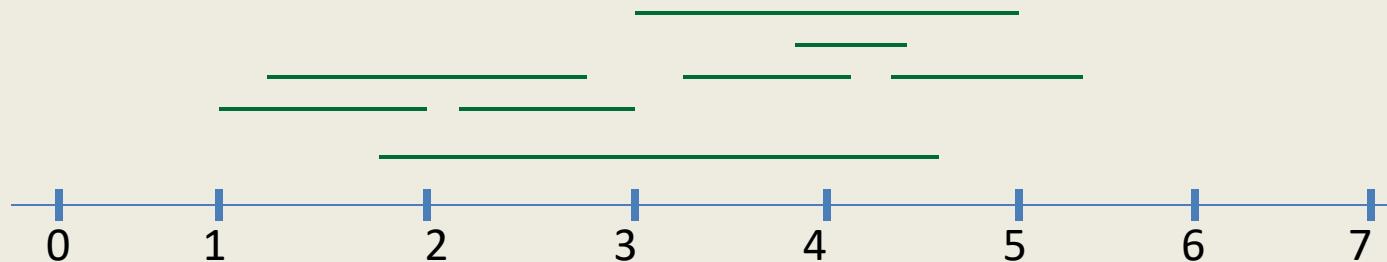
- Server can execute at most one job at any moment of time.
- Job  $j_i$ , if scheduled,

### Aim:

To select the **largest** subset of non-overlapping jobs which can be executed by the server.

## Example

INPUT:  $(1, 2)$ ,  $(1.2, 2.8)$ ,  $(1.8, 4.6)$ ,  $(2.1, 3)$ ,  $(3, 5)$ ,  $(3.3, 4.2)$ ,  $(3.9, 4.4)$ ,  $(4.3, 5.4)$



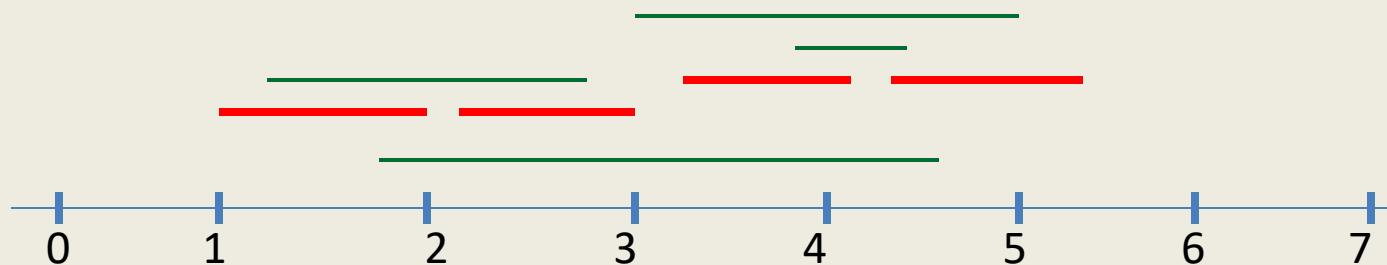
It makes sense to work with pictures than these numbers😊

job  $i$  is said to be **non-overlapping** with job  $k$  if

Try to find solution for the above example.

## Example

INPUT:  $(1, 2), (1.2, 2.8), (1.8, 4.6), (2.1, 3), (3, 5), (3.3, 4.2), (3.9, 4.4), (4.3, 5.4)$



job  $i$  is said to be **non-overlapping** if  $\text{f}(i) \cap \text{f}(k) = \emptyset$

What strategy come  
to your mind?

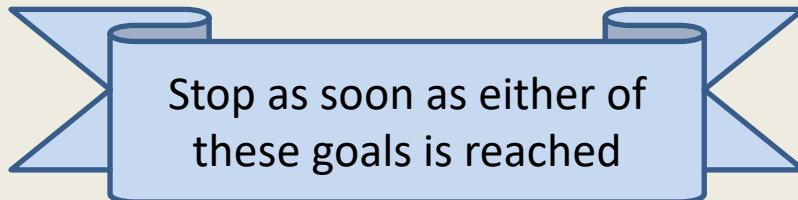
# Designing algorithm for any problem

1. Choose a strategy based on some intuition
2. Transform the strategy into an algorithm.

Try to prove  
correctness  
of the  
algorithm

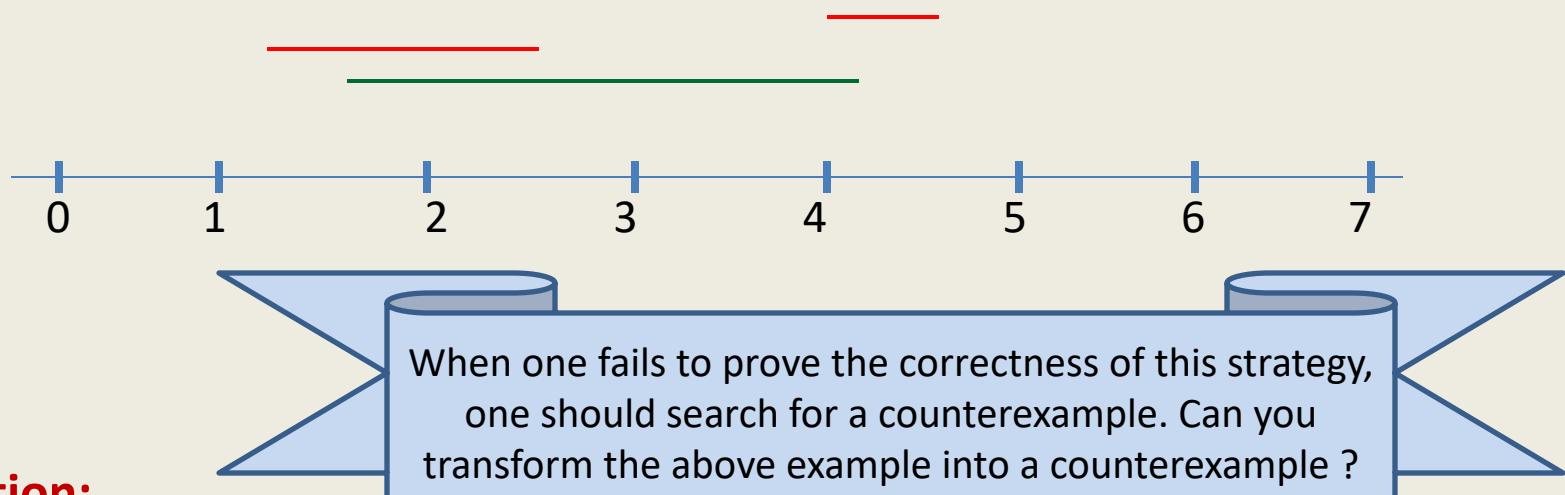


Try to design a  
conterexample



# Designing algorithm for the problem

Strategy 1: Select the earliest start time job

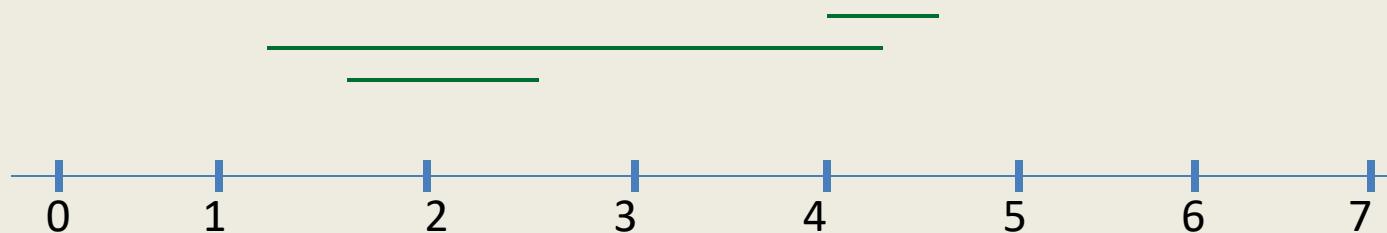


Intuition:

It might be better to assign jobs as early as possible so as to make optimum use of server.

# Designing algorithm for the problem

**Strategy 1:** Select the earliest start time job

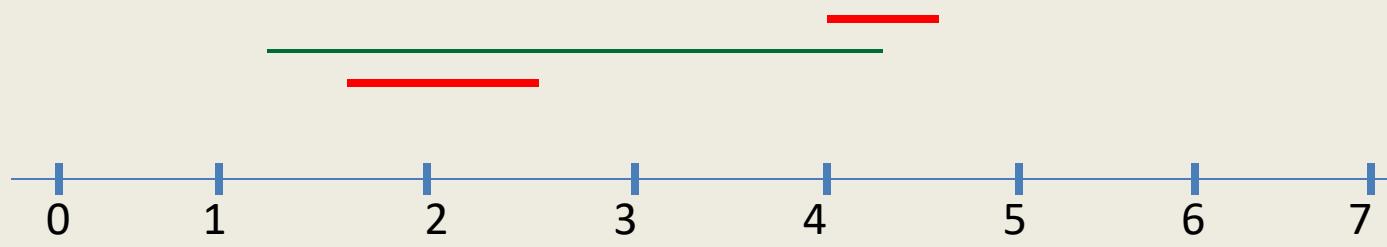


## Intuition:

It might be better to assign jobs as early as possible so as **to make optimum use of server**.

# Designing algorithm for the problem

**Strategy 1:** Select the earliest start time job



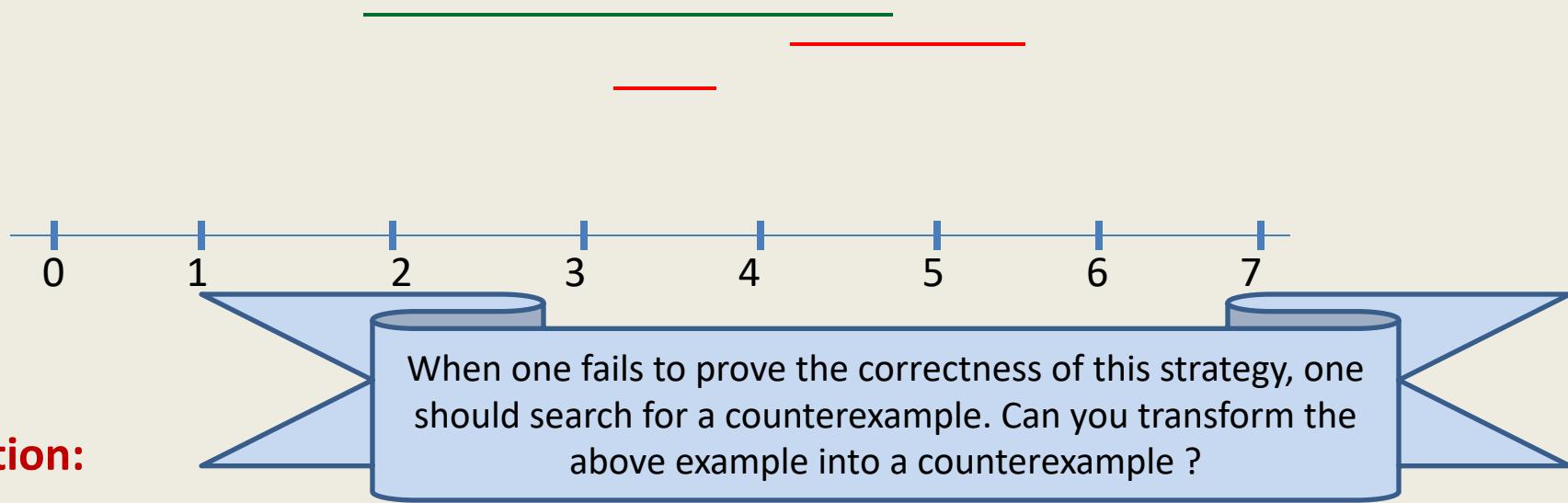
**Intuition:**

Instead of getting disappointed, try to realize that this counterexample points towards some other strategy which might work.

It might be better to assign jobs as early as possible so as to make optimum use of server.

# Designing algorithm for the problem

**Strategy 2:** Select the job with **smallest duration**



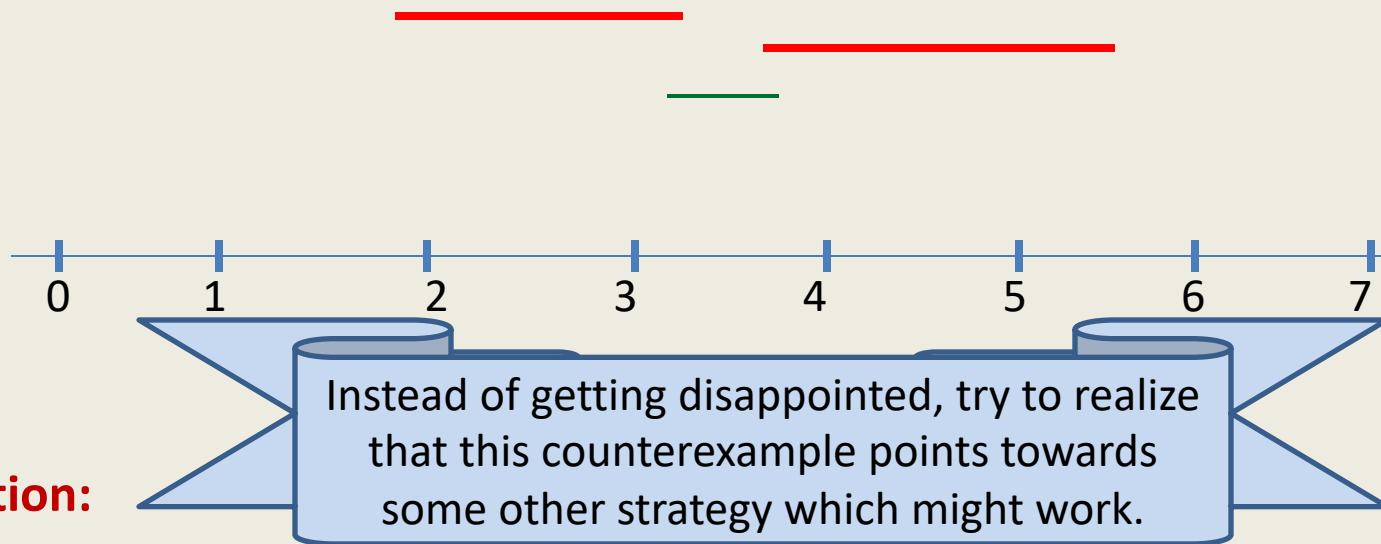
**Intuition:**

Such a job will make **least use of the server**

→ this might lead to larger number of jobs to be executed

# Designing algorithm for the problem

**Strategy 2:** Select the job with **smallest duration**



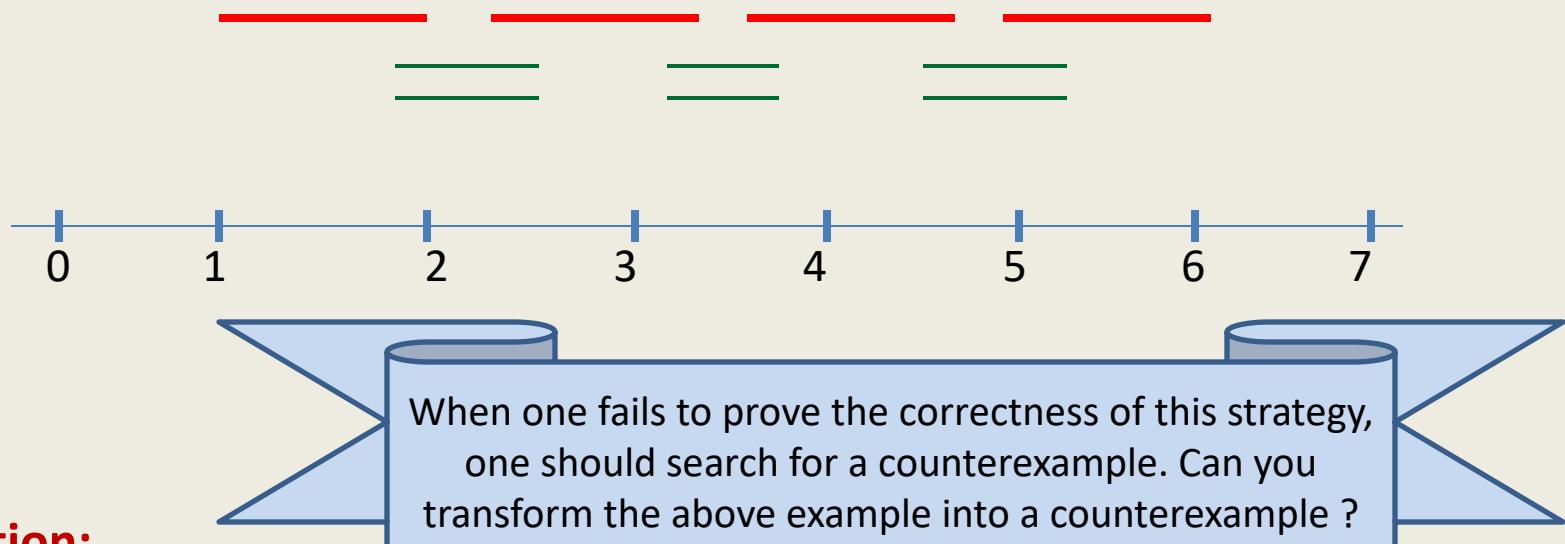
**Intuition:**

Such a job will make **least use of the server**

→ this might lead to larger number of jobs to be executed

# Designing algorithm for the problem

**Strategy 3:** Select the job with **smallest no. of overlaps**

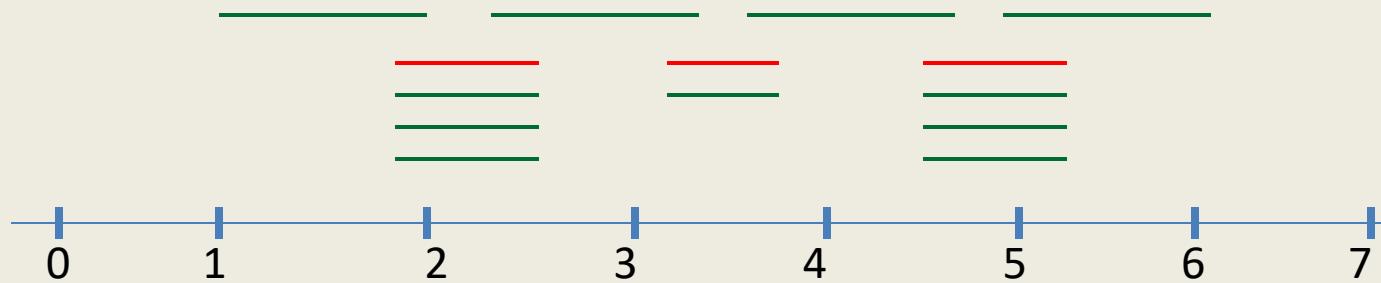


**Intuition:**

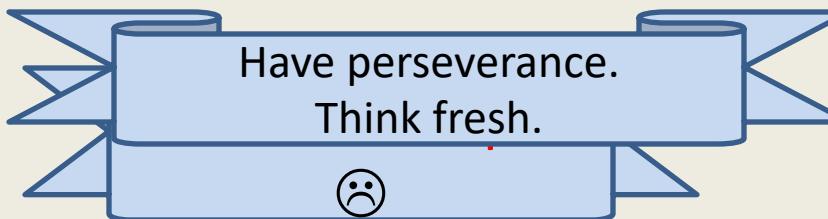
Selecting such a job will result in **least number of other jobs to be discarded**.

# Designing algorithm for the problem

**Strategy 3:** Select the job with **smallest no. of overlaps**



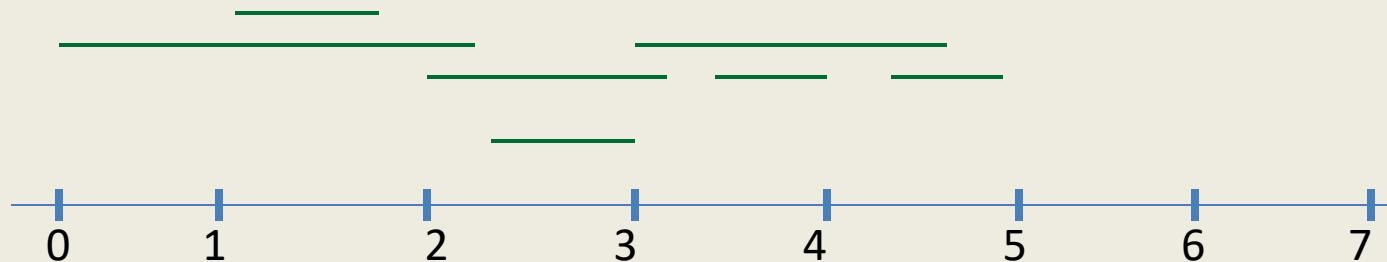
**Intuition:**



Selecting such a job will result in **least number of other jobs to be discarded**.

# Designing algorithm for the problem

**Strategy 4:** Select the job with **earliest finish time**



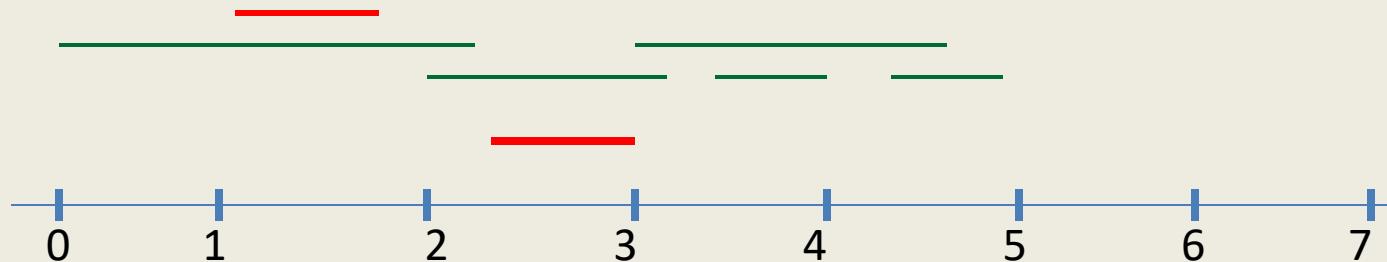
## Intuition:

Selecting such a job will **free** the server **earliest**

→ hence more no. of jobs might get scheduled.

# Designing algorithm for the problem

**Strategy 4:** Select the job with **earliest finish time**



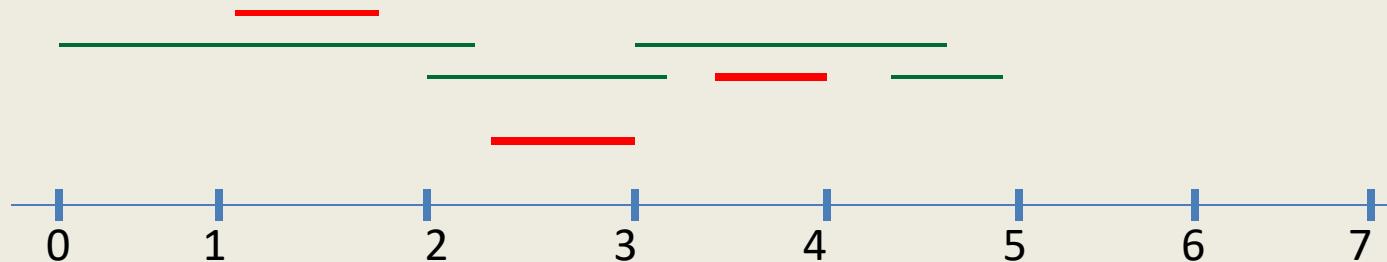
## Intuition:

Selecting such a job will **free** the server **earliest**

→ hence more no. of jobs might get scheduled.

# Designing algorithm for the problem

**Strategy 4:** Select the job with **earliest finish time**



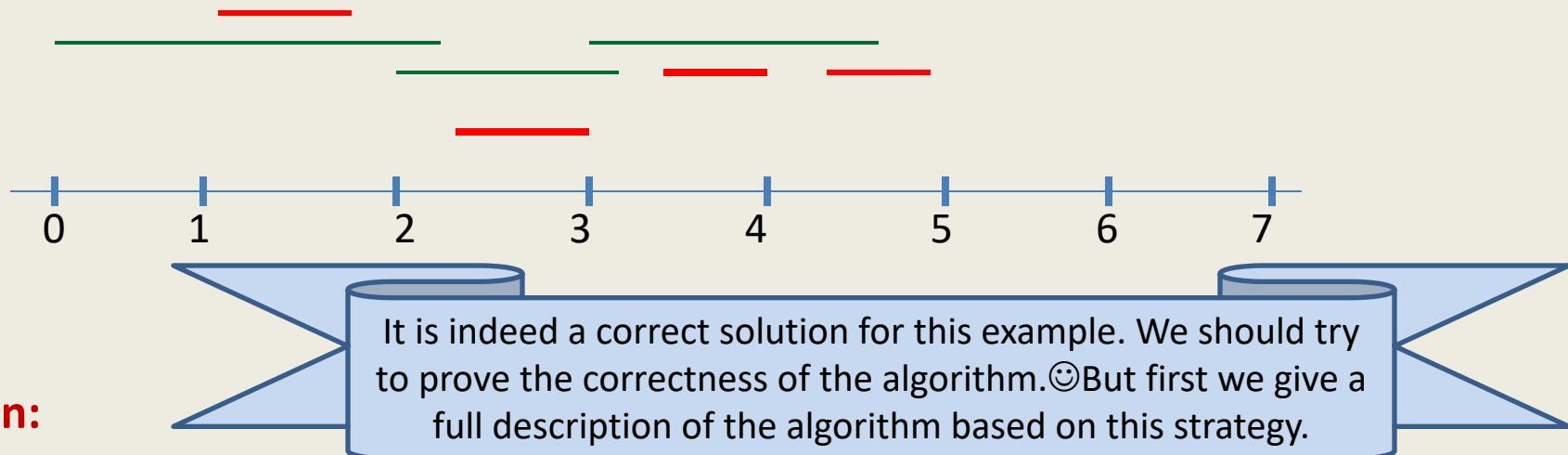
## Intuition:

Selecting such a job will **free** the server **earliest**

→ hence more no. of jobs might get scheduled.

# Designing algorithm for the problem

**Strategy 4:** Select the job with earliest finish time



**Intuition:**

Selecting such a job will **free** the server **earliest**

→ hence more no. of jobs might get scheduled.

# Algorithm “earliest finish time”

## Description

**Algorithm** (Input : set  $J$  of  $n$  jobs.)

1. Define  $A \leftarrow \emptyset$ ;
2. While  $J \neq \emptyset$  do
  - { Let  $x$  be the job from  $J$  with earliest finish time;  
 $A \leftarrow A \cup \{x\}$ ;  
Remove  $x$  and all jobs that overlap with  $x$  from set  $J$ ;
  - }
3. Return  $A$ ;

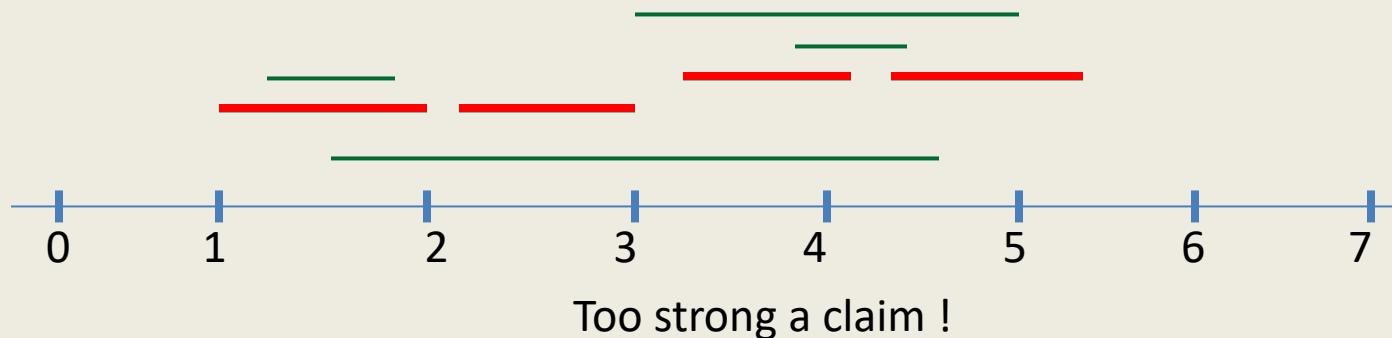
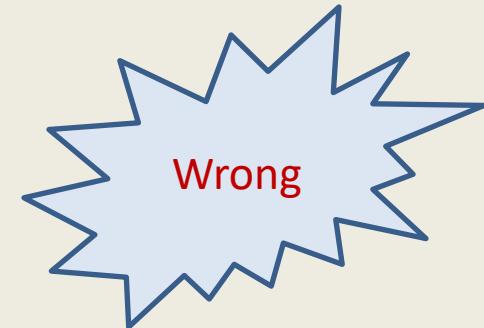
Running time for a trivial implementation of the above algorithm:  $O(n^2)$

# Algorithm “earliest finish time”

## Correctness

Let  $x$  be the job with earliest finish time.

**Lemma1:**  $x$  is present in the optimal solution for  $J$ .



# Algorithm “earliest finish time”

## Correctness

Let  $x$  be the job with earliest finish time.

**Lemma1:** There exists an optimal solution for  $J$  in which  $x$  is present.

**Proof:** Consider any optimal solution  $O$  for  $J$ .

Let  $y$  be the job from  $O$  with earliest finish time.

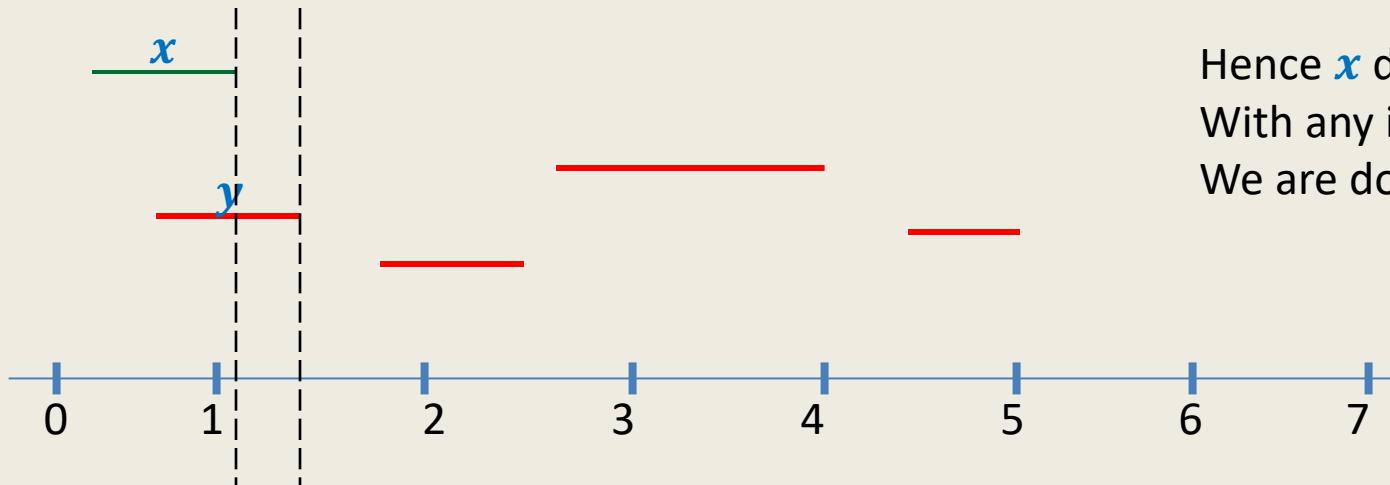
Let  $O' \leftarrow O \setminus \{y\}$

$$\rightarrow f(y) < s(z) \forall z \in O'$$

$$\rightarrow f(x) \leq f(y)$$

$O' \cup \{x\}$  is also an optimal solution.

**Reason:**  $O' \cup \{x\}$  has no overlapping intervals. Give arguments.



Hence  $x$  does not overlap  
With any interval of  $O'$ .  
We are done 😊

# Homework

Spend 30 minutes today on the following problems.

1. Use **Lemma1** to complete the proof of correctness of the algorithm.
2. Design an  **$O(n \log n)$**  implementation of the algorithm.

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 35**

- A new algorithm design paradigm: Greedy strategy  
part II

# **Continuing Problem from last class**

**JOB Scheduling**

**Largest subset of non-overlapping job**

# A job scheduling problem

## Formal Description

### INPUT:

- A set  $J$  of  $n$  jobs  $\{j_1, j_2, \dots, j_n\}$
- job  $j_i$  is specified by two real numbers
  - $s(i)$ : start time of job  $j_i$
  - $f(i)$ : finish time of job  $j_i$
- A single server

### Constraints:

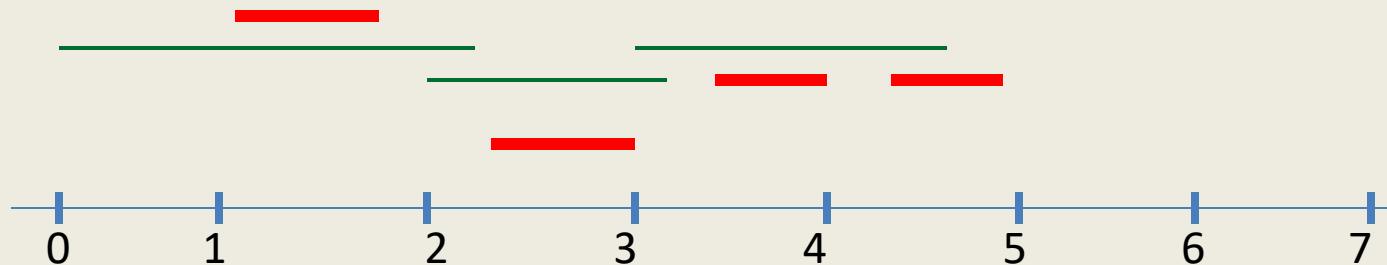
- Server can execute at most one job at any moment of time and a job.
- Job  $j_i$ , if scheduled, has to be scheduled during  $[s(i), f(i)]$  only.

### Aim:

To select the **largest** subset

# Designing algorithm for the problem

**Strategy 4:** Select the job with **earliest finish time**



## Intuition:

Selecting such a job will **free** the server **earliest**

→ hence more no. of jobs might get scheduled.

# Algorithm “earliest finish time”

Proof of correctness ?

Let  $x \in J$  be the job with earliest finish time.

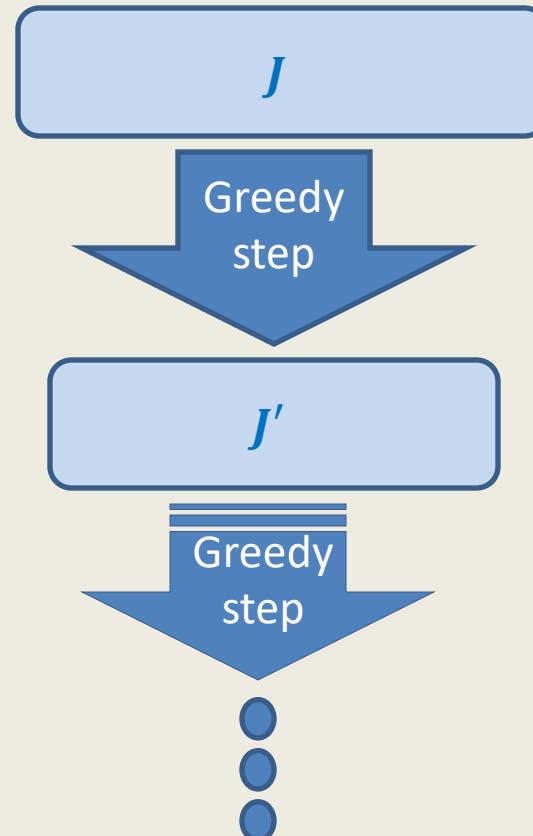
Let  $J' = J \setminus \text{Overlap}(x)$

**Algorithm** (Input : set  $J$  of  $n$  jobs.)

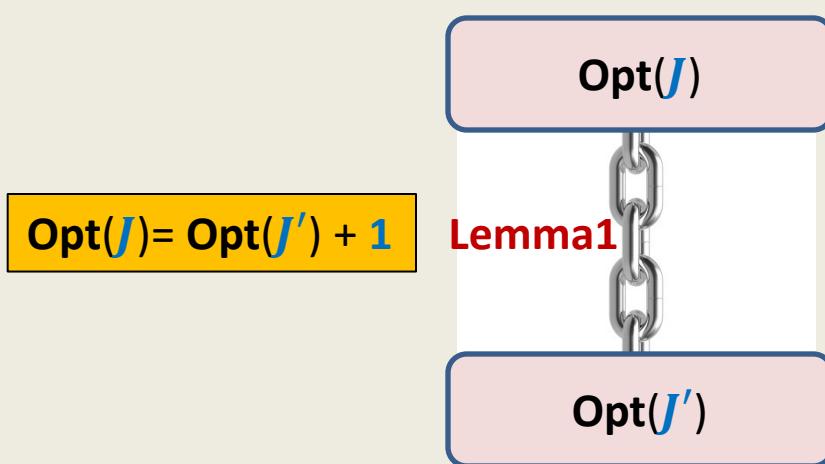
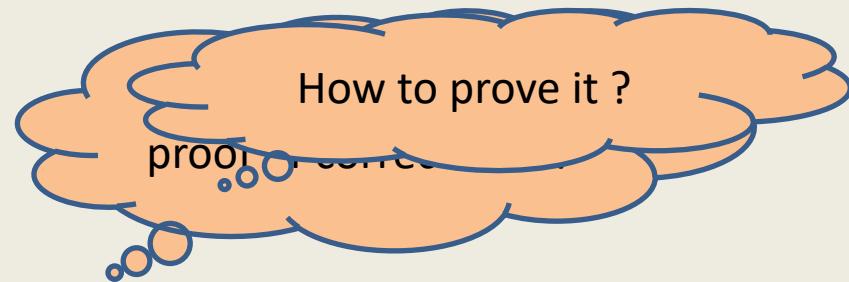
1. Define  $A \leftarrow \emptyset$ ;
2. While  $J \neq \emptyset$  do
  - { Let  $x \in J$  has earliest finish time;  
 $A \leftarrow A \cup \{x\}$ ;  
 $J \leftarrow J \setminus \text{Overlap}(x)$ ;
  - }
3. Return  $A$ ;

**Lemma1 (last class):**

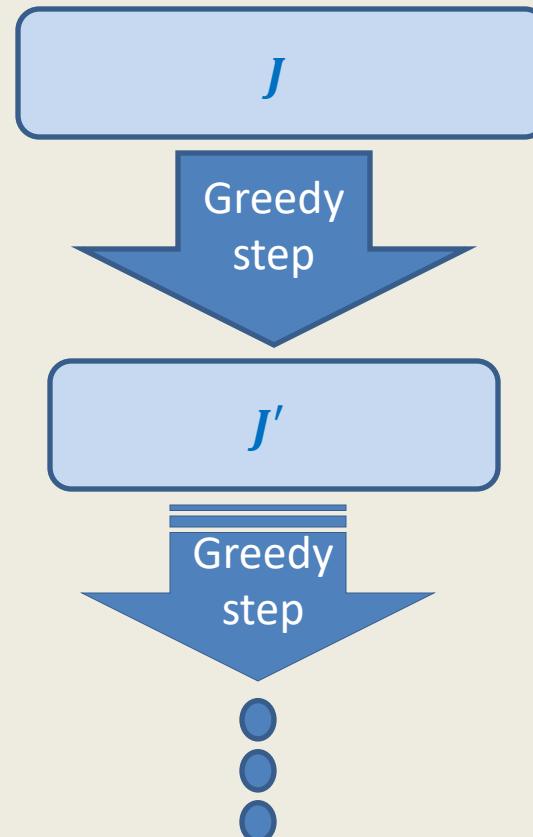
There exists an optimal solution for  $J$  containing the **earliest finish time** job.



# Algorithm “earliest finish time”



**Proof of correctness ?**  
Let  $x \in J$  be the job with earliest finish time.  
Let  $J' = J \setminus \text{Overlap}(x)$



**Notation:**

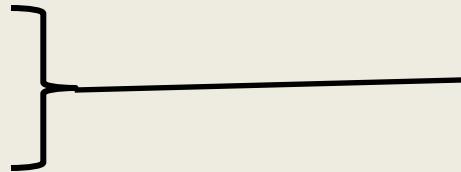
$\text{Opt}(J)$ :

# Theorem: $\text{Opt}(J) = \text{Opt}(J') + 1$ .

- Proof has two parts

$$\text{Opt}(J) \geq \text{Opt}(J') + 1$$

$$\text{Opt}(J') \geq \text{Opt}(J) - 1$$



Try to give a physical interpretation to these inequalities.

- Proof for each part is a proof **by construction**

# Algorithm “earliest finish time”

Proving  $\text{Opt}(J) \geq \text{Opt}(J') + 1$ .

**Observation:** start time of every job in  $J'$  is greater than finish time of  $x$ .

Let  $O'$  be any optimal solution for  $J'$ .

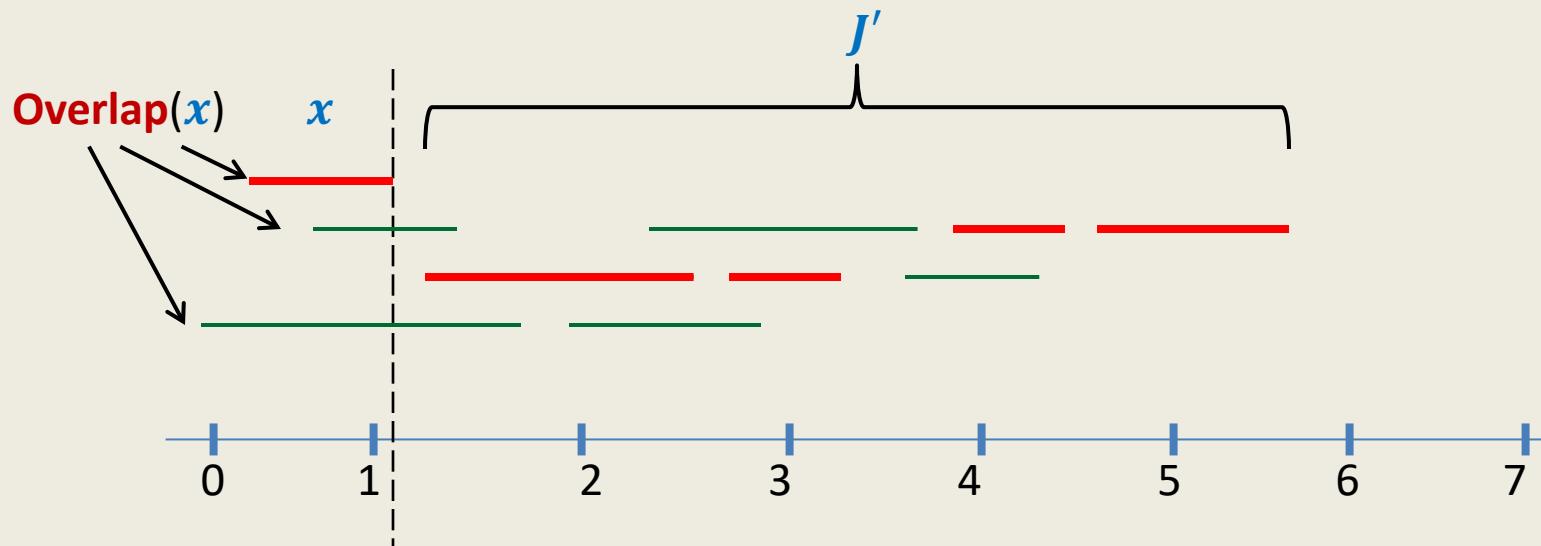
None of the jobs in

From an optimal solution of  $J'$

Hence  $O' \cup \{x\}$  is a

can you derive a solution for  $J$  with one extra job?

Therefore  $\text{Opt}(J) \geq |O'| + 1$



# Algorithm “earliest finish time”

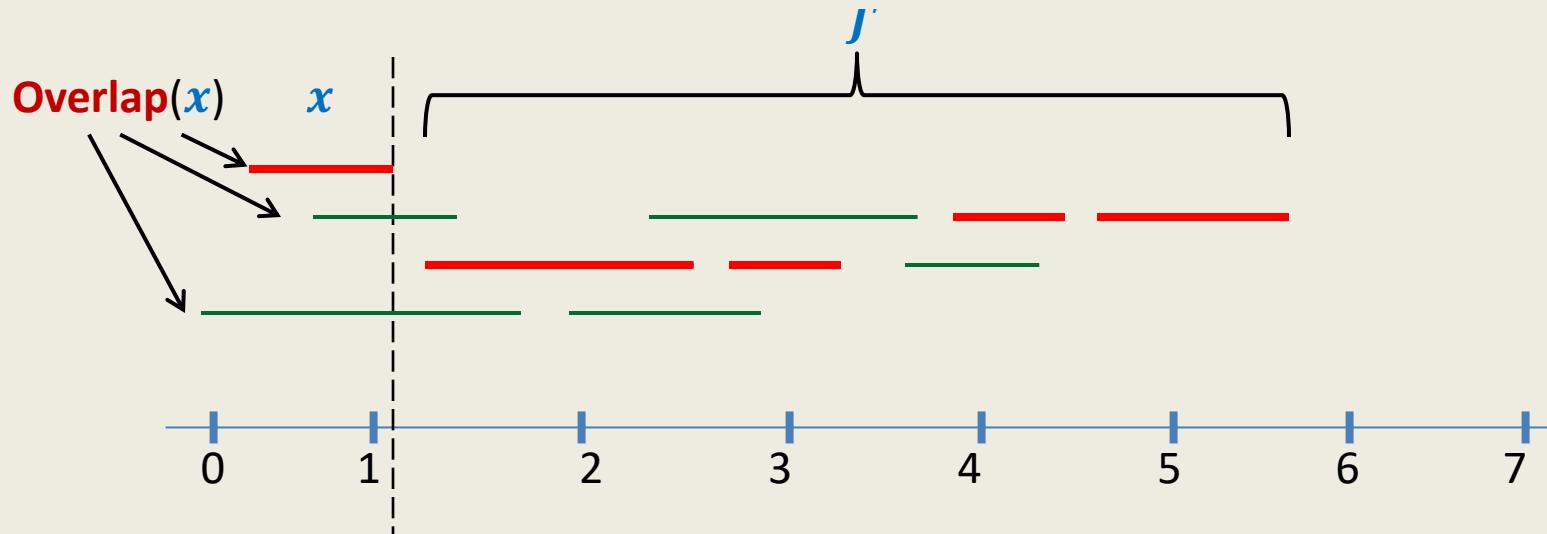
**Lemma1 (last class):** There exists an optimal solution for  $J$  in which  $x$  is present.

Let  $O$  be an optimal solution for  $J$  containing  $x$ .

None of the jobs in  $O$  overlaps with  $x$ .

→ Every job from  $O$  can you derive a solution for  $J'$  with one job less?  
Hence  $O \setminus \{x\}$  is a subset of non-overlapping jobs from  $J$ .

Therefore  $\text{Opt}(J') \geq |O| - 1$ :



## Theorem:

Given any set  $J$  of  $n$  jobs,

the algorithm based on “earliest finish time” approach  
computes the largest subset of non-overlapping job.

# $O(n \log n)$ implementation of the Algorithm

This is not the only way to achieve  $O(n \log n)$  time. It can be done other ways as well.

**Algorithm (Input : set  $J$  of  $n$  jobs.)**

1. Define  $A \leftarrow \emptyset$ ;
2. While  $J \neq \emptyset$  do

{ Let  $x \in J$  have earliest finish time;

$A \leftarrow A \cup \{x\}$ ;

$J \leftarrow J \setminus \text{Overlap}(x)$ ;

}

3. Return  $A$ ;

Maintain a **binary min-heap** for  $J$  based on **finish time** as the key.

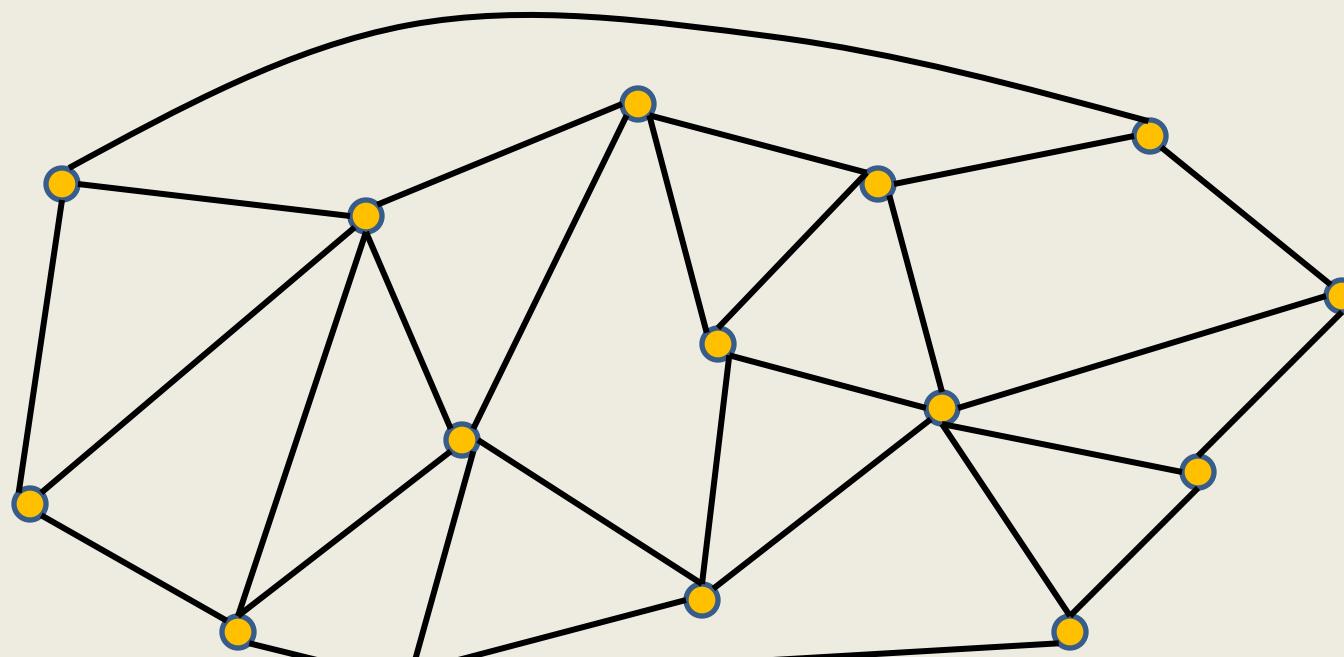
Sort  $J$  in increasing order of **start time**.

→  $O(n^2)$  time complexity is obvious

# **Problem 2**

**First we shall give motivation.**

# Motivation: A road or telecommunication network



Suppose there is a collection of possible links/roads that can be laid.  
But laying down each possible link/road is costly.

**Aim:** To lay down **least number** of links/roads to ensure **connectivity** between each pair of nodes/cities.

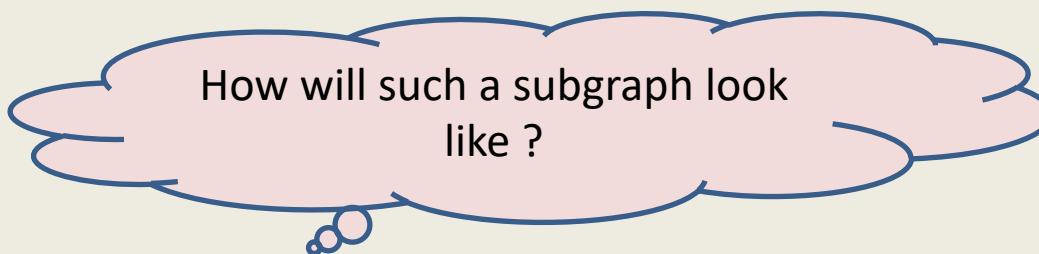
# Motivation

## Formal description of the problem

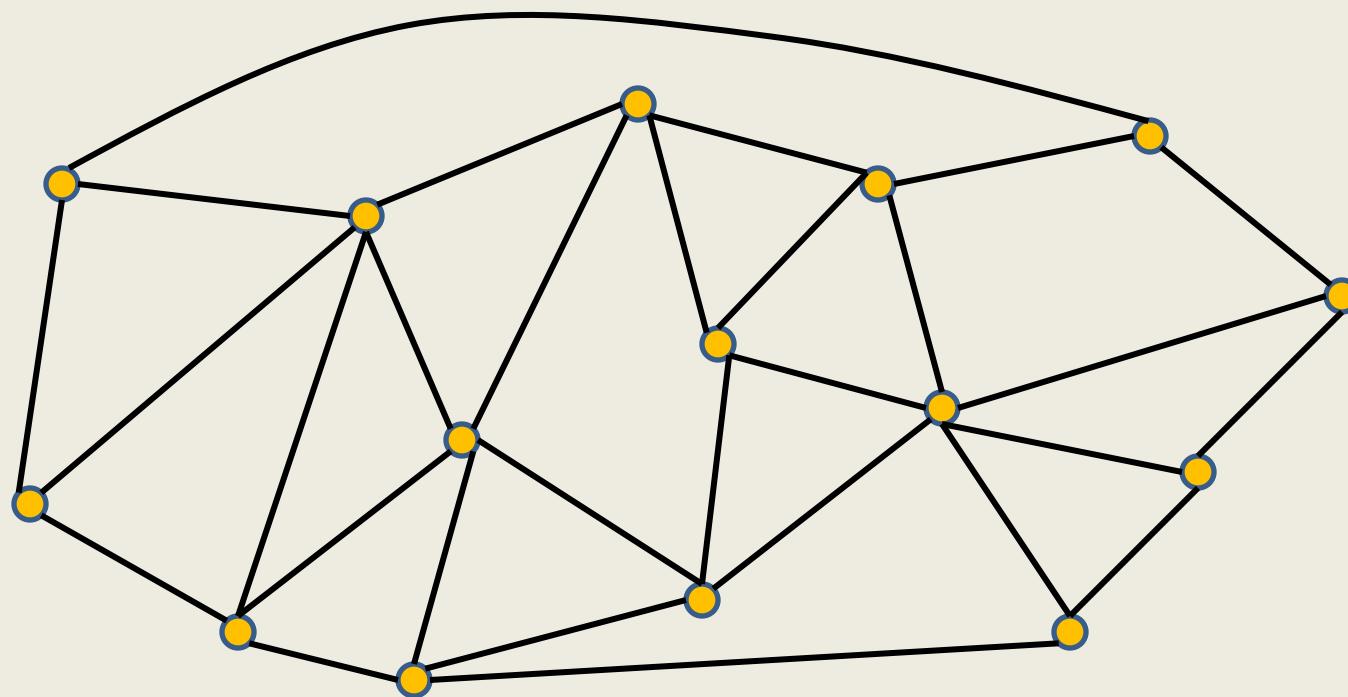
**Input:** an undirected graph  $\mathbf{G}=(\mathbf{V},\mathbf{E})$ .

**Aim:** compute a **subgraph**  $(\mathbf{V}',\mathbf{E}')$ ,  $\mathbf{E}' \subseteq \mathbf{E}$  such that

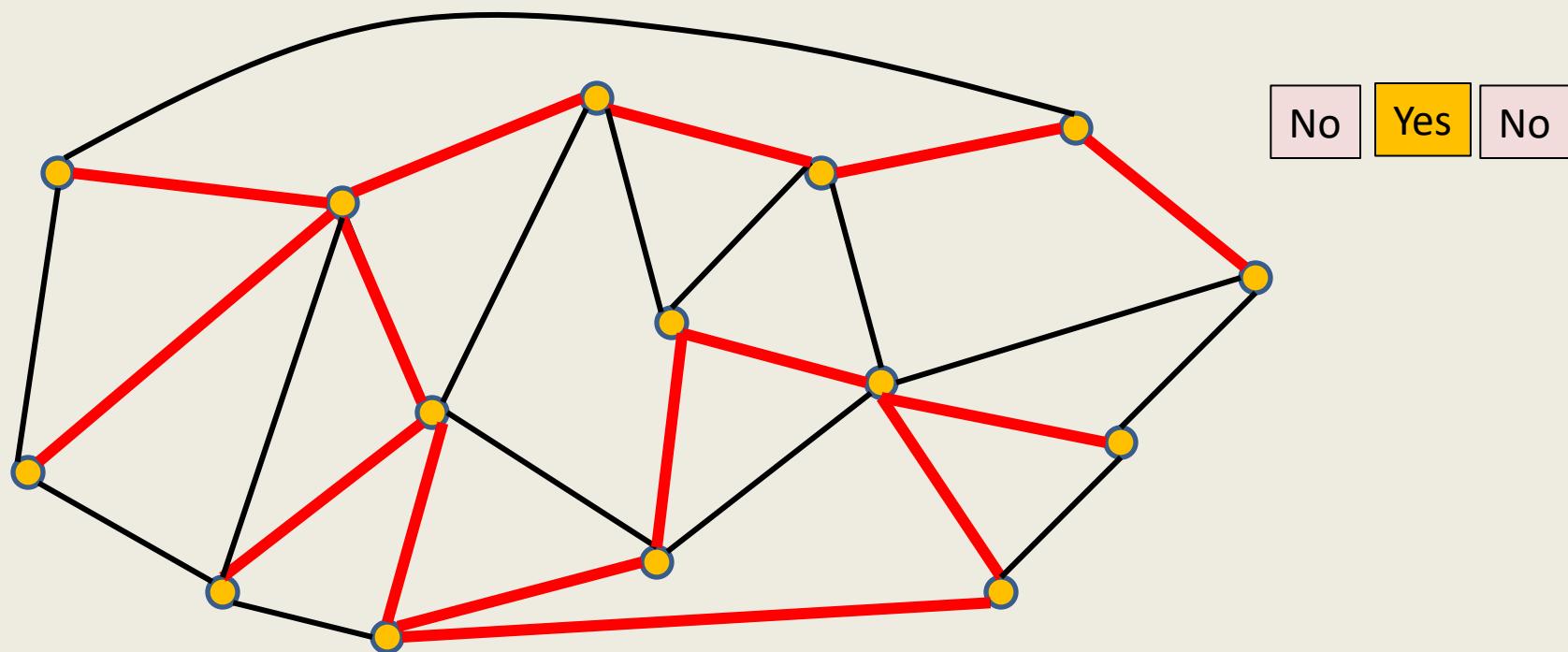
- **Connectivity** among all  $\mathbf{V}$  is guaranteed in the **subgraph**.
- $|\mathbf{E}'|$  is **minimum**.



# A road or telecommunication network



# A road or telecommunication network



Is this subgraph meeting our requirement ?

# A tree

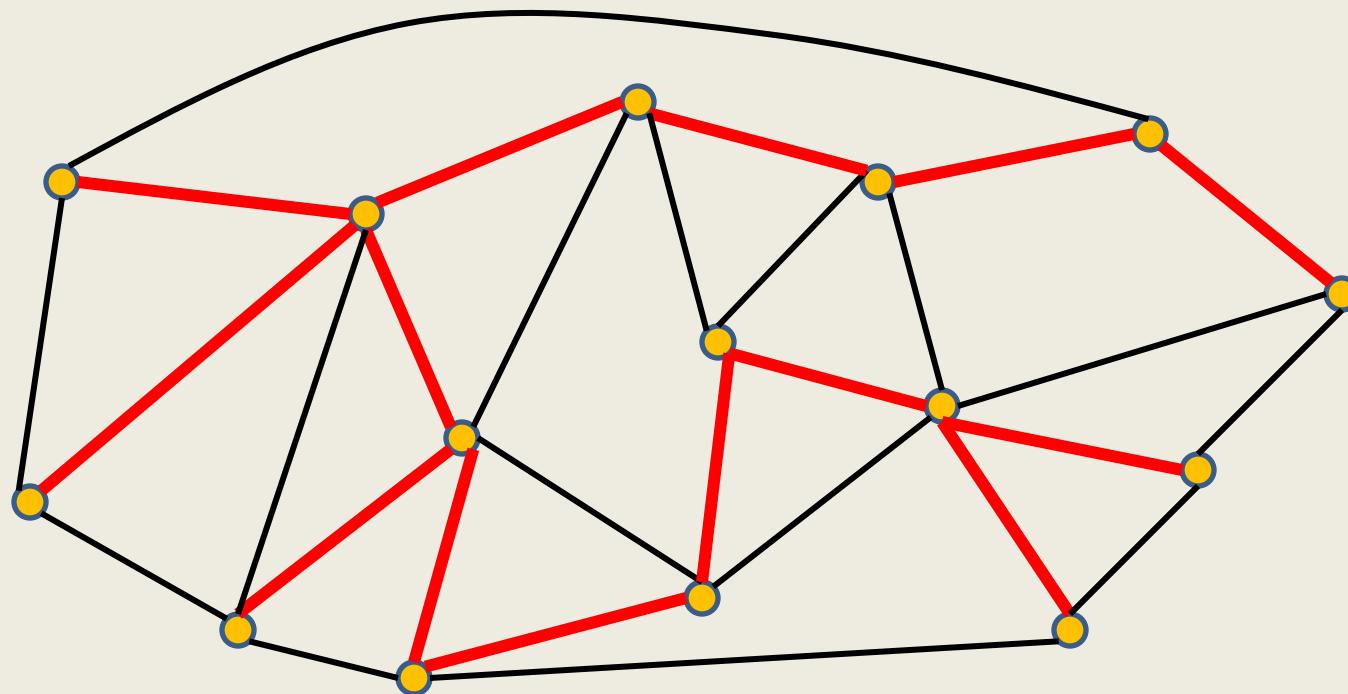
The following definitions are **equivalent**.

- An undirected graph which is **connected**
- An undirected graph where each pair of vertices has
- An undirected **connected** graph on  $n$  vertices and  $n - 1$  edges
- An undirected graph on  $n$  vertices and  $n - 1$  edges and

# A Spanning tree

**Definition:** For an undirected graph  $(V, E)$ ,

A spanning tree is a **subgraph**  $(V, E')$ ,  $E' \subseteq E$  which is a tree.

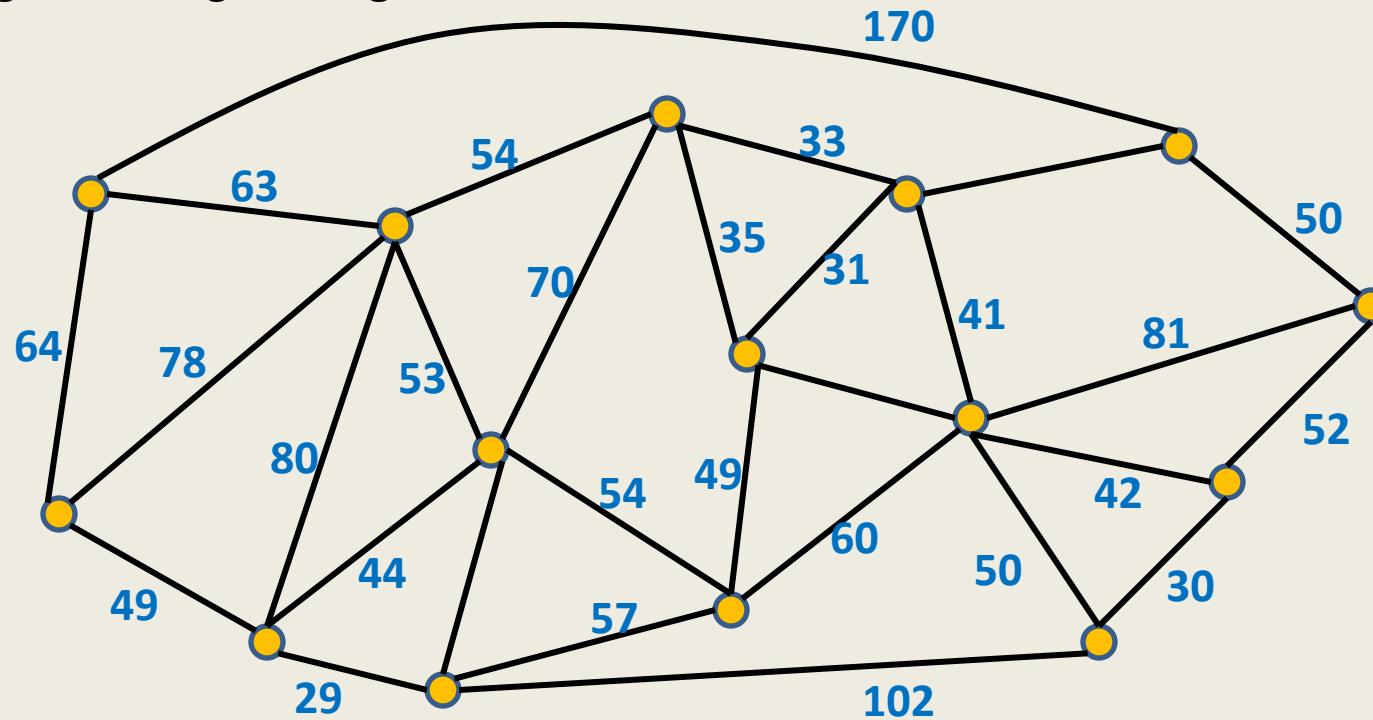


**Observation:** Given a spanning tree  $T$  of a graph  $G$ , adding a nontree edge  $e$  to  $T$  creates a unique cycle.

There will be total  $m - n + 1$  such cycles. These are called **fundamental cycles** in  $G$  induced by the spanning tree  $T$ .

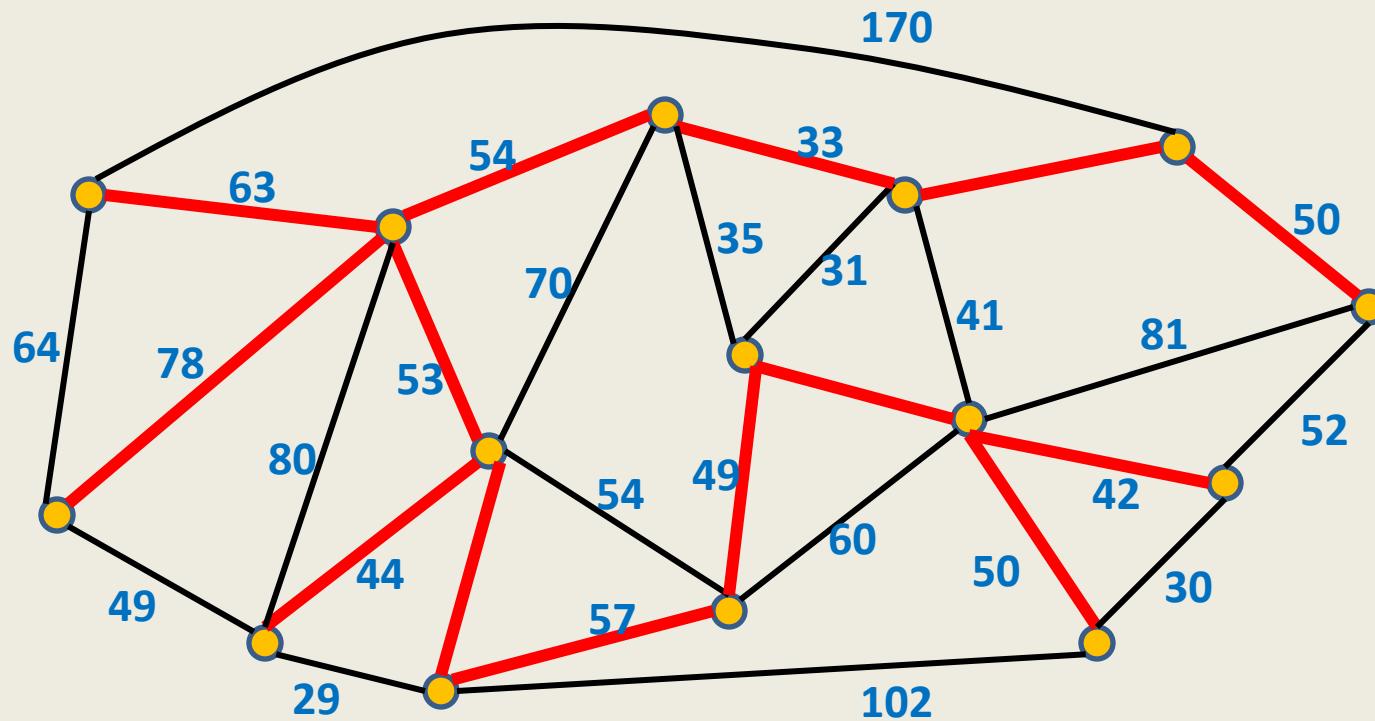
# A road or telecommunication network

Assign each edge a **weight/cost**.



Adding more reality to the problem

# A road or telecommunication network



Any arbitrary spanning tree (like the one shown above) will not serve our goal 😞.

We need to select the spanning tree with **least weight/cost**.

# **Problem 2**

**Minimum spanning tree**

# Problem Description

**Input:** an undirected graph  $\mathbf{G}=(\mathbf{V}, \mathbf{E})$  with  $w: \mathbf{E} \rightarrow \mathbb{R}$ ,

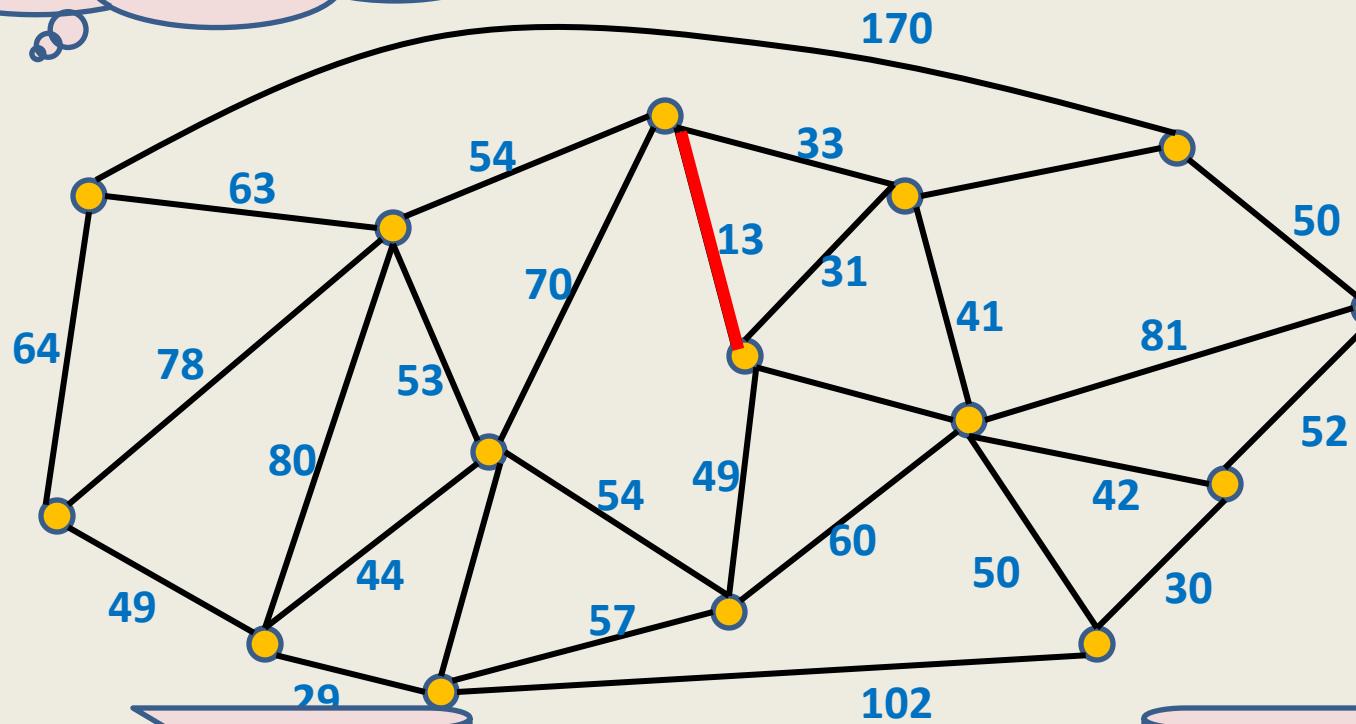
**Aim:** compute a **spanning tree**  $(\mathbf{V}, \mathbf{E}')$ ,  $\mathbf{E}' \subseteq \mathbf{E}$  such that

$\sum_{e \in \mathbf{E}'} w(e)$  is **minimum**.



# How to compute a MST ?

The least weight edge  
should be in MST.  
But why ?



Look at the graph. Is there any edge for which you feel strongly to be present in MST ?

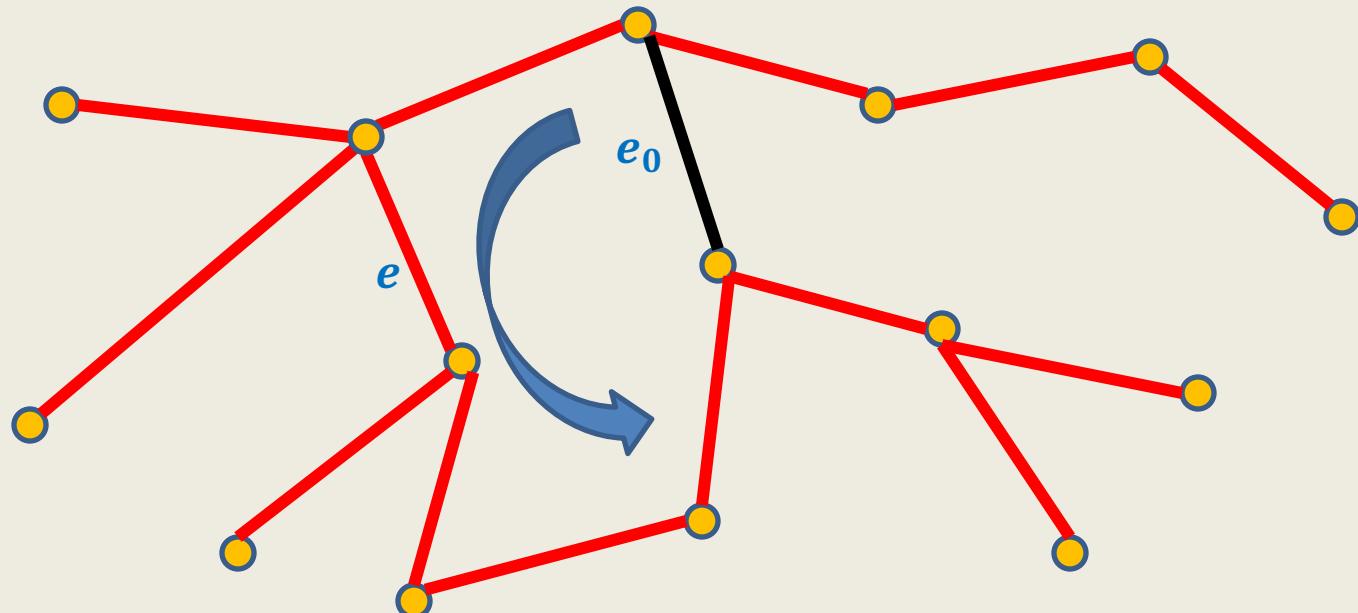
Let  $e_0 \in E$  be the edge of least weight in the given graph.

**Lemma2:** There is a MST  $T$  containing  $e_0$ .

**Proof:** Consider any MST  $T$ . Let  $e_0 \notin T$ .

Consider the fundamental cycle  $C$  defined by  $e_0$  in  $T$ .

Swap  $e_0$  with any edge  $e \in T$  present in  $C$ .



Let  $e_0 \in E$  be the edge of least weight in the given graph.

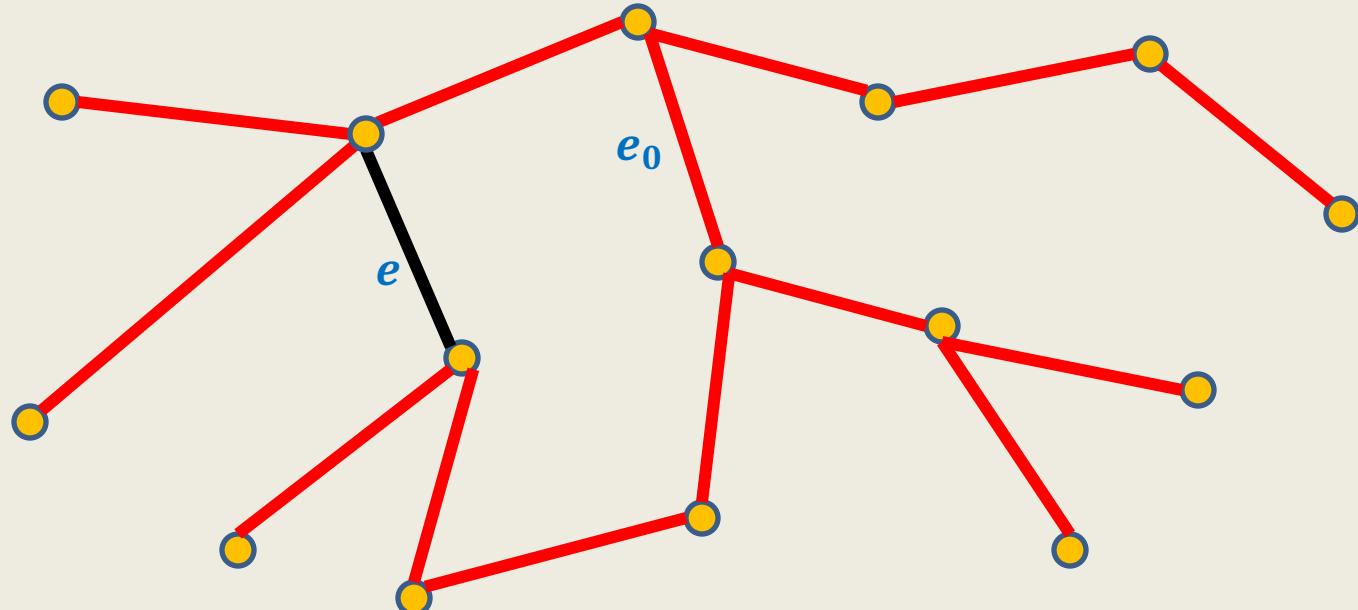
**Lemma2:** There is a MST  $T$  containing  $e_0$ .

**Proof:** Consider any MST  $T$ . Let  $e_0 \notin T$ .

Consider the fundamental cycle  $C$  defined by  $e_0$  in  $T$ .

Swap  $e_0$  with any edge  $e \in T$  present in  $C$ .

We get a spanning tree of weight  $\leq w(T)$ .



Try to translate Lemma2 to an algorithm for MST ?

with **inspiration** from the job scheduling problem ☺

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 36**

- A new algorithm design paradigm: Greedy strategy  
part III

# **Continuing Problem from last class**

**Minimum spanning tree**

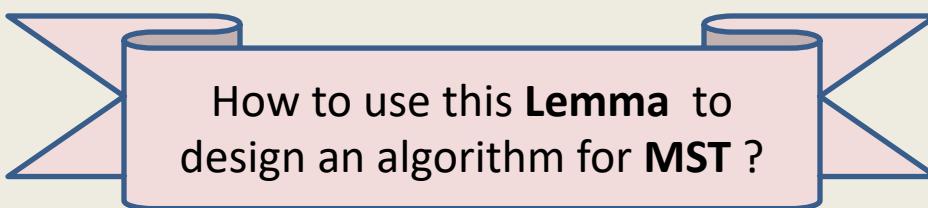
# Problem Description

**Input:** an undirected graph  $G=(V,E)$

**Aim:** compute a **spanning tree**  $(V, E')$

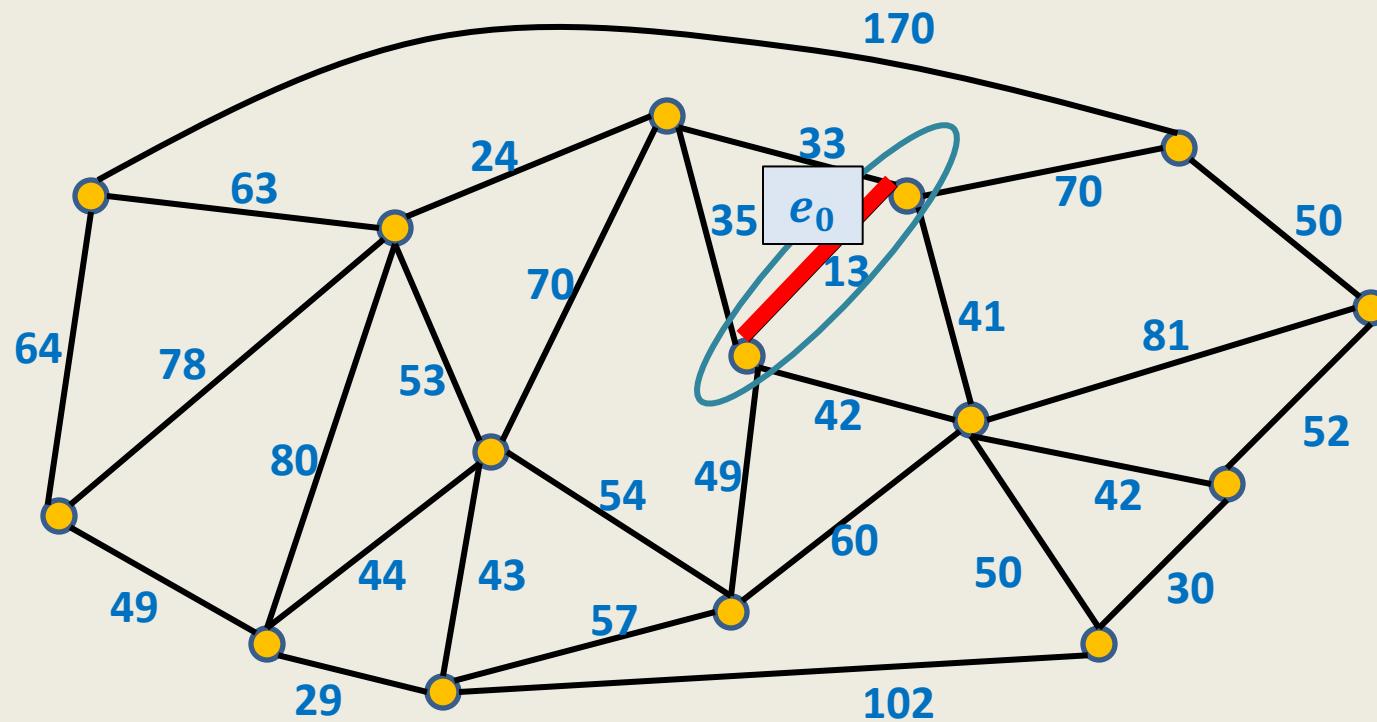
**Lemma (proved in last class):**

If  $e_0 \in E$  is the edge of **least weight** in  $G$ ,  
then there is a **MST**  $T$  containing  $e_0$ .

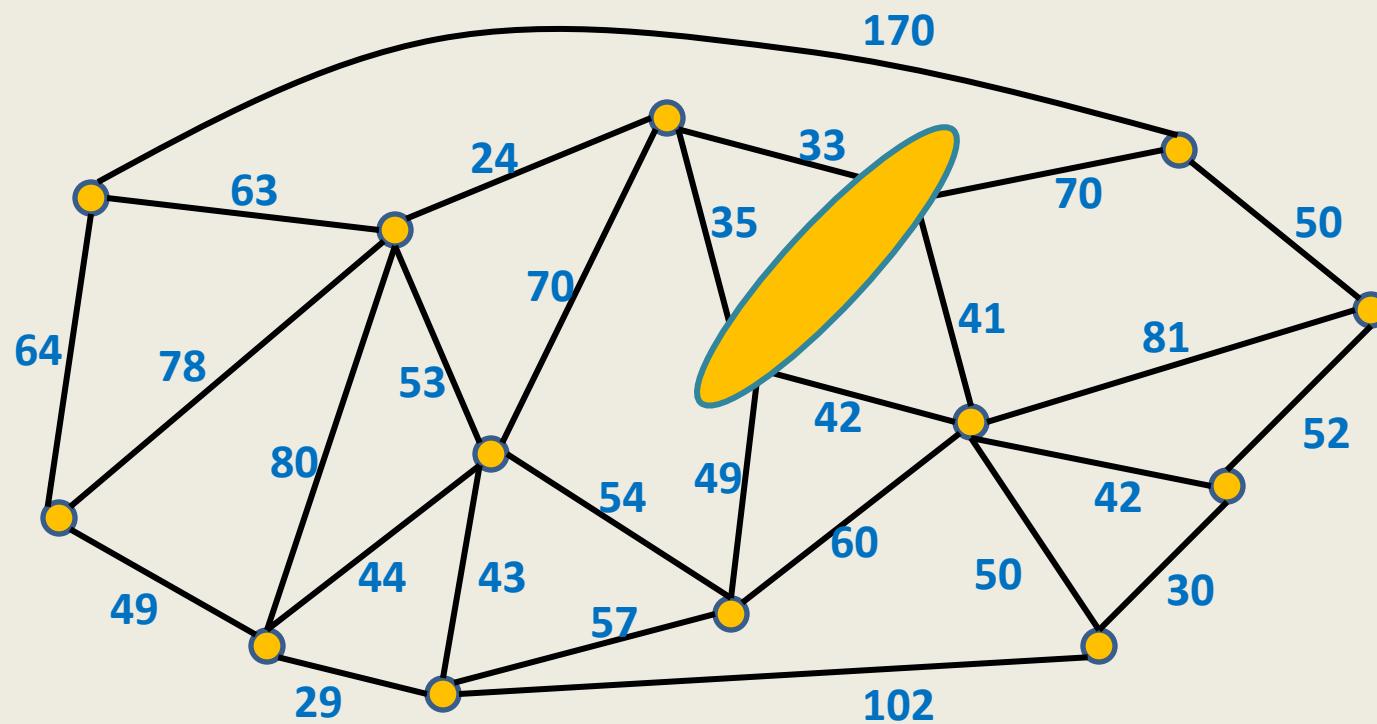


How to use this **Lemma** to  
design an algorithm for **MST** ?

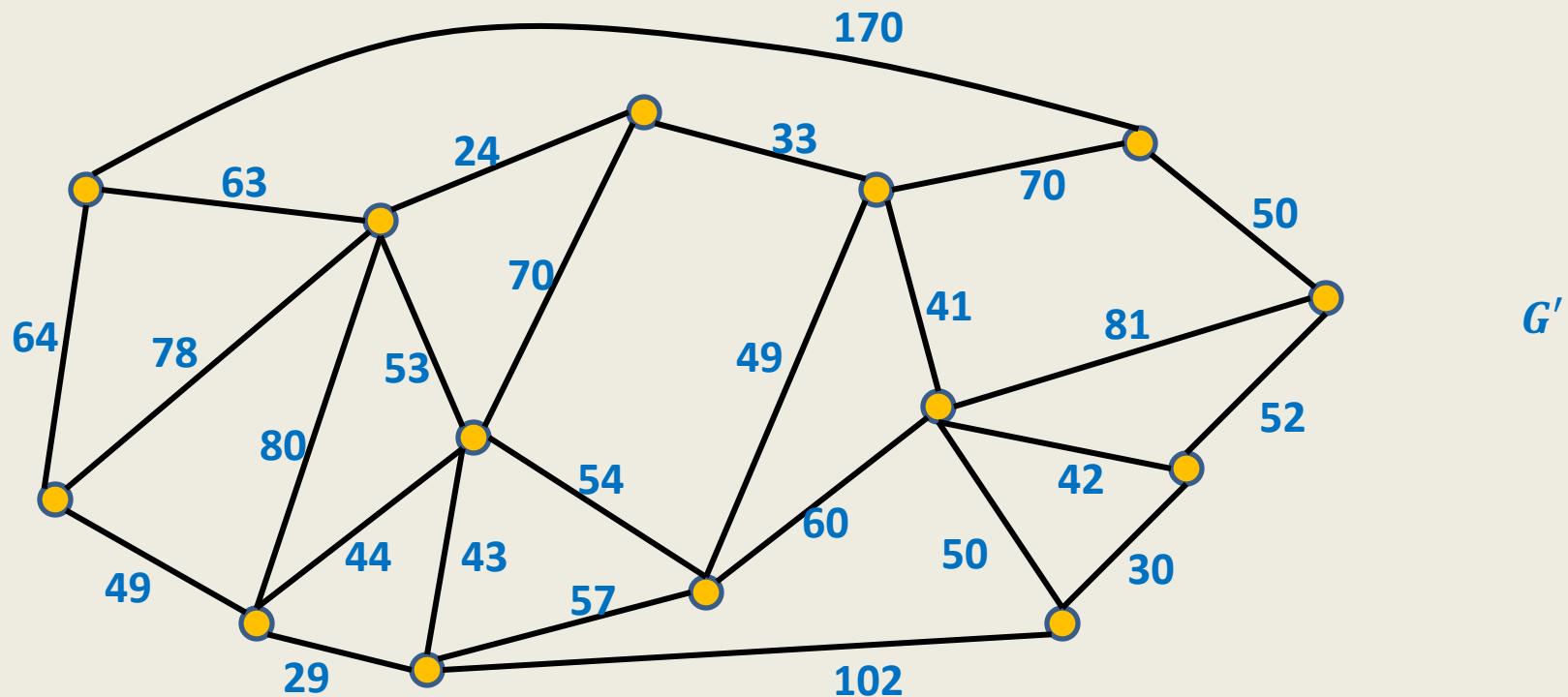
# Minimum Spanning Tree (MST)



# Minimum Spanning Tree (MST)



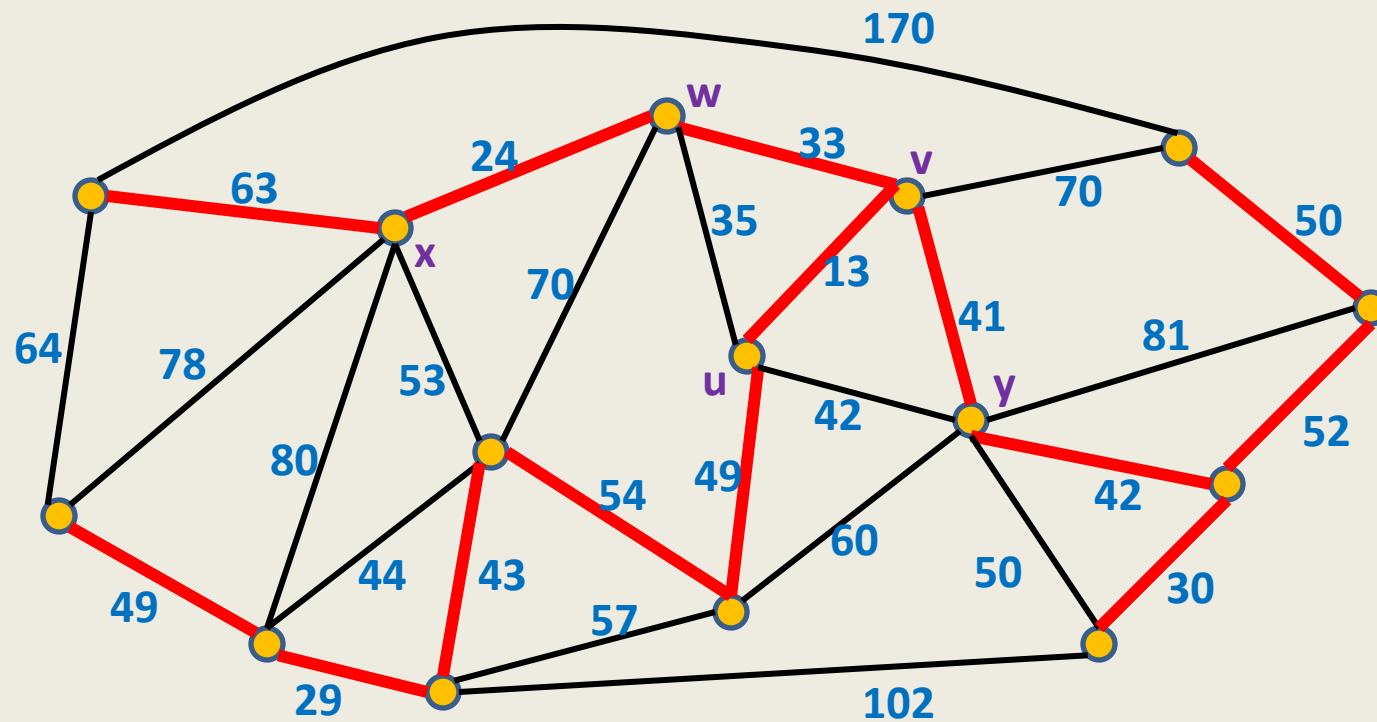
# Minimum Spanning Tree (MST)



Theorem:

$$w(\text{MST}(G)) =$$

# Minimum Spanning Tree (MST)



# A useful lesson for design of a graph algorithm

If you have a complicated algorithm for a graph problem, ...

- search for **some graph theoretic property**

to design simpler and more efficient algorithm

# Two graph theoretic properties of MST

- Cut property
- Cycle property



Every algorithm till date is based on  
one of these properties!

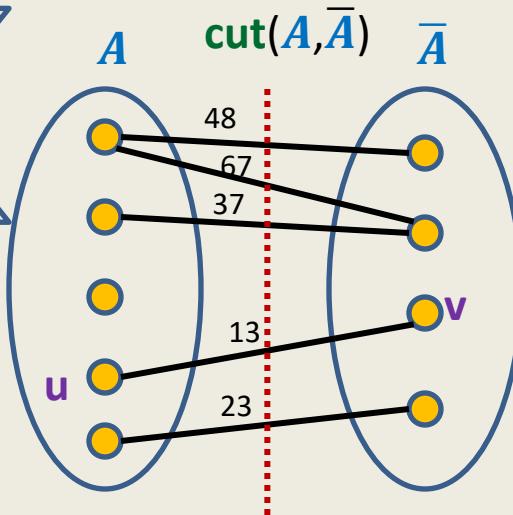
# Cut Property

# Cut Property

**Definition:** For any subset  $A \subseteq V$

$$\text{cut}(A, \bar{A}) = \{ (u, v) \in E \mid$$

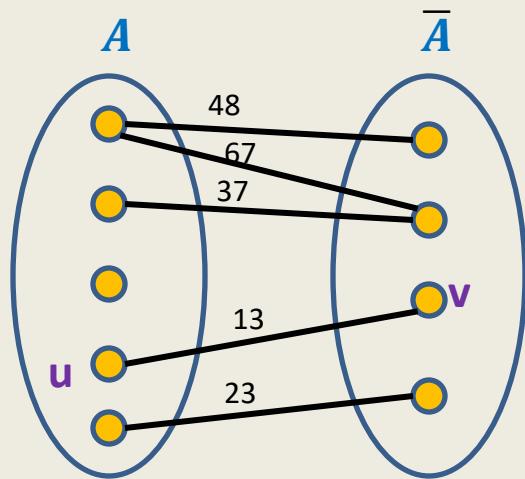
Pursuing **greedy strategy** to minimize weight of MST, what can we say about the edges of  $\text{cut}(A, \bar{A})$  ?



**Cut-property:**

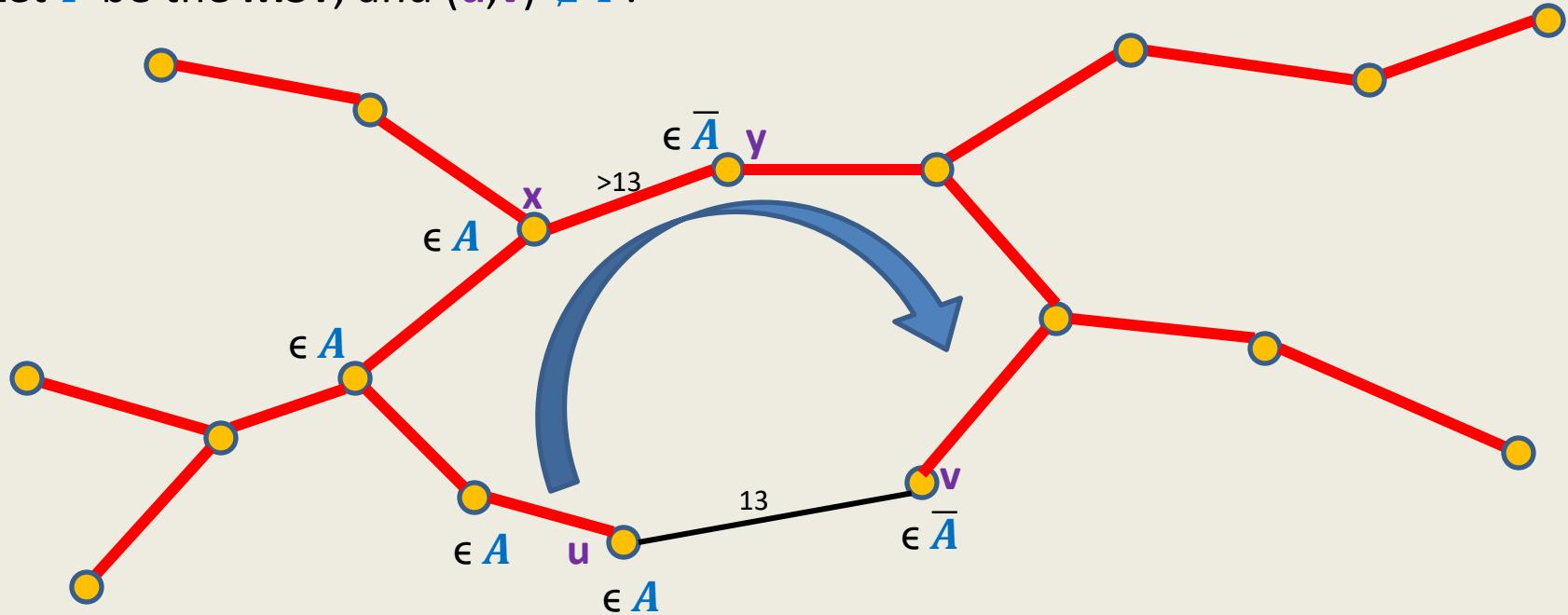
The **least weight edge** of a  $\text{cut}(A, \bar{A})$  must be in **MST**.

# Proof of cut-property



# Proof of cut-property

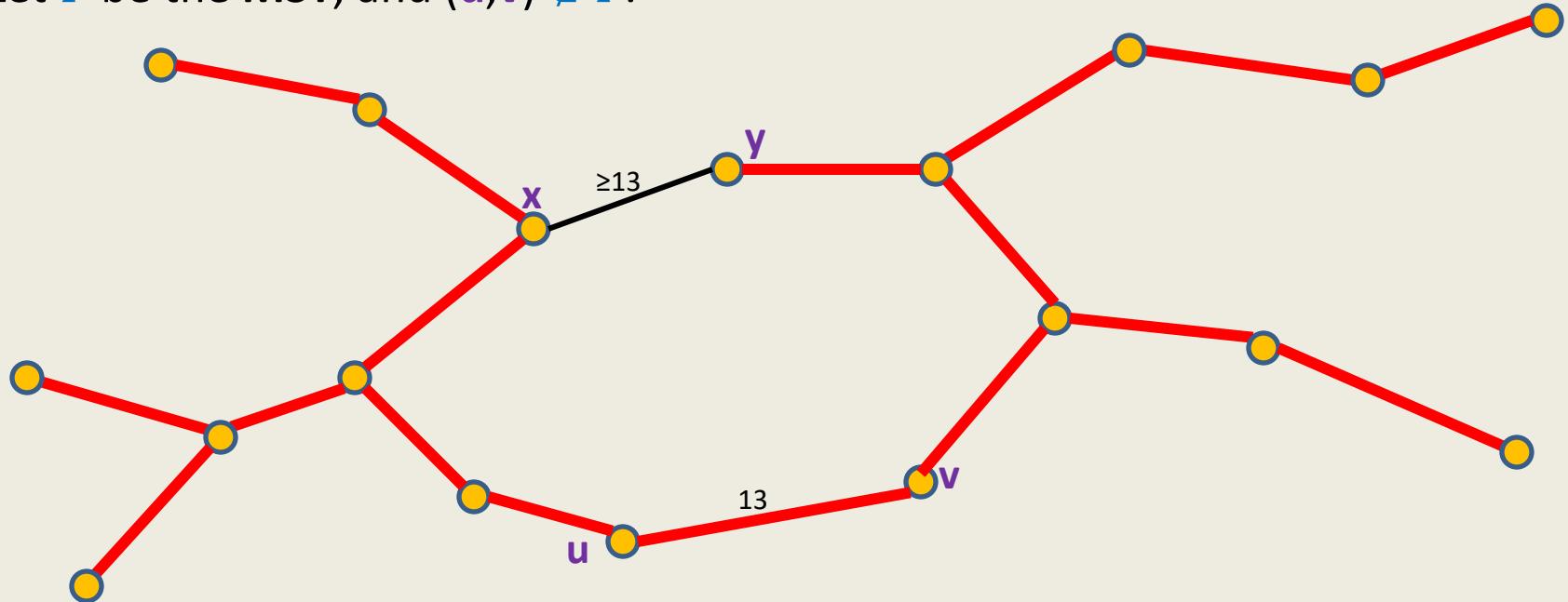
Let  $T$  be the MST, and  $(u, v) \notin T$ .



**Question:** What happens if we remove  $(x, y)$  from  $T$ , and add  $(u, v)$  to  $T$ .

# Proof of cut-property

Let  $T$  be the MST, and  $(u, v) \notin T$ .



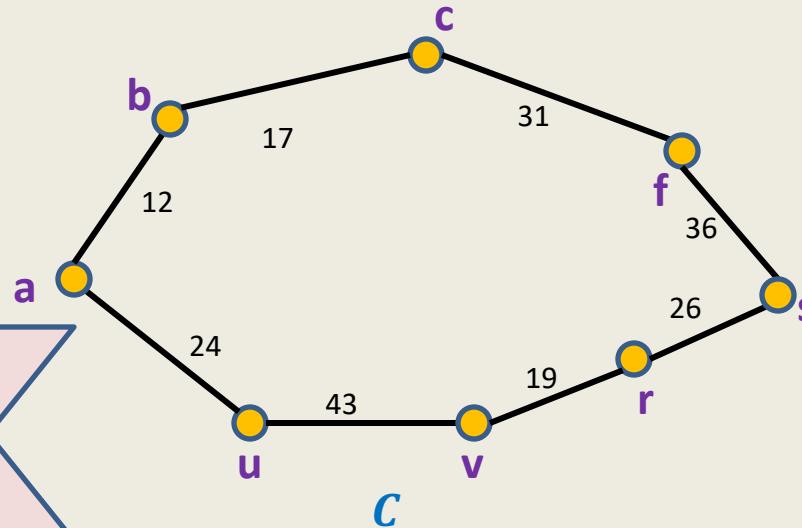
**Question:** What happens if we remove  $(x, y)$  from  $T$ , and add  $(u, v)$  to  $T$ .

We get a spanning tree  $T'$  with weight < weight( $T$ )  
**A contradiction !**

# Cycle Property

# Cycle Property

Let  $\mathcal{C}$  be any cycle in  $G$ .

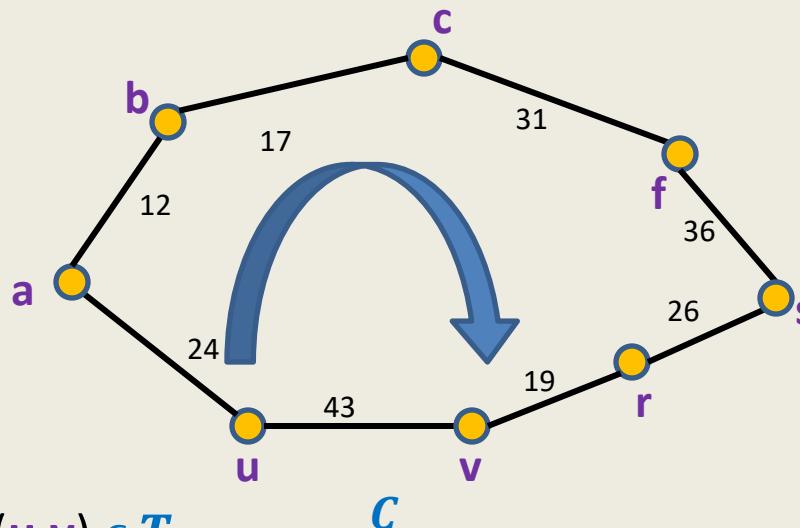


Pursuing **greedy strategy** to minimize weight of MST, what can we say about the edges of cycle  $\mathcal{C}$  ?

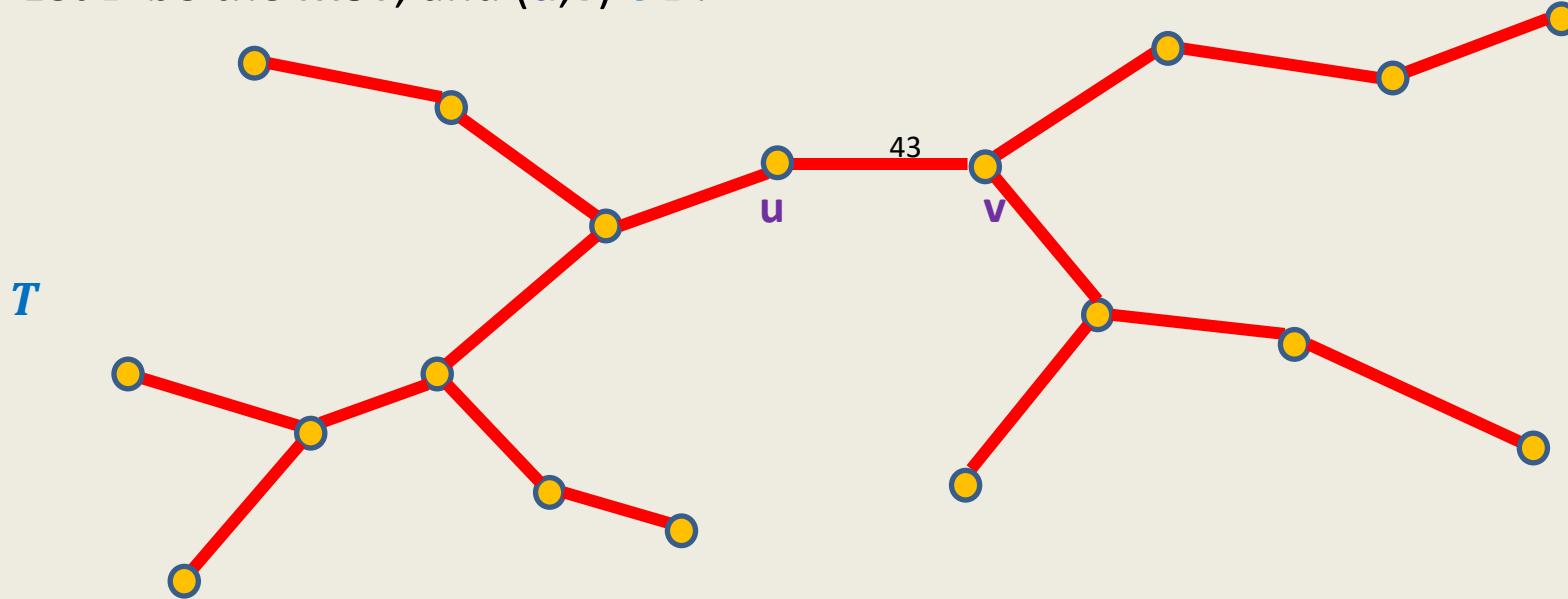
**Cycle-property:**

**Maximum weight** edge of any cycle  $\mathcal{C}$  **can not** be present in **MST**.

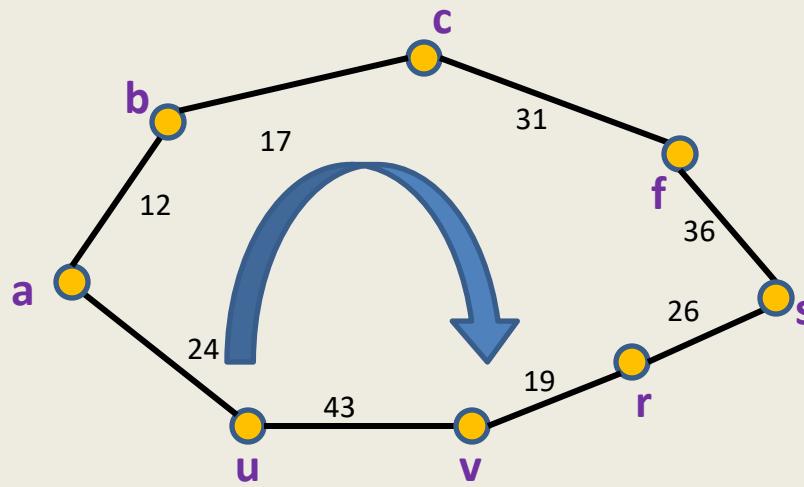
# Proof of Cycle property



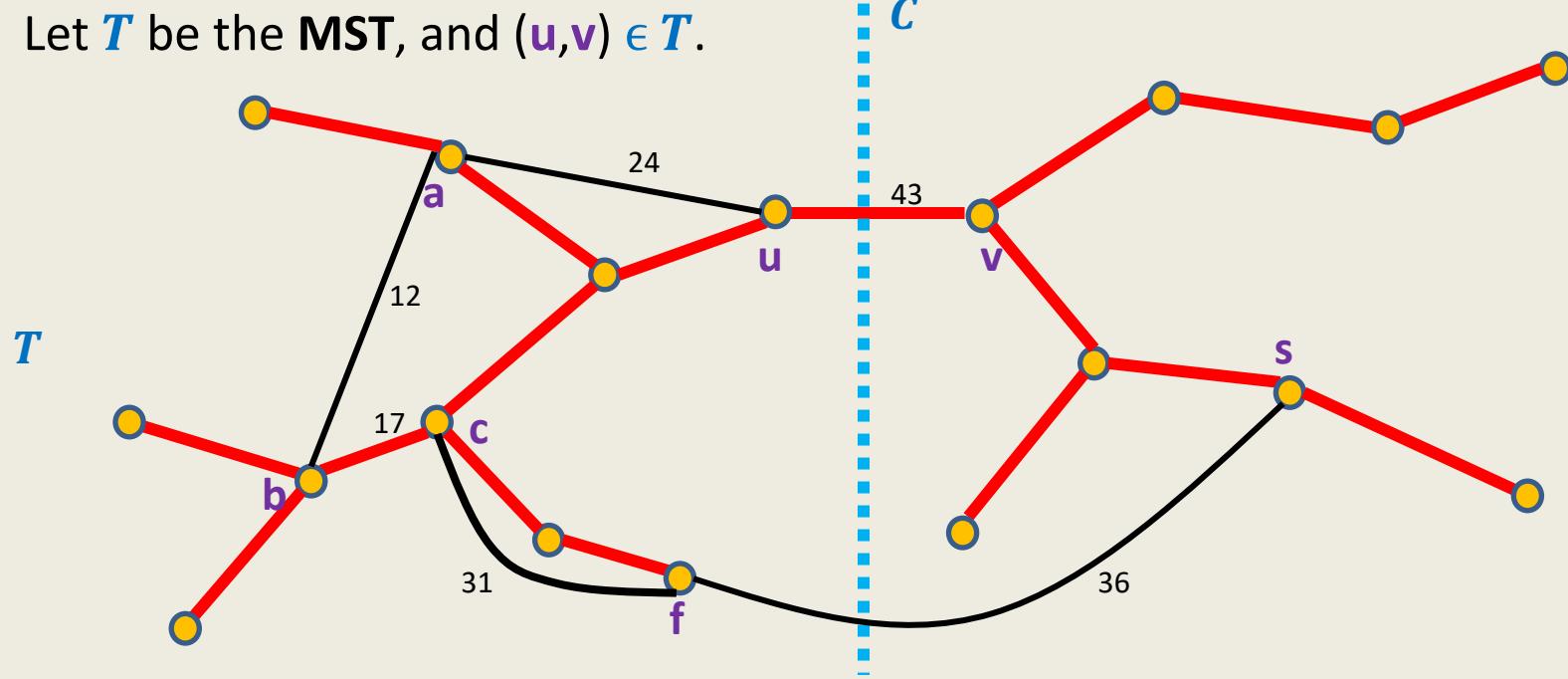
Let  $T$  be the MST, and  $(u,v) \in T$ .



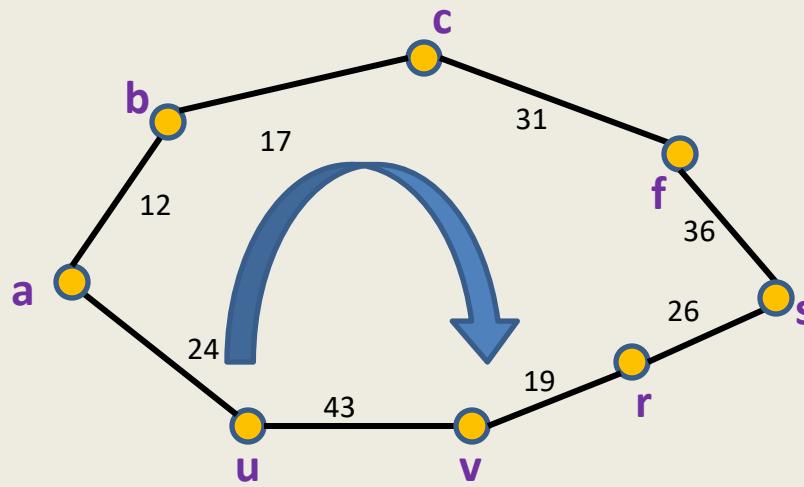
# Proof of Cycle property



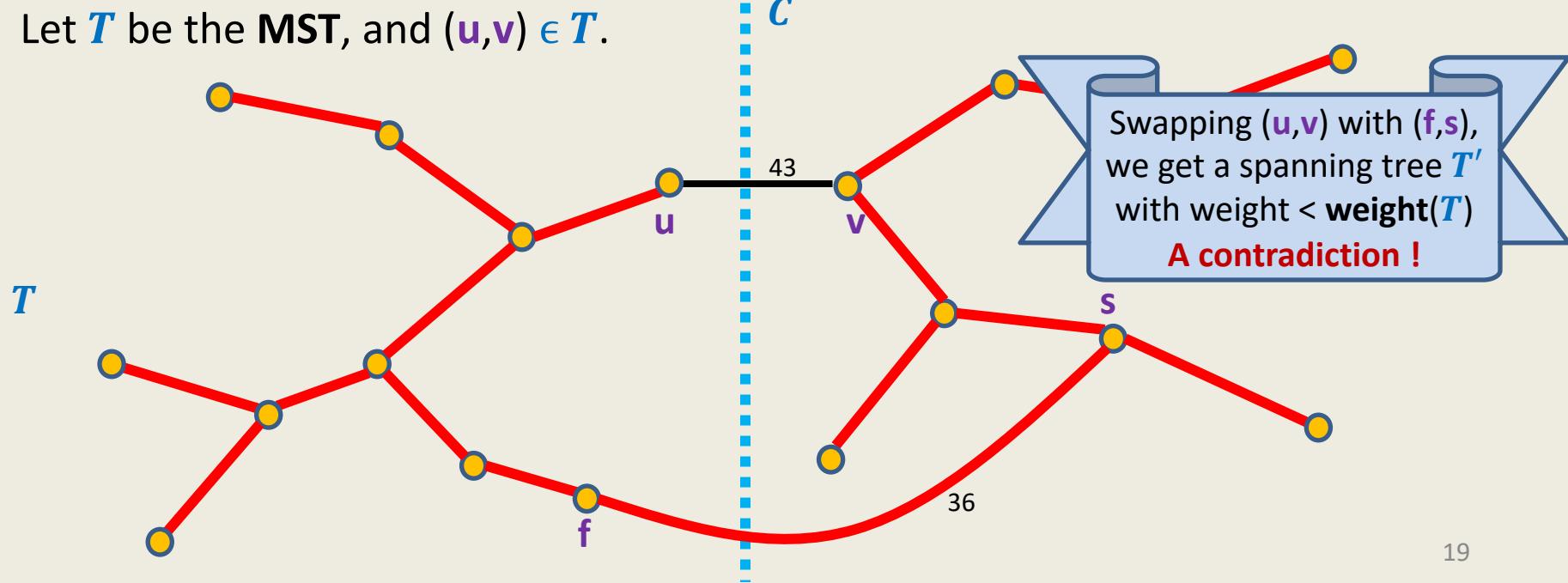
Let  $T$  be the MST, and  $(u,v) \in T$ .



# Proof of Cycle property

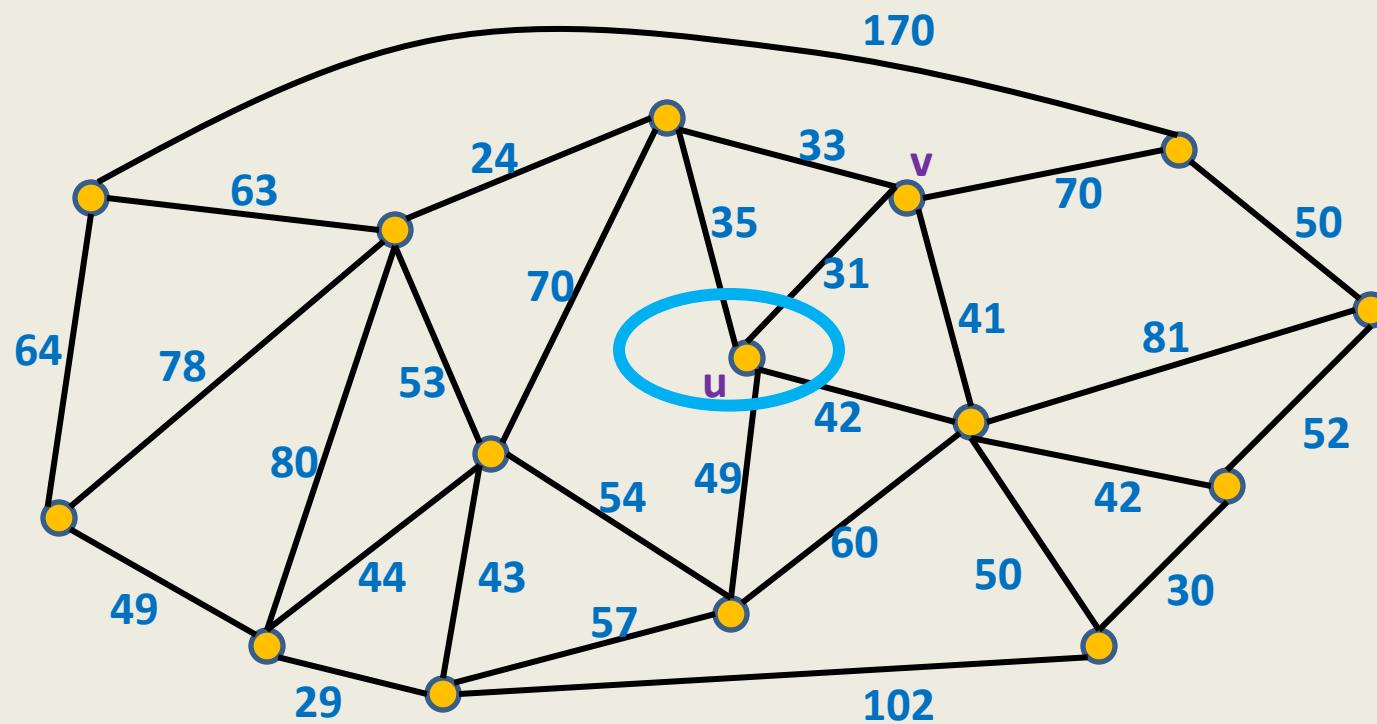


Let  $T$  be the MST, and  $(u,v) \in T$ .

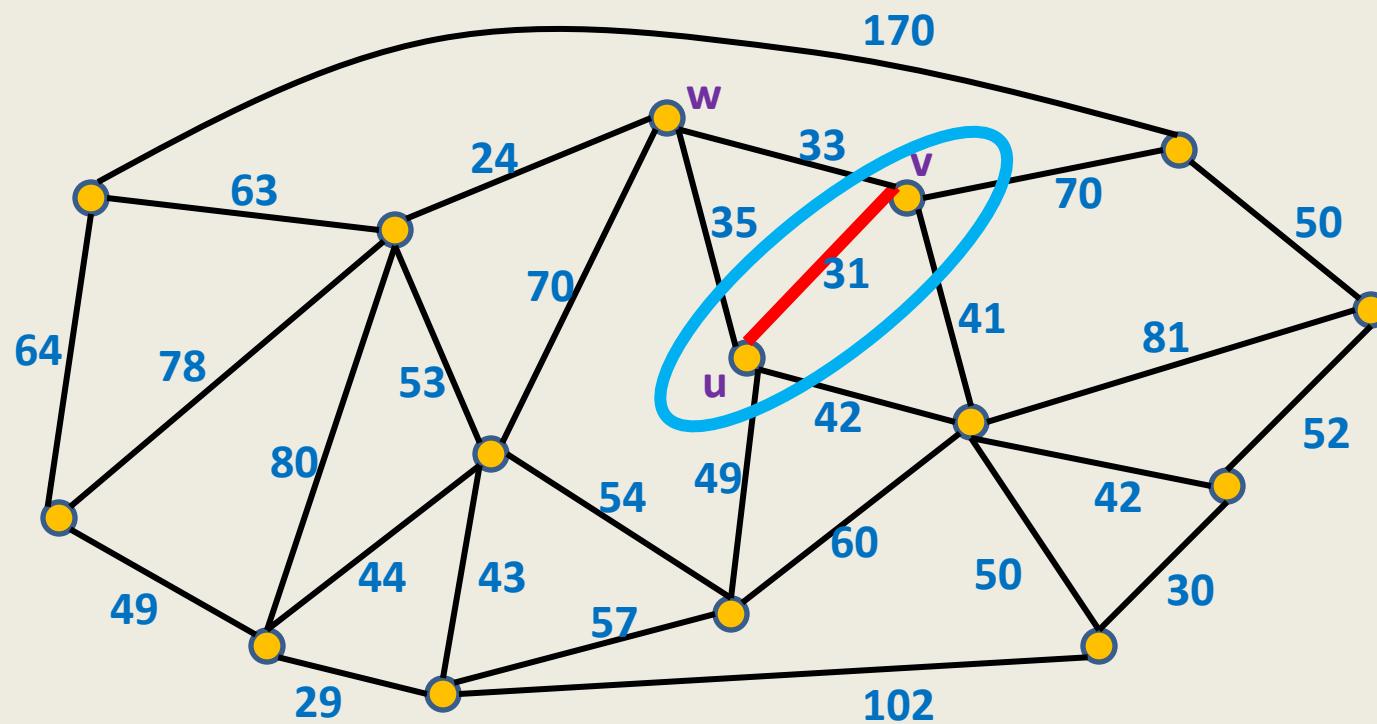


# Algorithms based on cut Property

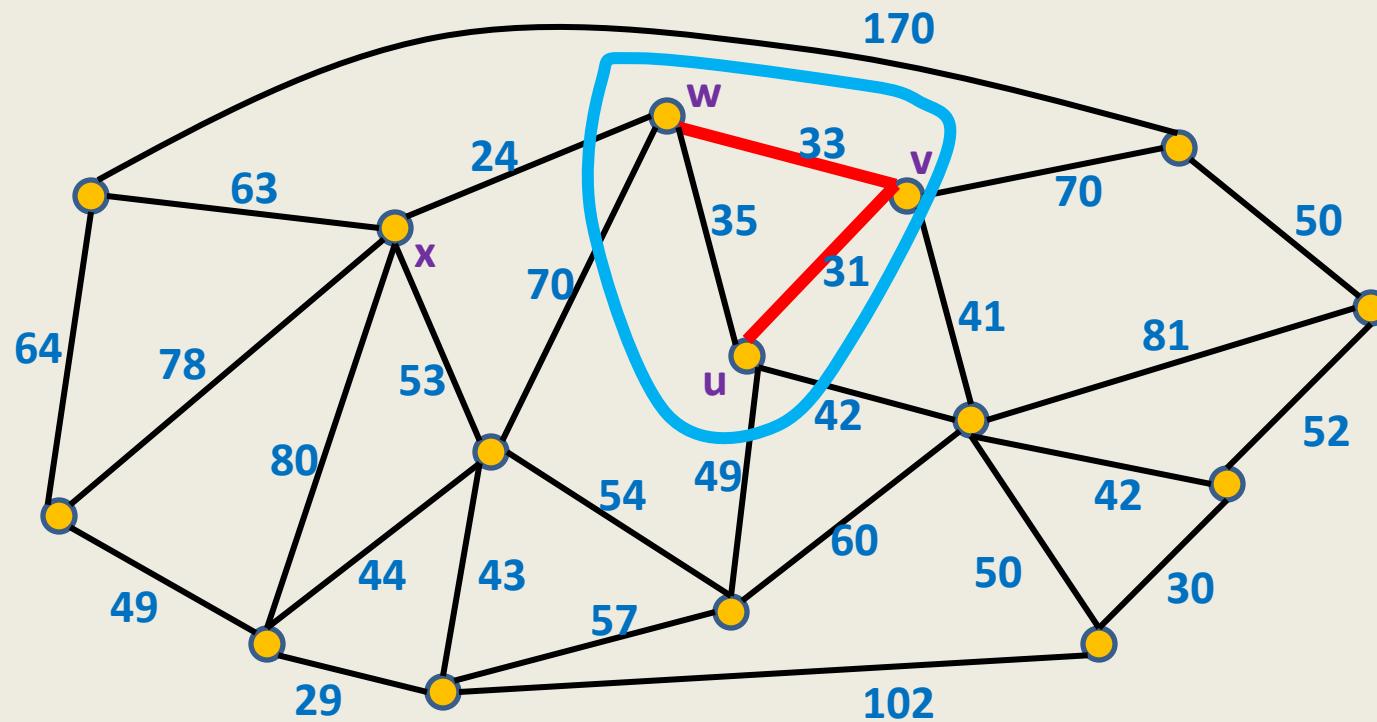
# How to use cut property to compute a MST ?



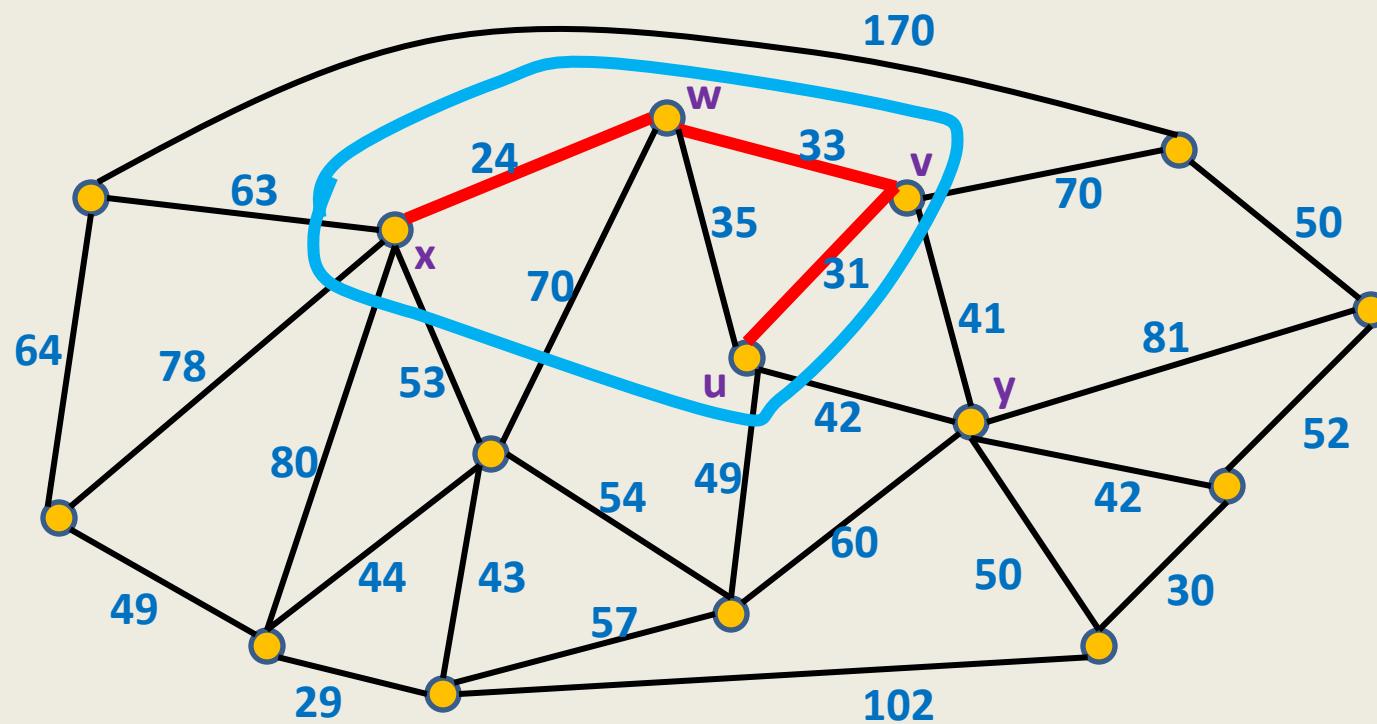
# How to use cut property to compute a MST ?



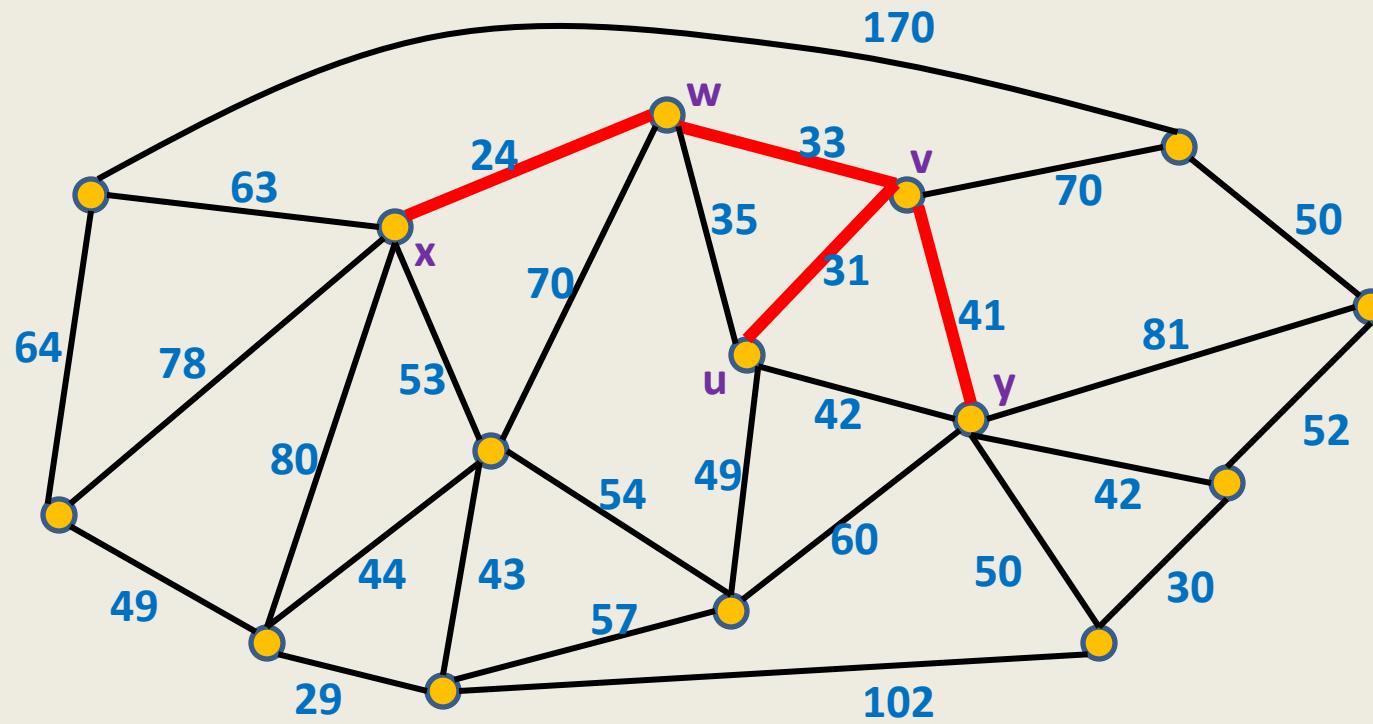
# How to use cut property to compute a MST ?



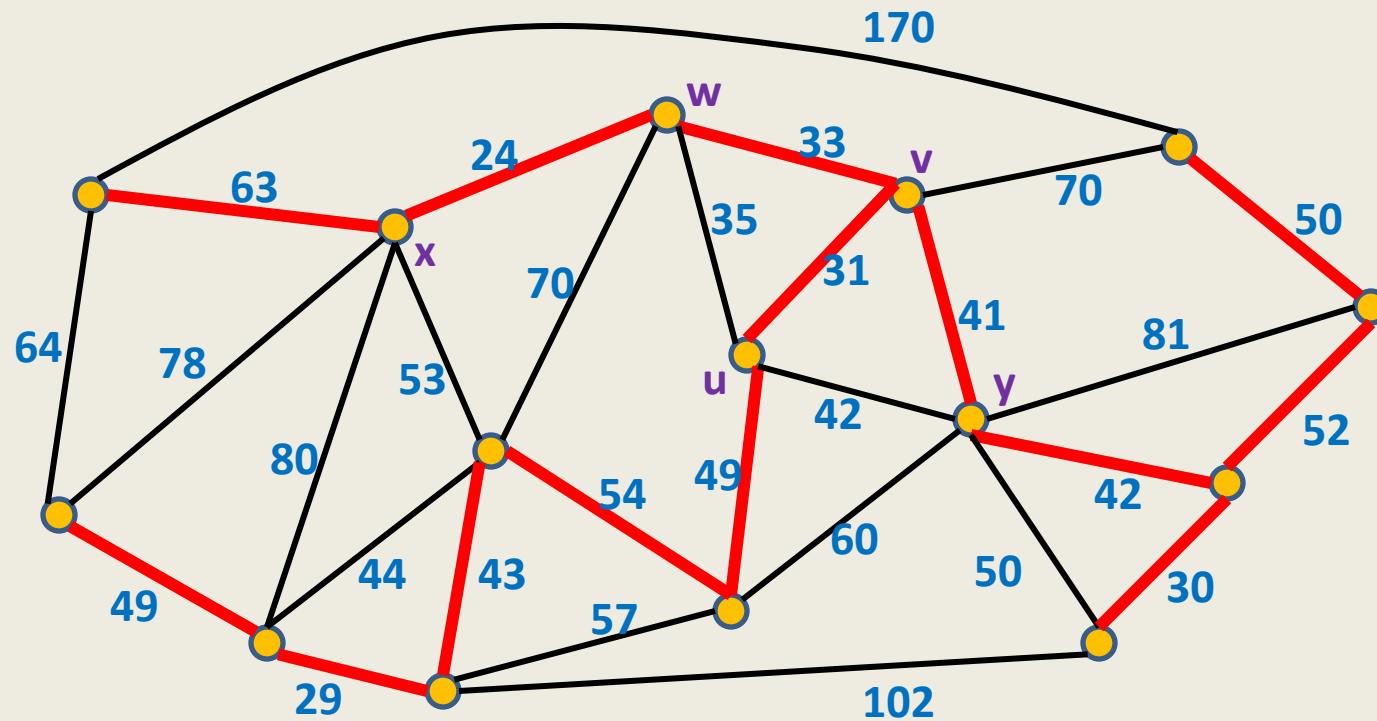
# How to use cut property to compute a MST ?



# How to use cut property to compute a MST ?



# How to use cut property to compute a MST ?



# An Algorithm based on **cut** property

**Algorithm** (Input: graph  $G = (V, E)$  with **weights** on edges)

$T \leftarrow \emptyset;$

$A \leftarrow \{u\};$

**While** ( $A <> V$ ) **do**

{ Compute the least weight edge from  $\text{cut}(A, \bar{A})$ ;

Let this edge be  $(x, y)$ , with  $x \in A, y \in \bar{A}$ ;

$T \leftarrow T \cup \{(x, y)\};$

$A \leftarrow A \cup \{y\};$

}

Return  $T$ ;

---

Number of iterations of the **While** loop :  $n - 1$

Time spent in one iteration of While loop:  $O(m)$

→ Running time of the algorithm:  $O(mn)$

# Algorithm based on cycle Property

# An Algorithm based on cycle property

## Description

**Algorithm** (Input: graph  $G = (V, E)$  with weights on edges)

While ( $E$  has any cycle) do

{ Compute any cycle  $C$ ;

Let  $(u, v)$  be the maximum weight edge of the cycle  $C$ ;

Remove  $(u, v)$  from  $E$ ;

}

Return  $E$ ;

---

Number of iterations of the While loop :  $m - n + 1$

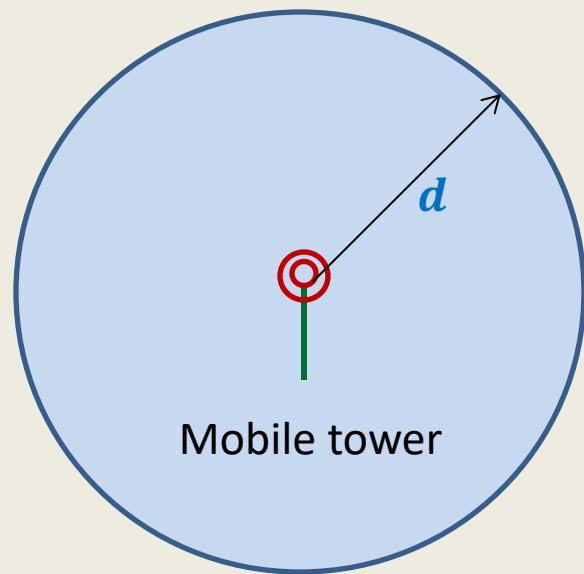
Time spent in one iteration of While loop:  $O(n)$

→ Running time of the algorithm:  $O(mn)$

# **Problem 3**

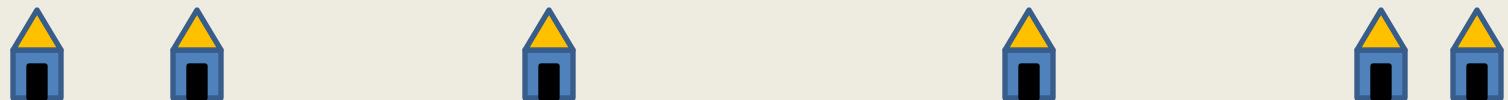
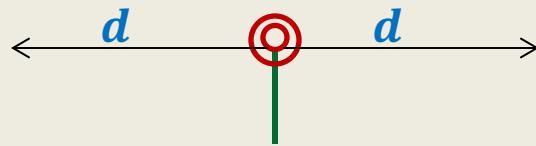
**Mobile towers on a road**

# Mobile towers on a road



A mobile tower can cover any cell phone within radius  $d$ .

# Mobile towers on a road



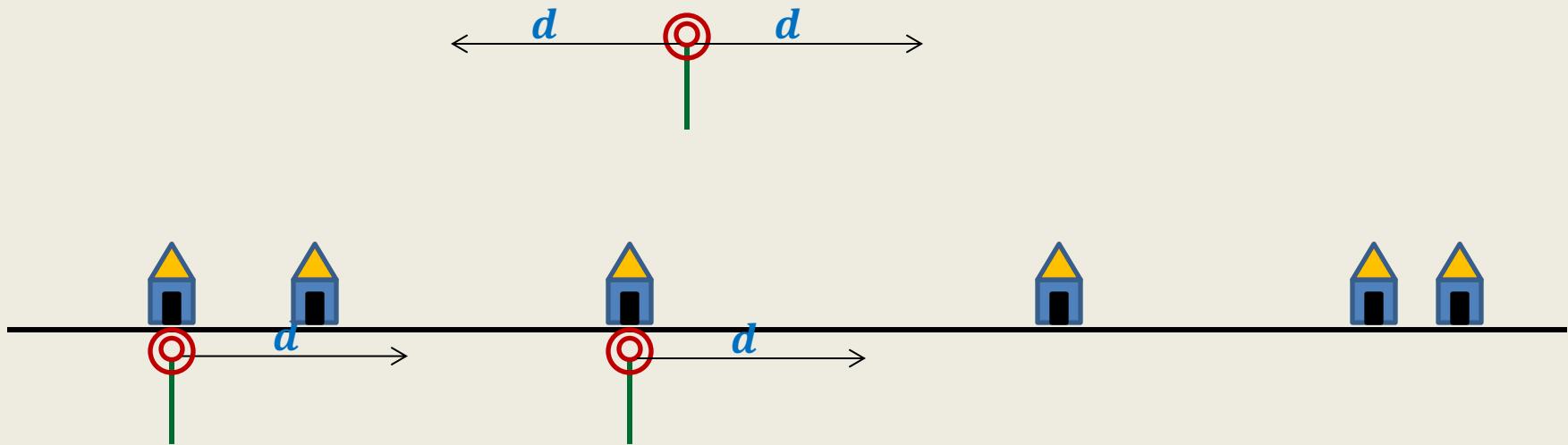
## Problem statement:

There are  $n$  houses located along a road.

We want to place mobile towers such that

- Each house is covered by at least one mobile tower.
- The number of mobile towers used is **least** possible.

# Mobile towers on a road



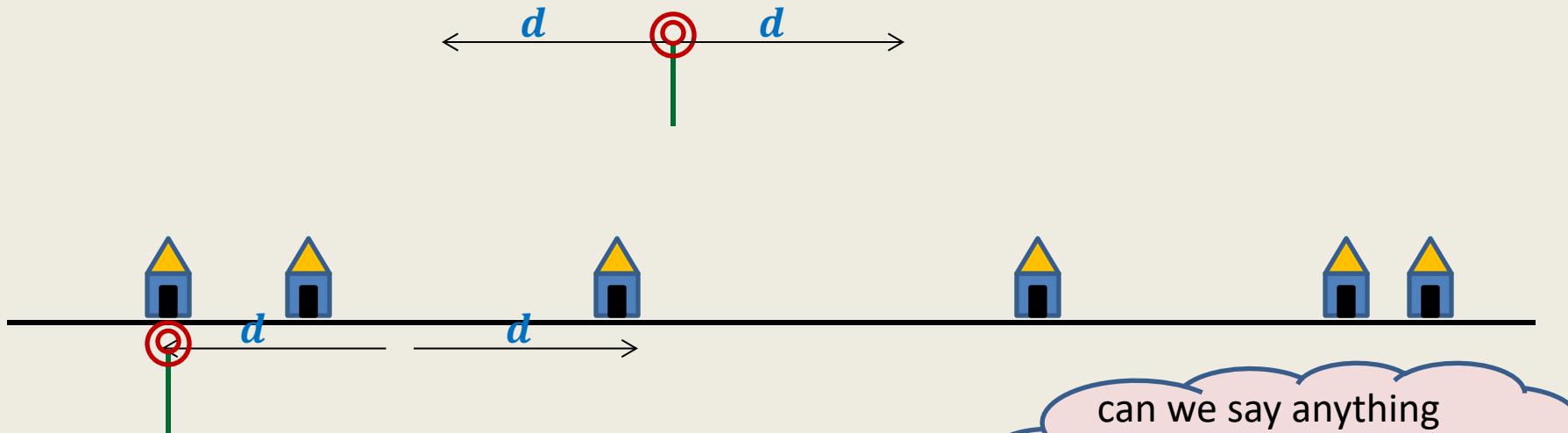
## Strategy 1:

Place tower at first house,

Remove all houses covered by this tower.

Proceed to the next uncovered house ...

# Mobile towers on a road



## Strategy 2:

Place tower at distance  $d$  to the right of the first house;

Remove all houses covered by this tower;

Proceed to the next uncovered house along the road...

**Lemma:** There is an optimal solution for the problem in which the leftmost tower is placed at distance  $d$  to the right of the first house

# Homework ...

Ponder over the following questions before coming for the next class

- Use **cycle property** and/or **cut property** to design a **new algorithm for MST**
- Use some data structure to improve the running time of the algorithms discussed in this class to  **$O(m \log n)$**

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 37**

- A new algorithm design paradigm: Greedy strategy  
part **IV**

# Problems solved till now

1. Job Scheduling Problem

2. Mobile Tower Problem

3. MST



Did you notice  
anything common in  
their solutions ?

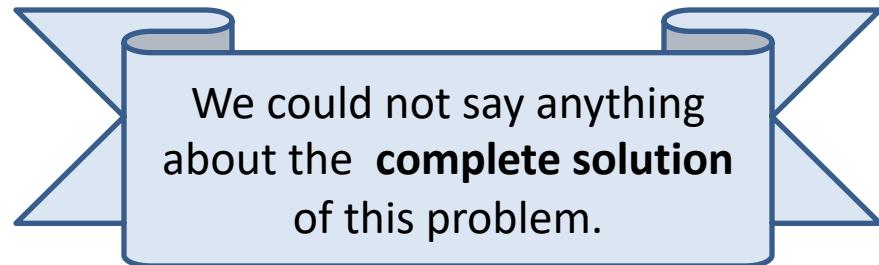
Ponder over this question before moving ahead

# Problem 1

## Job scheduling Problem

### INPUT:

- A set  $J$  of  $n$  jobs  $\{j_1, j_2, \dots, j_n\}$
- job  $j_i$  is specified by two real numbers
  - $s(i)$ : start time of job  $j_i$
  - $f(i)$ : finish time of job  $j_i$
- A single server



### Constraints:

- Server can execute at most one job at any moment of time and a job.
- **Job  $j_i$** , if scheduled, has to be scheduled during  $[s(i), f(i)]$  only.

### Aim:

To select the **largest** subset of non-overlapping jobs which can be executed by the server.

# Problem 1

## Job scheduling Problem

### INPUT:

- A set  $J$  of  $n$  jobs  $\{j_1, j_2, \dots, j_n\}$
- job  $j_i$  is specified by two real numbers
  - $s(i)$ : start time of job  $j_i$
  - $f(i)$ : finish time of job  $j_i$
- A single server

### Constraints:

- Server can execute at most one job at any moment of time and a job.
- **Job  $j_i$** , if scheduled, has to be scheduled during  $[s(i), f(i)]$  only.

### Aim:

To select the **largest** subset of non-overlapping jobs which can be executed by the server.

# All that we could do was to make a local observation

Let  $x \in J$  be the job with earliest finish time.

**Lemma1 :** There exists an optimal solution for  $J$  in which  $x$  is present.

Let  $J' = J \setminus \text{Overlap}(x)$

**Lemma 1** gives very small information  
about the optimal solution ☹  
How to use it to compute this solution ?

$J$ (original instance)

$\text{Opt}(J)$

$$\text{Opt}(J) = \text{Opt}(J') + 1 \quad \text{-- (i)}$$

Greedy  
step

$J'$  (smaller instance)

**Lemma1**

$\text{Opt}(J')$

Equation (i) hints at recursive solution of the problem ☺

**Theorem:**  $\text{Opt}(J) = \text{Opt}(J') + 1$ .

- Proof has two parts

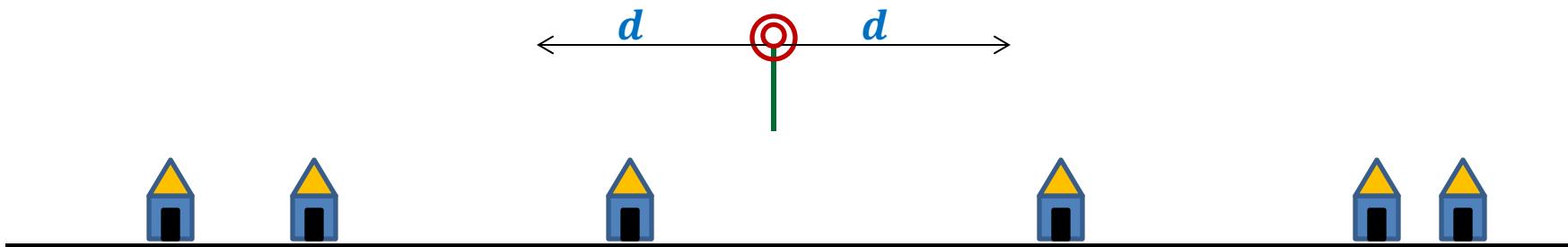
$$\text{Opt}(J) \geq \text{Opt}(J') + 1$$

$$\text{Opt}(J') \geq \text{Opt}(J) - 1$$

- Proof for each part is a proof **by construction**

# Problem 2

## Mobile Tower Problem

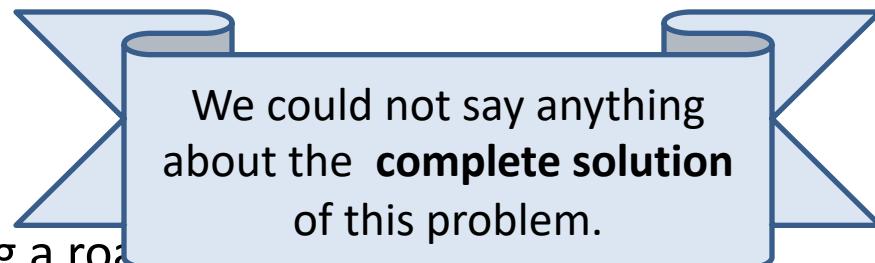


### Problem statement:

There is a set  $H$  of  $n$  houses located along a road.

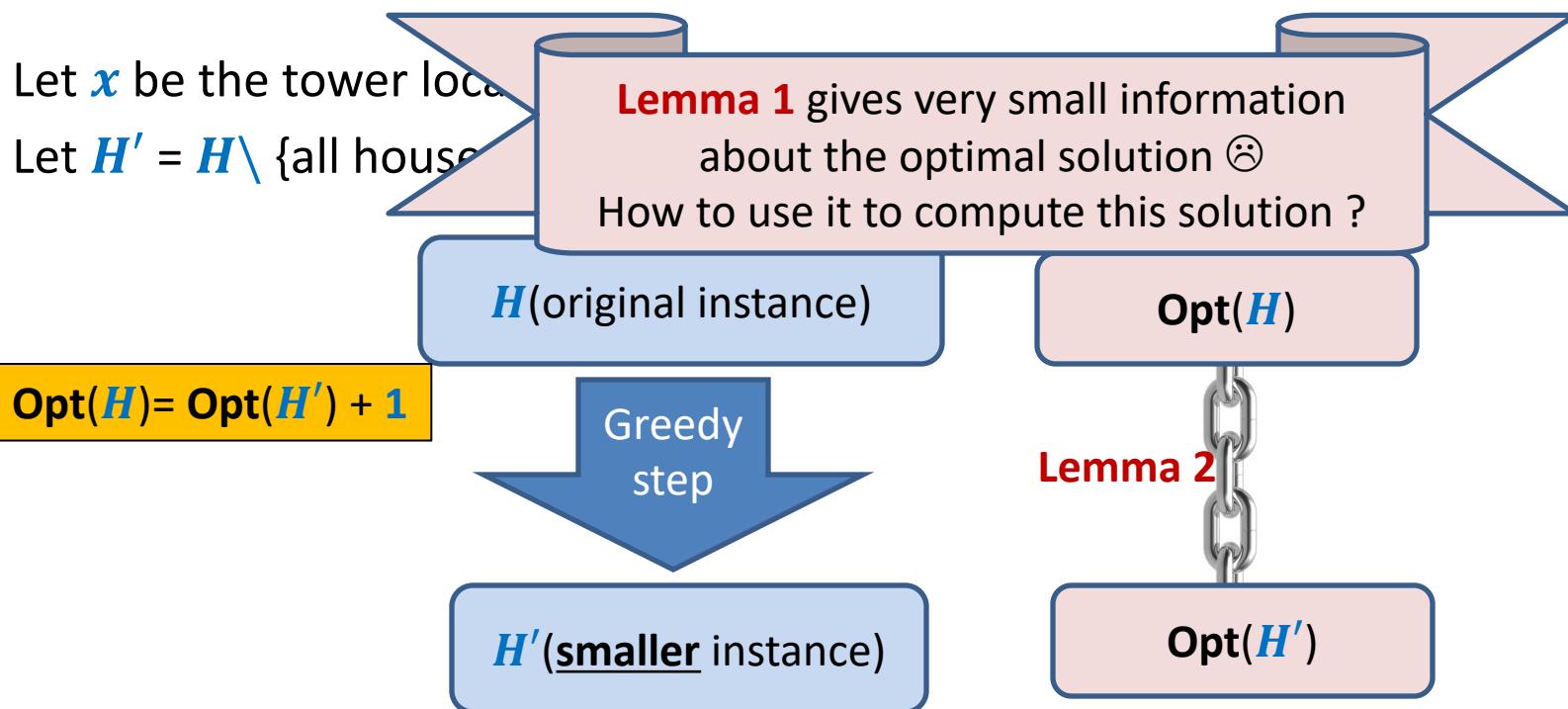
We want to place mobile towers such that

- Each house is covered by at least one mobile tower.
- The number of mobile towers used is **least** possible.



# All that we could do was to make a local observation

**Lemma 2:** There is an optimal solution for the problem in which the leftmost tower is placed at distance  $d$  to the right of the first house.



Equation (i) hints at recursive solution of the problem 😊

# What is a **greedy strategy** ?

A strategy that is

- Based on some **local** approach
- With the **objective to optimize** some function.

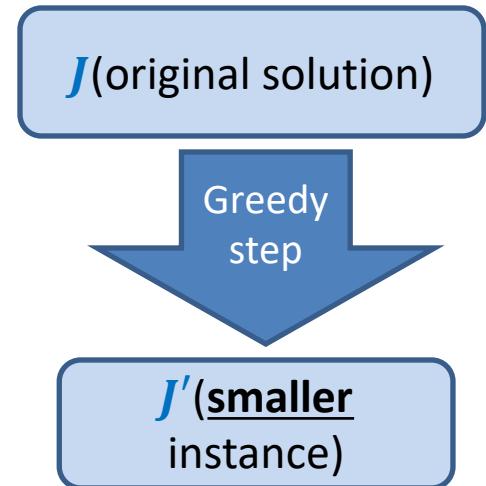
## Note:

Recall that the divide and conquer strategy takes a **global approach**.

# Design of a greedy algorithm

Let  $\mathbf{A}$  be an instance of an optimization problem.

1. Make a **local observation** about the solution.
2. Use this observation to express optimal solution of  $\mathbf{A}$  in terms of
  - Optimal solution of a smaller instance  $\mathbf{A}'$
  - Local step
3. This gives a recursive solution.
4. Transform it into iterative one.



# MST

**Input:** an undirected graph  $G=(V,E)$  with  $w: E \rightarrow \mathbb{R}$ ,

**Aim:** compute a **spanning tree**  $(V, E')$ ,  $E' \subseteq E$  such that  $\sum_{e \in E'} w(e)$  is **minimum**.

**Lemma 2**

If you have understood a generic way to design a greedy algorithm, then try to solve the MST problem.

If  $e_0 \in E$  is the edge of least weight in  $G$ , then there is a MST  $T$  containing  $e_0$ .

How to use this **Lemma** to design an algorithm for **MST** ?

# Problem 4

## Overlapping Intervals

The aim of this problem is to make you realize that it is sometime very nontrivial to design a greedy algorithm. In particular, it is quite challenging to design the smaller instance.

# **Problem 4**

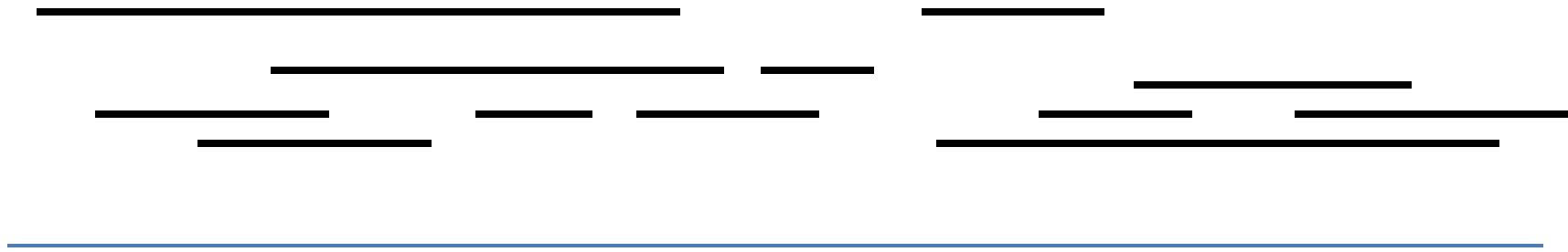
## **Overlapping Intervals**

# Overlapping Intervals

## Problem statement:

Given a set **A** of  $n$  intervals, compute smallest set **B** of intervals so that for every interval **I** in  $A \setminus B$ , there is some interval in **B** which overlaps/intersects with **I**.

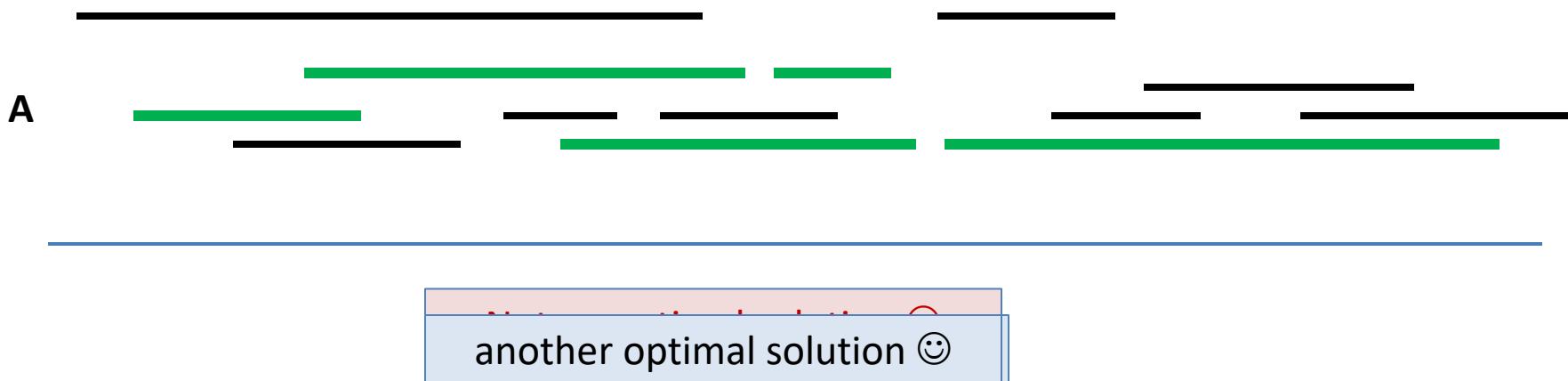
**A**



# Overlapping Intervals

## Problem statement:

Given a set **A** of  $n$  intervals, compute smallest set **B** of intervals so that for every interval **I** in  $A \setminus B$ , there is some interval in **B** which overlaps/intersects with **I**.



# Overlapping Intervals

## Strategy 1

Interval with maximum length should be there in optimal solution



## Intuition:

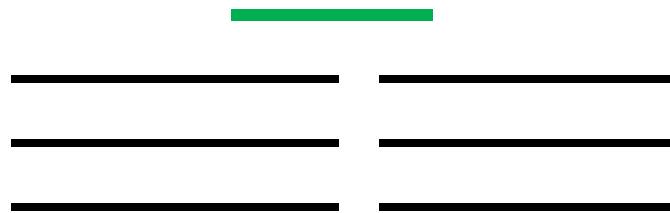
Selecting such an interval will **cover maximum** no. of other intervals

There is a counter example 😞

# Overlapping Intervals

## Strategy 1

Interval with maximum length should be there in optimal solution



## Intuition:

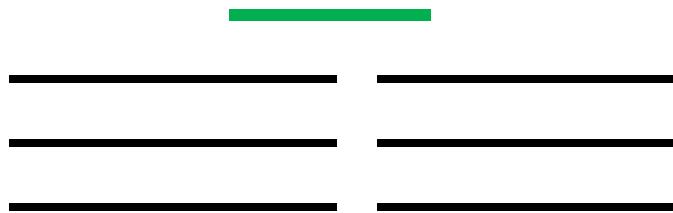
Selecting such an interval will **cover maximum** no. of other intervals

There is a counter example 😔

# Overlapping Intervals

## Strategy 2

Interval that overlaps maximum no. of intervals should be there in optimal solution



### Intuition:

Selecting such an interval will **cover maximum** no. of other intervals

There is a counter example 😞

# Overlapping Intervals

## Strategy 2

Interval that overlaps maximum no. of intervals should be there in optimal solution



# Overlapping Intervals

## Strategy 2

Interval that overlaps maximum no. of intervals should be there in optimal solution



Not an optimal solution 😞

# Overlapping Intervals

## Strategy 2

Interval that overlaps maximum no. of intervals should be there in optimal solution



An optimal solution has size 2.

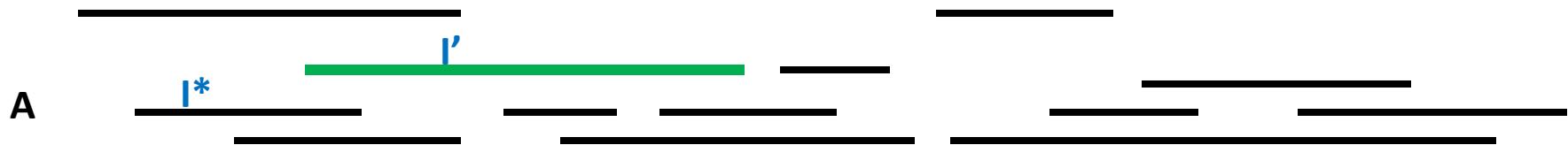
Think for a while :

After failure of two strategies, how to proceed to design the algorithm.

# Overlapping Intervals

Let  $I^*$  be the interval with earliest finish time.

Let  $I'$  be the interval with **maximum** finish time overlapping  $I^*$ .



---

**Lemma1:** There is an optimal solution for set A that contains  $I'$ .

**Proof:(sketch) :**

If  $I^*$  is overlapped by any other interval in the optimal solution, say  $I^\wedge$ ,

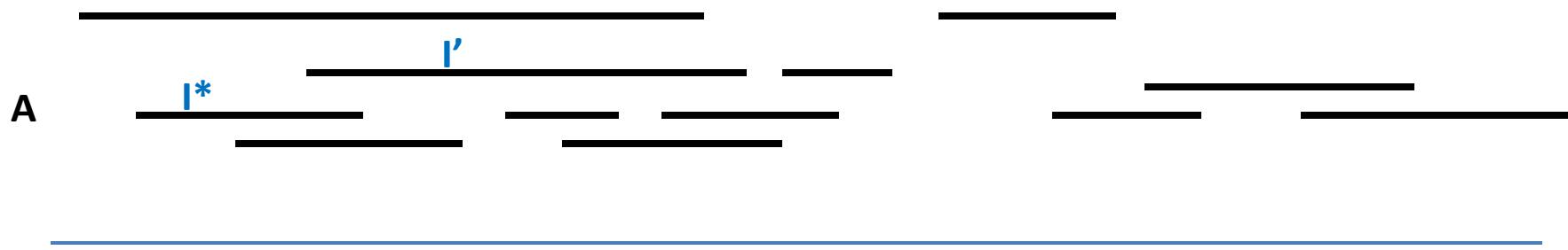
$I'$  will surely overlap all intervals that are overlapped by  $I^\wedge$ .

→ Swapping  $I^\wedge$  by  $I'$  will still give an optimal solution.

Exploit the fact that  
 $I^*$  has earliest finish  
time for this claim.

# Overlapping Intervals

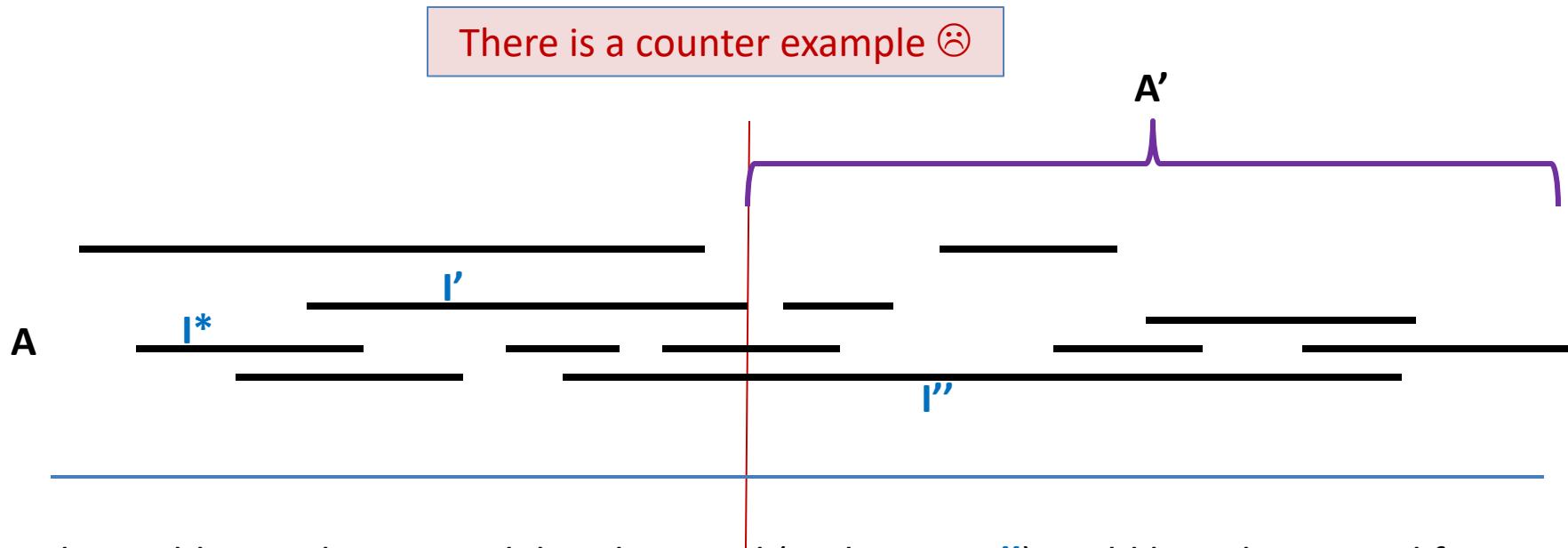
**Question:** How to obtain smaller instance  $A'$  using **Lemma 1** ?



# Overlapping Intervals

**Question:** How to obtain smaller instance  $A'$  using **Lemma 1** ?

**Naive approach :** remove from  $A$  all intervals which overlap with  $I'$ . This is  $A'$ .



The problem is that some deleted interval (in this case  $I''$ ) could have been used for intersecting many intervals if it were not deleted. But deleting it from the instance disallows it to be selected in the solution.

# Homework

- How will you form the smaller instance ?
- Design an algorithm for the problem.
- Give a neat, concise, and formal proof of correctness of the algorithm.

# Data Structures and Algorithms

(ESO207)

## Lecture 38

- An interesting problem:  
shortest path from a **source** to **destination**

# **SHORTEST PATHS IN A GRAPH**

**A fundamental problem**

# Notations and Terminologies

A directed graph  $G = (V, E)$

- $\omega: E \rightarrow R^+$
- Represented as **Adjacency lists** or **Adjacency matrix**
- $n = |V|$  ,  $m = |E|$

**Question:** what is a path in  $G$ ?

Answer: A sequence  $v_1, v_2, \dots, v_k$  such that  $(v_i, v_{i+1}) \in E$  for all  $1 \leq i < k$ .



Length of a path  $P = \sum_{e \in P} \omega(e)$

# Notations and Terminologies

## Definition:

The path from  $u$  to  $v$  of minimum length is called the **shortest path** from  $u$  to  $v$

**Definition:** Distance from  $u$  to  $v$  is the length of the shortest path from  $u$  to  $v$ .

## Notations:

$\delta(u, v)$  : distance from  $u$  to  $v$ .

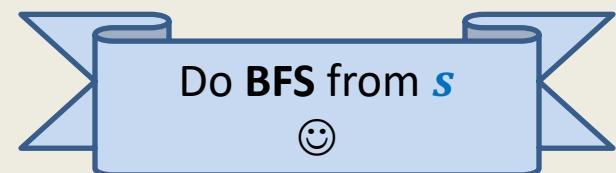
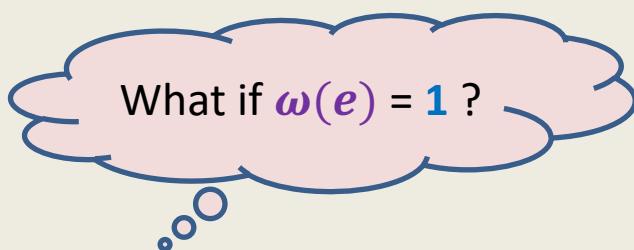
$P(u, v)$  : The shortest path from  $u$  to  $v$ .

# Problem Definition

**Input:** A directed graph  $G = (V, E)$  with  $\omega: E \rightarrow R^+$  and a source vertex  $s \in V$

## Aim:

- Compute  $\delta(s, v)$  for all  $v \in V \setminus \{s\}$
- Compute  $P(s, v)$  for all  $v \in V \setminus \{s\}$



# Problem Definition

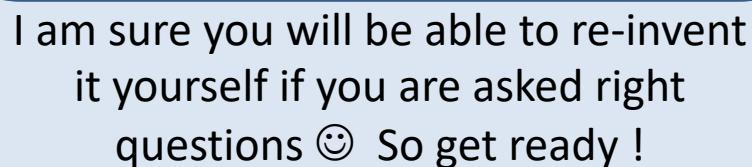
**Input:** A directed graph  $G = (V, E)$  with  $\omega: E \rightarrow R^+$  and a source vertex  $s \in V$

**Aim:**

- Compute  $\delta(s, v)$  for all  $v \in V \setminus \{s\}$
- Compute  $P(s, v)$  for all  $v \in V \setminus \{s\}$

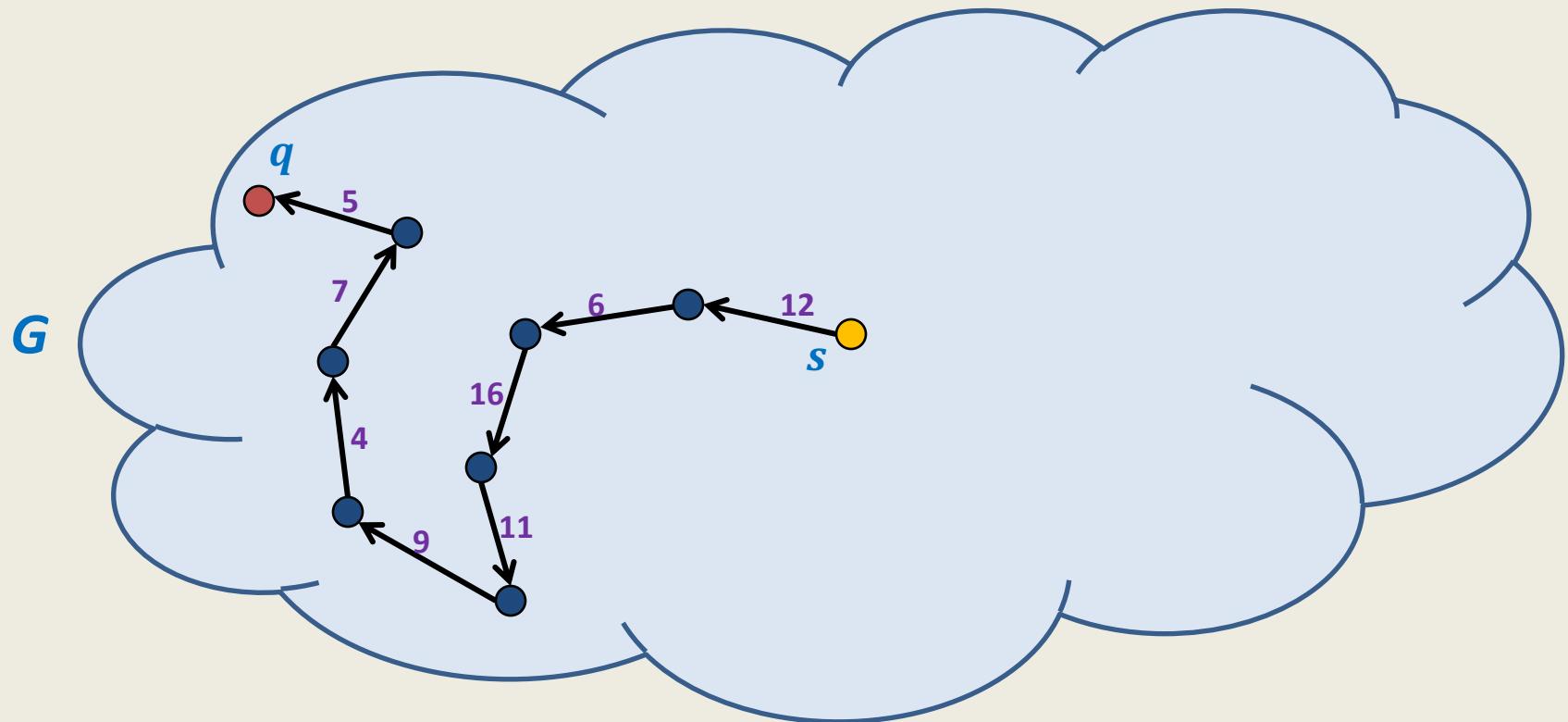
**First algorithm :** by Edsger Dijkstra in 1956

And still the best ...



I am sure you will be able to re-invent  
it yourself if you are asked right  
questions ☺ So get ready !

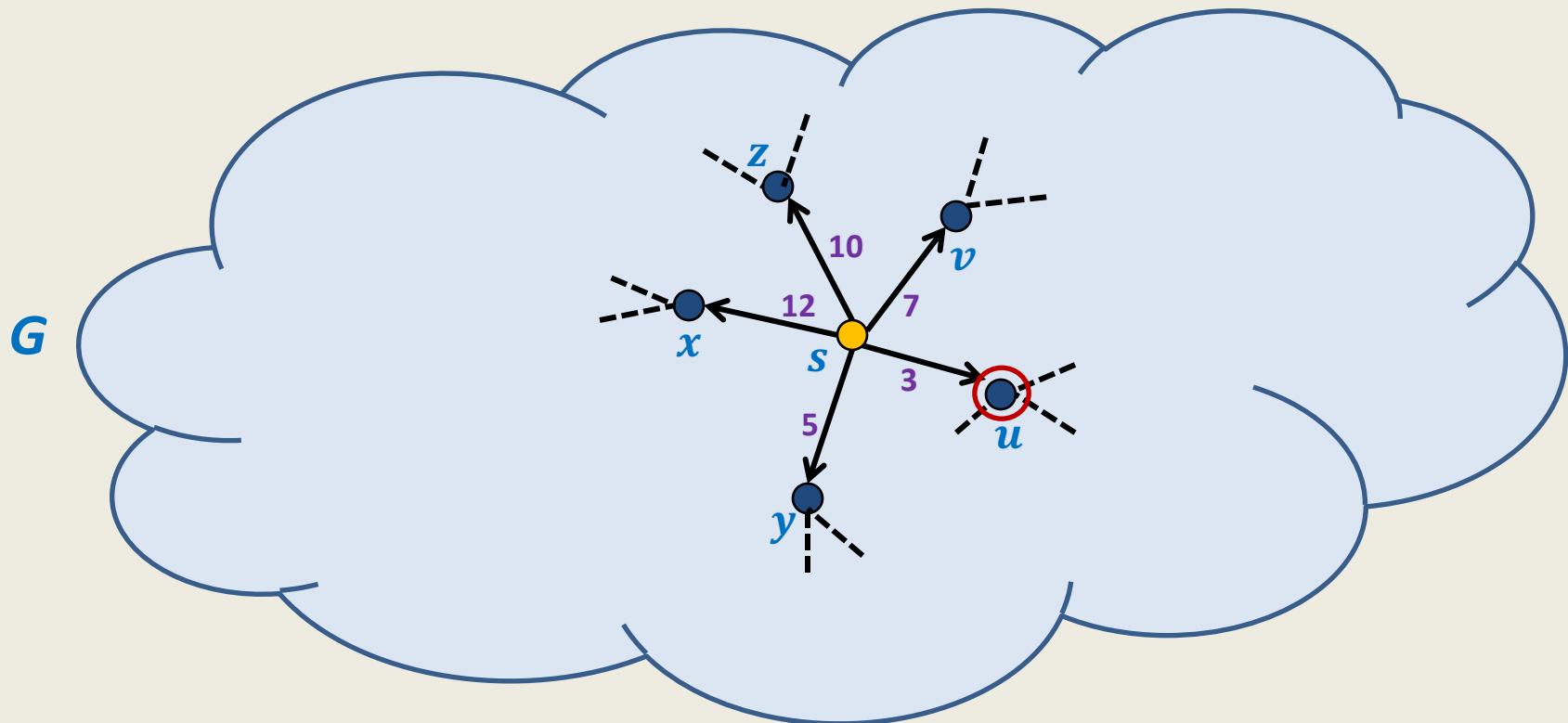
# An example to get an insight into this problem



**Inference:**

The distance to any vertex depends upon global parameters.

# An example to get an insight into this problem

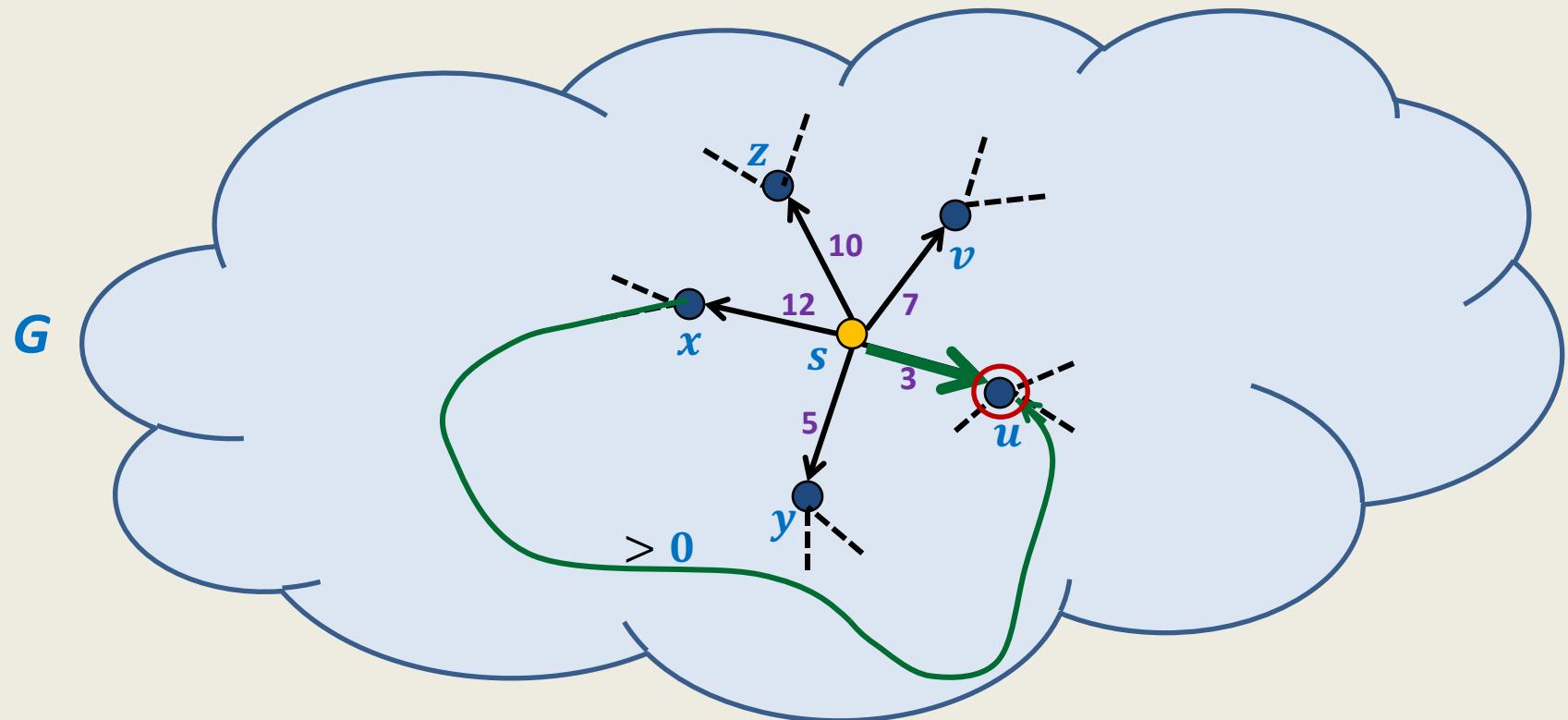


**Question:** Is there any vertex in this picture for which you are certain about the distance from  $s$  ?

**Answer:** vertex  $u$ .

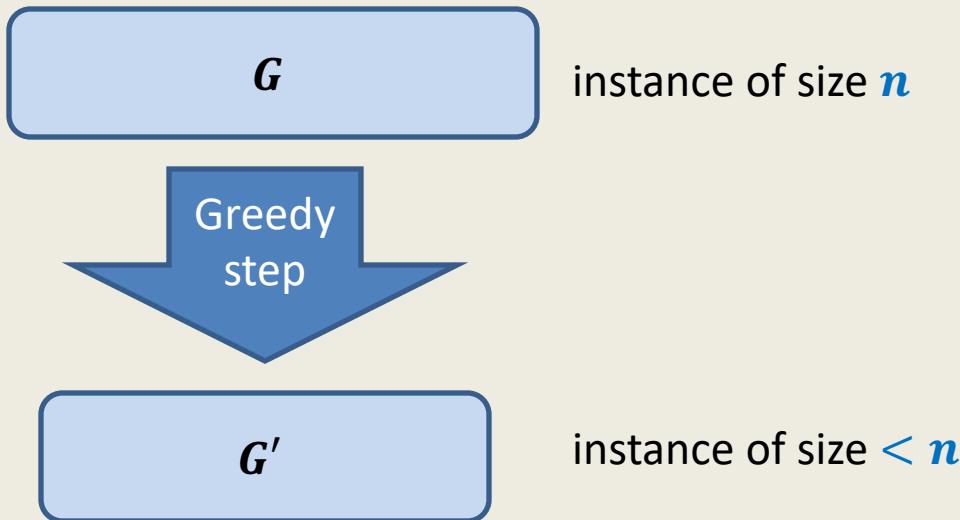
Give reasons.

# An example to get an insight into this problem



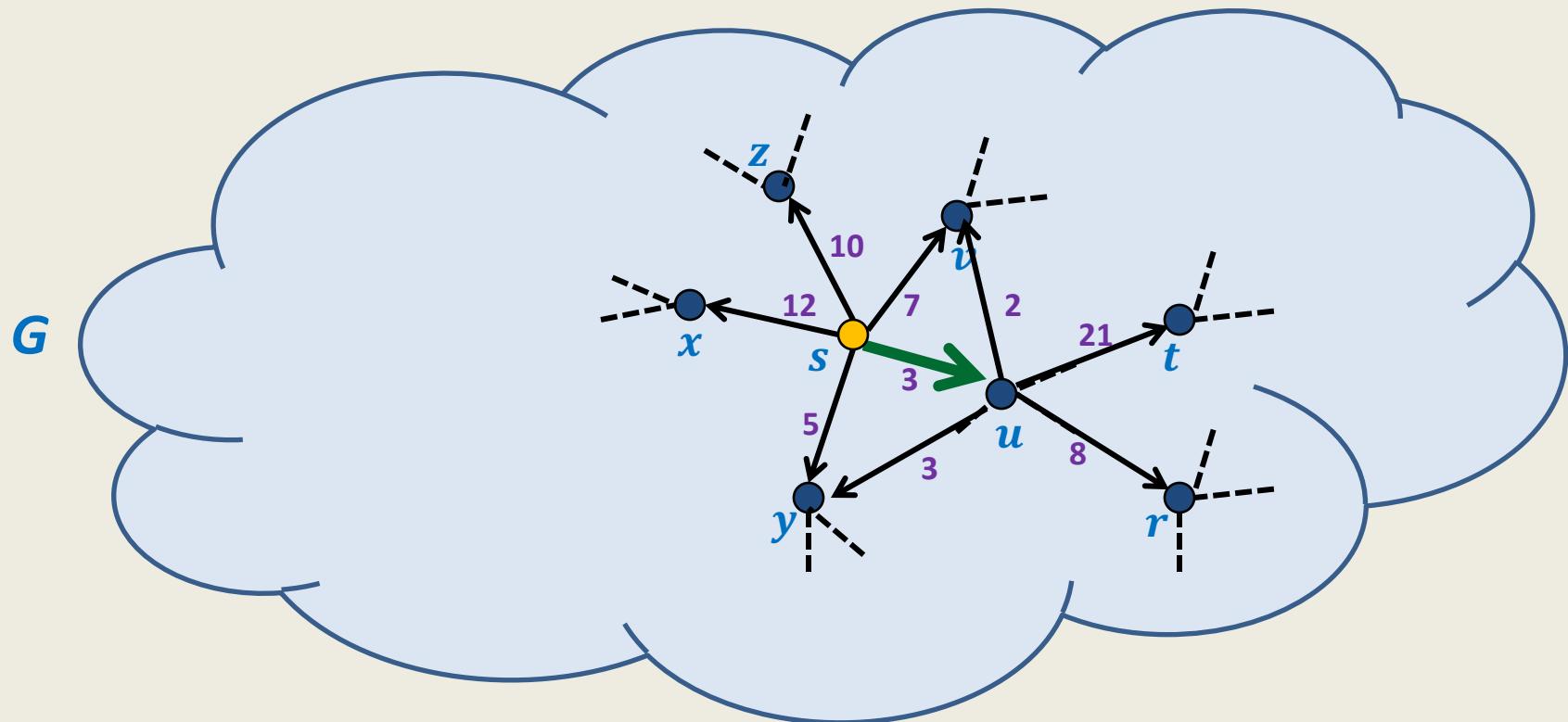
→ The shortest path to vertex  $u$  is edge  $(s, u)$ .

# Designing a greedy algorithm for shortest paths



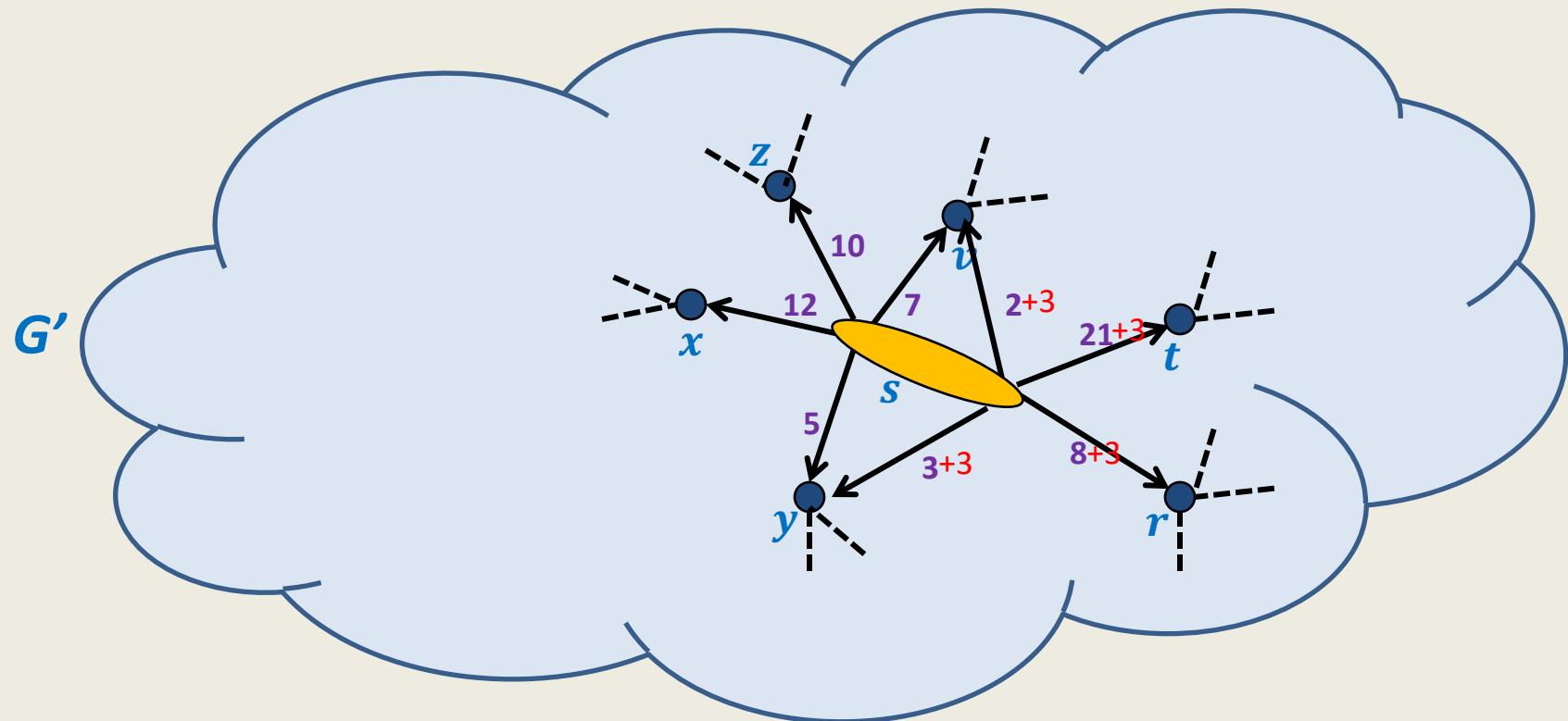
1. Establish a relation between  
and  
**shortest paths in  $G'$**

# An example to get an insight into this problem

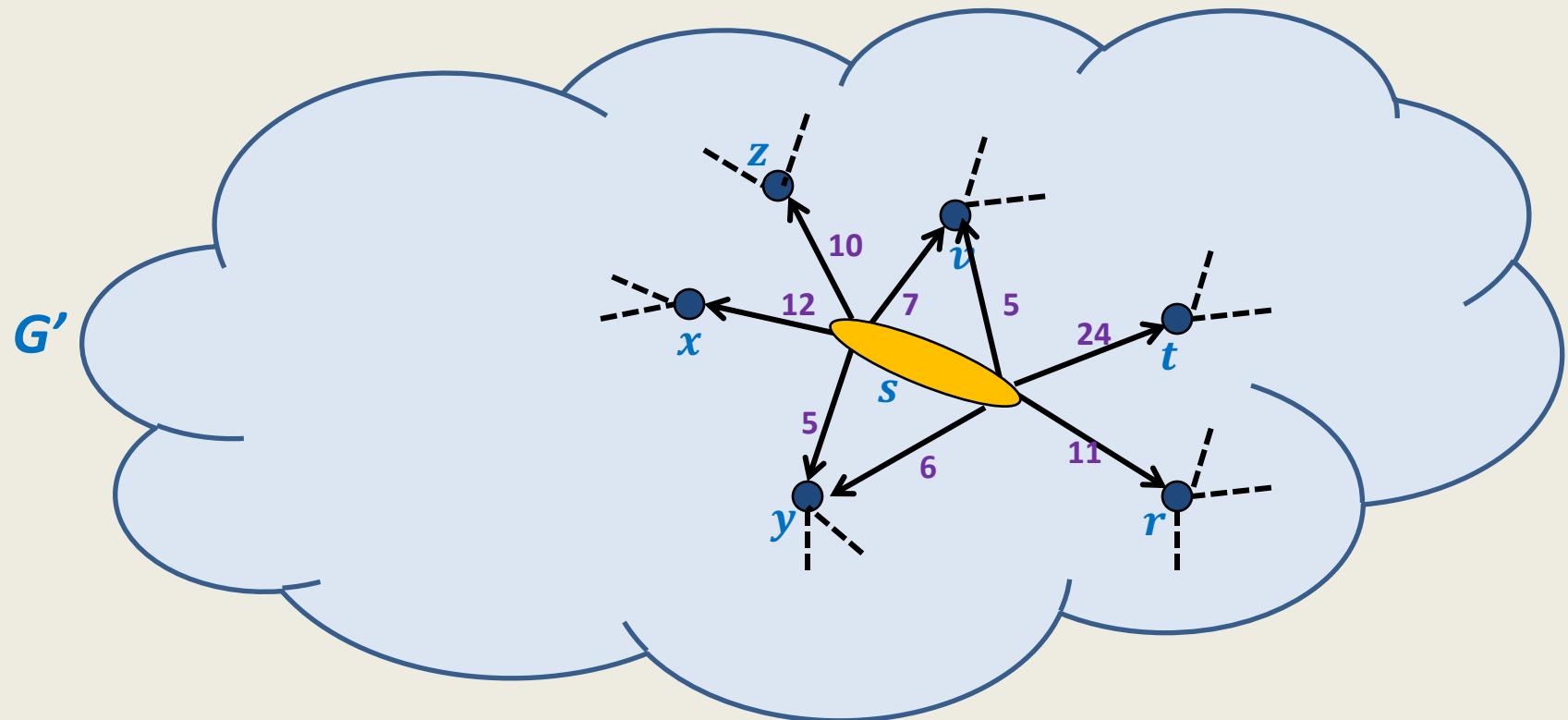


**Question:** Can you remove vertex  $u$  without affecting the distance from  $s$  ?

# An example to get an insight into this problem



# An example to get an insight into this problem



# How to compute instance $G'$

Let  $(s, u)$  be the **least weight edge** from  $s$  in  $G = (V, E)$ .

Transform  $G$  into  $G'$  as follows.

1. For each edge  $(u, x) \in E$ ,

add edge  $(s, x)$ ;

$$\omega(s, x) \leftarrow \omega(s, u) + \omega(u, x);$$

2. In case of two edges from  $s$  to any vertex  $x$

3. Remove vertex  $u$ .



**Theorem:** For each  $v \in V \setminus \{s, u\}$ ,

$$\delta_G(s, v) = \delta_{G'}(s, v)$$

→ an algorithm for **distances** from  $s$

Can you see some **negative points** of this algorithm ?

# Shortcomings of the algorithm

- **No insight** into the (beautiful) structure of shortest paths.
- **Just convinces** that we can solve the shortest paths problem in polynomial time.
- **Very few options** to improve the time complexity.

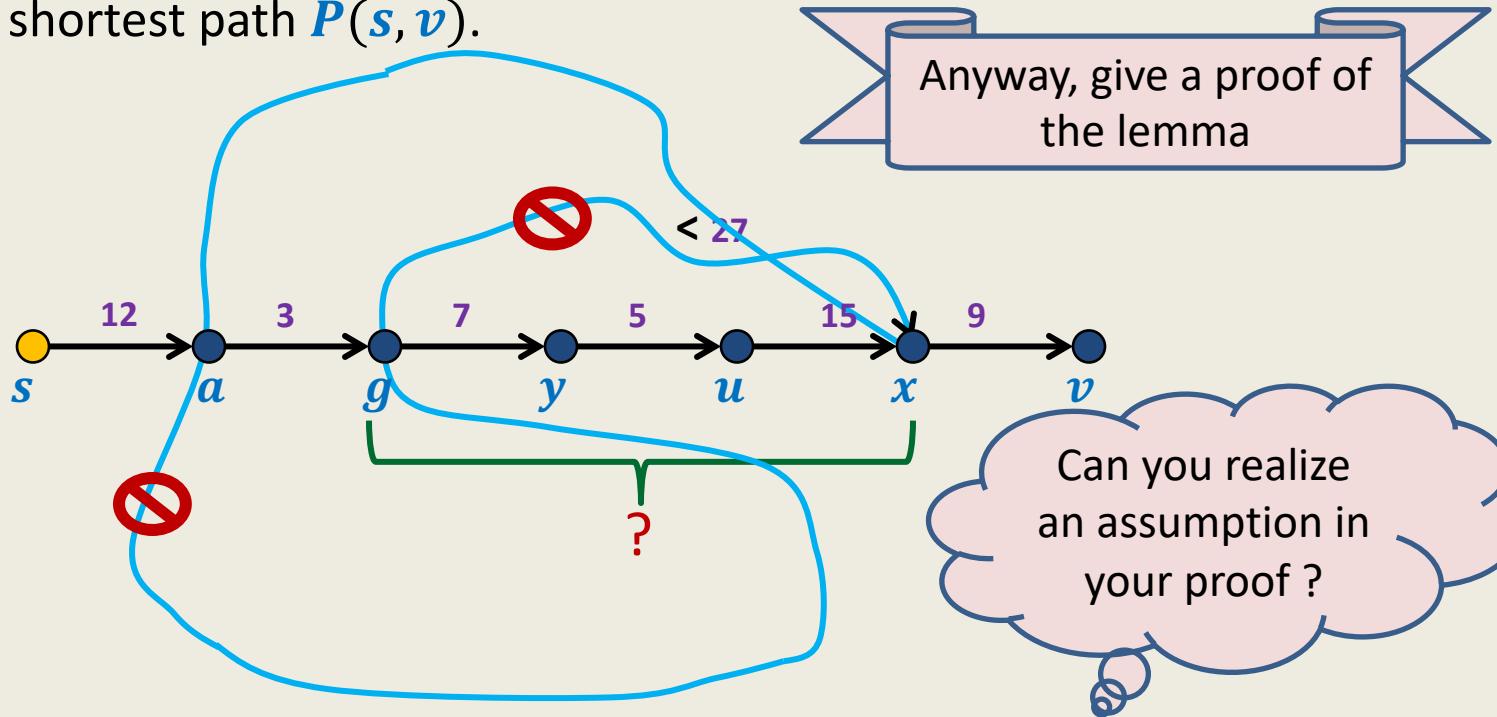
We shall now design a very **insightful** algorithm based on **properties** of shortest paths.

# PROPERTY OF A SHORTEST PATH

# Optimal subpath property

Consider any shortest path  $P(s, v)$ .

$$\delta(s, v) = 51$$



**Lemma 1:** Every **subpath** of a shortest path is also a shortest path.

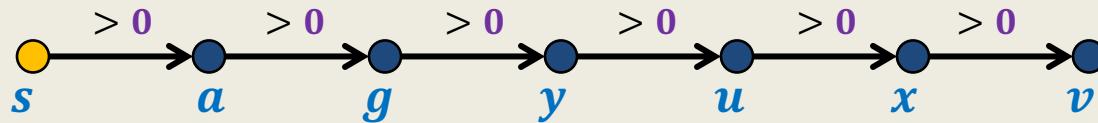
**NOTE:** Does the lemma use the fact that the edge weights are positive?

If yes, can you locate the exact place where it used it?

**Homework:** Write a complete and formal proof for **Lemma 1**

# Exploiting the positive weight on edges

Consider once again a shortest path  $P(s, v)$ .



$$\rightarrow \delta(s, a) < \delta(s, g) < \delta(s, y) < \delta(s, u) < \delta(s, x) < \delta(s, v) \quad (1)$$

→ The first nearest vertex of  $s$  must be its **neighbor**.

# More insights ...

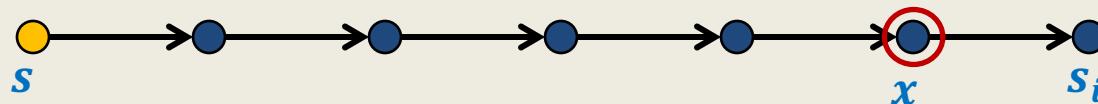
Let  $s_i$  :

$$s_0 = s.$$

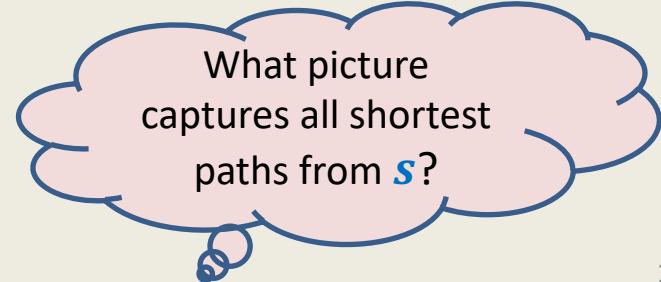
Consider the shortest path  $P(s, s_i)$ .



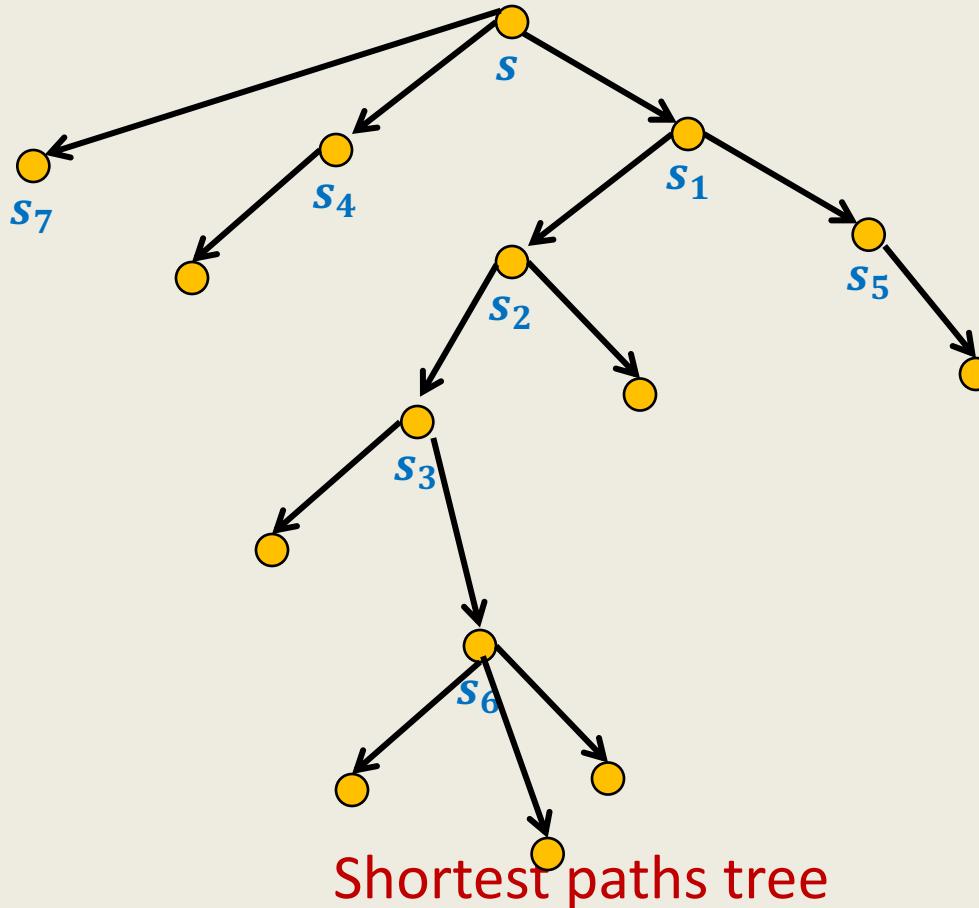
$x$  must be  $s_j$  for some  $j < i$ .



**Lemma 2:**  $P(s, s_i)$  must be of the form



# Complete picture of all shortest paths ?



# Designing the algorithm ...

**Lemma 2:**  $P(s, s_i)$  must be of the form  $s \rightsquigarrow s_j \rightarrow s_i$  for some  $j < i$ .

shortest

**Question:** Can we use **Lemma 2** to design an algorithm ?

**Incremental way** to compute shortest paths.

Ponder over it before going to the next slide 😊

# Designing the algorithm ...

**Lemma 2:**  $P(s, s_i)$  must be of the form  $s \rightsquigarrow s_j \rightarrow s_i$  for some  $j < i$ .

shortest

**Question:** Can we use **Lemma 2** to design an algorithm ?

**Incremental way** to compute shortest paths.

The next slide explains it precisely.

If shortest paths to  $s_j, j < i$  is known, we can compute  $s_i$ .

But how ?

All we know is that

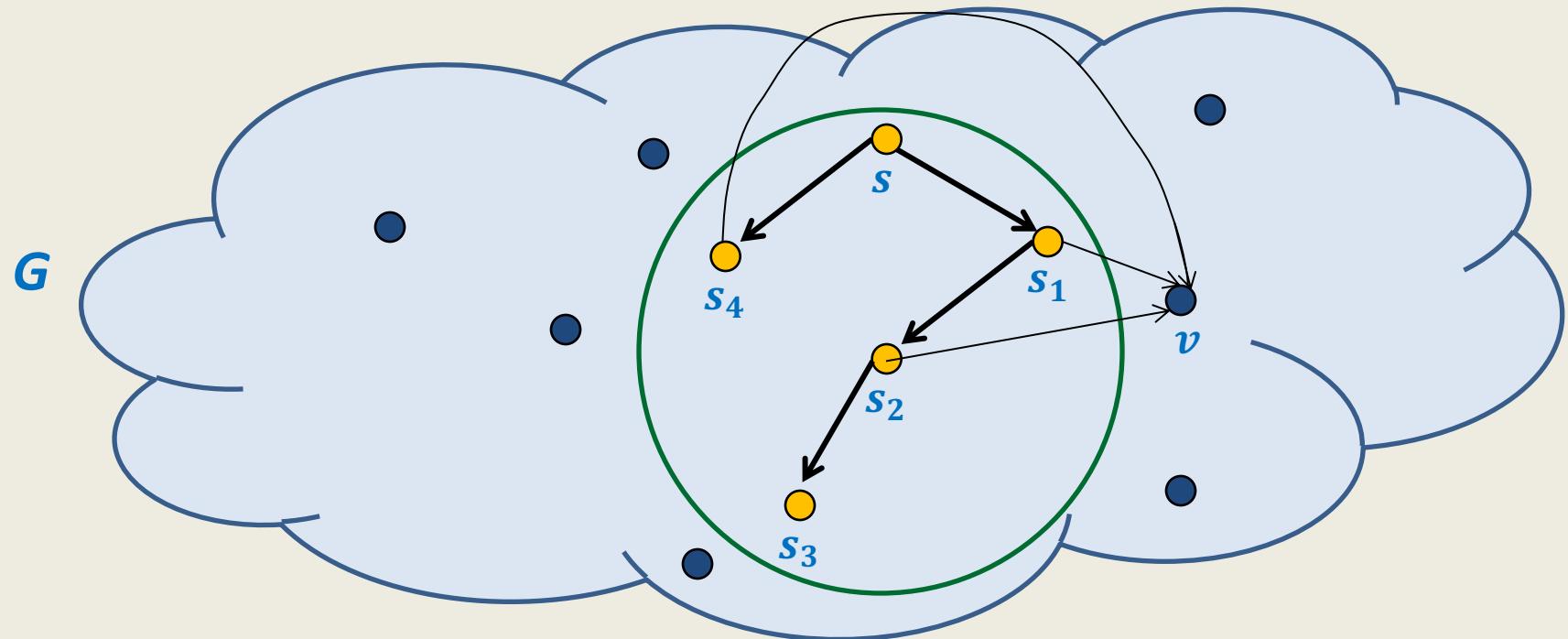
But which neighbor ?

Hint:

For each neighbor, compute some *label* based on  
**Lemma 2** s.t.  $s_i$  is the neighbour with least label.

Suppose we have computed  $s_1, s_2, \dots, s_{i-1}$ .

We can compute  $s_i$  as follows.



For each  $v \in V \setminus \{s_1, s_2, \dots, s_{i-1}\}$

$$L(v) = \min_{(s_j, v) \in E} (\delta(s, s_j) + \omega(s_j, v))$$

$s_i$  is the vertex with **minimum** value of  $L$ .

# Dijkstra's algorithm

Dijkstra-algo( $s, G$ )

{  $U \leftarrow V \setminus \{s\}$  ;

$S \leftarrow \{s\}$ ;

For  $i = 1$  to  $n - 1$  do

{     For each  $v \in U$  do

{          $L(v) \leftarrow \infty$ ;

        For each  $(x, v) \in E$  with  $x \in S$  do

$L(v) \leftarrow \min(L(v), \delta(s, x) + \omega(x, v))$

}

$y \leftarrow$  vertex from  $U$  with minimum value of  $L$ ;

$\delta(s, y) \leftarrow L(y)$ ;

    move  $y$  from  $U$  to  $S$ ;

}

In this algorithm, we first compute  $L(v)$  for each  $v \in U$  and then find the vertex with the least  $L$  value.

Try to rearrange its statements so that in the beginning of each iteration, we have  $L$  values **computed already**.

This rearrangement will be helpful for improving the running time.

So please try it on your own first before viewing the next slide.

# Dijkstra's algorithm

Dijkstra-algo( $s, G$ )

{  $U \leftarrow V$ ;  $L(v) \leftarrow \infty$  for all  $v \in U$ ;

$L(s) \leftarrow 0$ ;

For  $i = 0$  to  $n - 1$  do

{  $y \leftarrow$  vertex from  $U$  with minimum value of  $L$ ;

$\delta(s, y) \leftarrow L(y)$ ;

move  $y$  from  $U$  to  $S$ ;

For each  $v \in U$  do

{  $L(v) \leftarrow \infty$ ;

For each  $(x, v) \in E$  with  $x \in S$  do

$L(v) \leftarrow \min(L(v), \delta(s, x) + \omega(x, v))$

}

}

a lot of re-computation

# Dijkstra's algorithm

Dijkstra-algo( $s, G$ )

{  $U \leftarrow V$ ;  $L(v) \leftarrow \infty$  for all  $v \in U$ ;

$L(s) \leftarrow 0$ ;

For  $i = 0$  to  $n - 1$  do

{  $y \leftarrow$  vertex from  $U$  with minimum value of  $L$ ;

$\delta(s, y) \leftarrow L(y)$ ;

move  $y$  from  $U$  to  $S$ ;

For each  $v \in U$  do

{

For each  $(x, v) \in E$  with  $x \in S$  do

$L(v) \leftarrow \min(L(v), \delta(s, x) + \omega(x, v))$

}

}

Only neighbors of  $y$

What are the vertices whose  $L$  value may change in this iteration ?

# Dijkstra's algorithm

Dijkstra-algo( $s, G$ )

{  $U \leftarrow V$ ;  $L(v) \leftarrow \infty$  for all  $v \in U$ ;

$L(s) \leftarrow 0$ ;

For  $i = 0$  to  $n - 1$  do

{  $y \leftarrow$  vertex from  $U$  with minimum value of  $L$ ;

$\delta(s, y) \leftarrow L(y)$ ;

move  $y$  from  $U$  to  $S$ ;

For each  $(y, v) \in E$  with  $v \in U$  do

{

$L(v) \leftarrow \min(L(v), \delta(s, y) + \omega(y, v))$

}

}

1 extract-min operation

deg( $y$ ) Decrease-key operations

# Time complexity of Dijkstra's algorithm

Total number of **extract-min** operation :  $n$

Total **Decrease-key** operations :  $m$

Using **Binary heap** to maintain the set  $U$ , the time complexity:  $O(m \log n)$

**Theorem:** Given a directed graph with positive weights on edges, we can compute all shortest paths from a given vertex in  $O(m \log n)$  time.

Fibonacci heap supports **Decrease-key** in  $O(1)$  time and **extract-min** in  $O(\log n)$ .

→ Total time complexity using Fibonacci heap:  $O(m + n \log n)$

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 39**

- Integer sorting
- Counting Sort and Radix Sort

# Integer sorting

# Algorithms for Sorting $n$ elements

- **Insertion** sort:  $\mathcal{O}(n^2)$
- **Selection** sort:  $\mathcal{O}(n^2)$
- **Bubble** sort:  $\mathcal{O}(n^2)$
- **Merge** sort:  $\mathcal{O}(n \log n)$
- **Quick** sort: worst case  $\mathcal{O}(n^2)$ , average case  $\mathcal{O}(n \log n)$
- **Heap** sort:  $\mathcal{O}(n \log n)$

**Question:** What is common among these algorithms ?

**Answer:** All of them use only **comparison** operation to perform sorting.

**Theorem (we will not prove it in this course):**

Every comparison based sorting algorithm

must perform at least  $n \log n$  comparisons in the worst case.

# **Question:** Can we sort in $O(n)$ time ?

The answer depends upon

- the model of computation.
- the domain of input.

# Integer sorting

# Counting sort: algorithm for sorting integers

**Input:** An array  $A$  storing  $n$  integers in the range  $[0 \dots k - 1]$ .

$$k = O(n)$$

**Output:** Sorted array  $A$ .

**Running time:**  $O(n + k)$  in word RAM model of computation.

**Extra space:**  $O(k)$

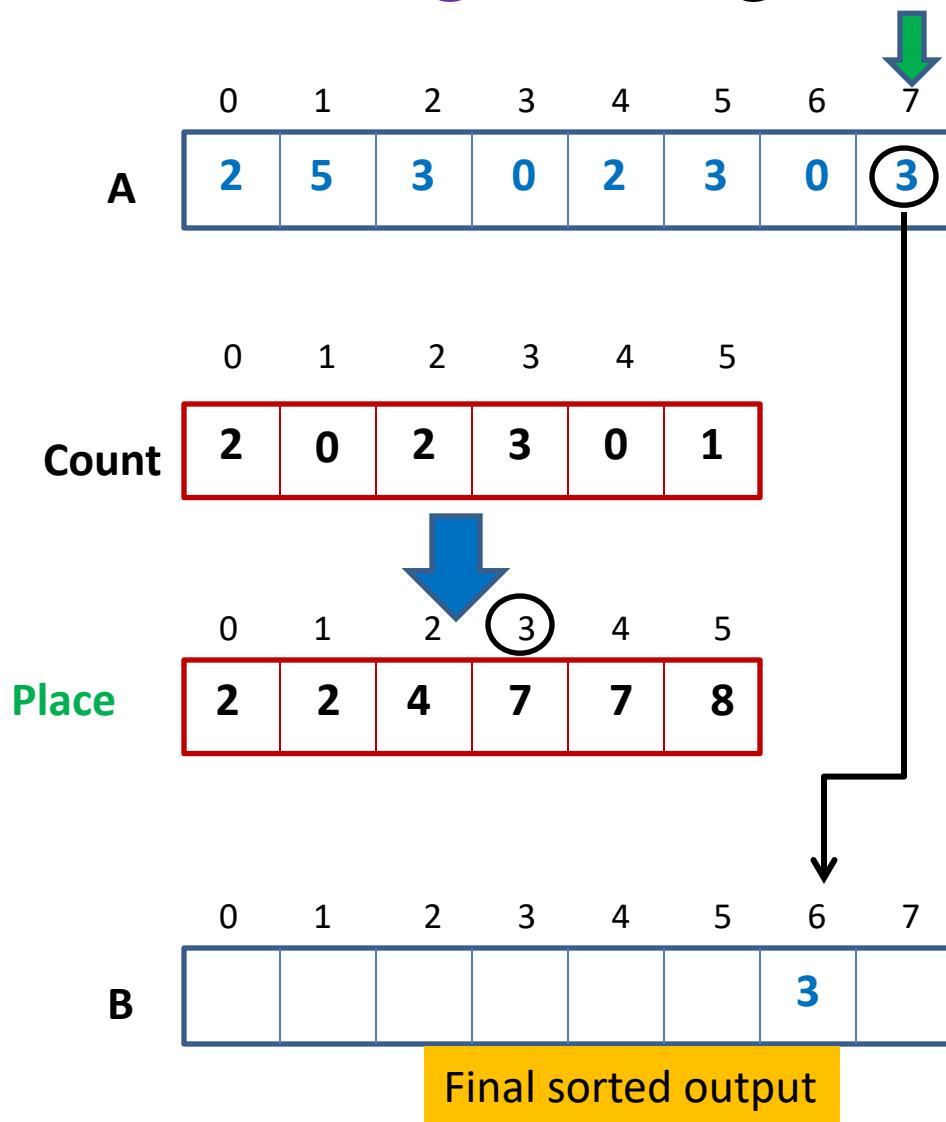
## Motivating example: Indian railways

There are 13 lacs employees.

**Aim :** To sort them list according to DOB (date of birth)

**Observation:** There are only 14600 different date of births possible.

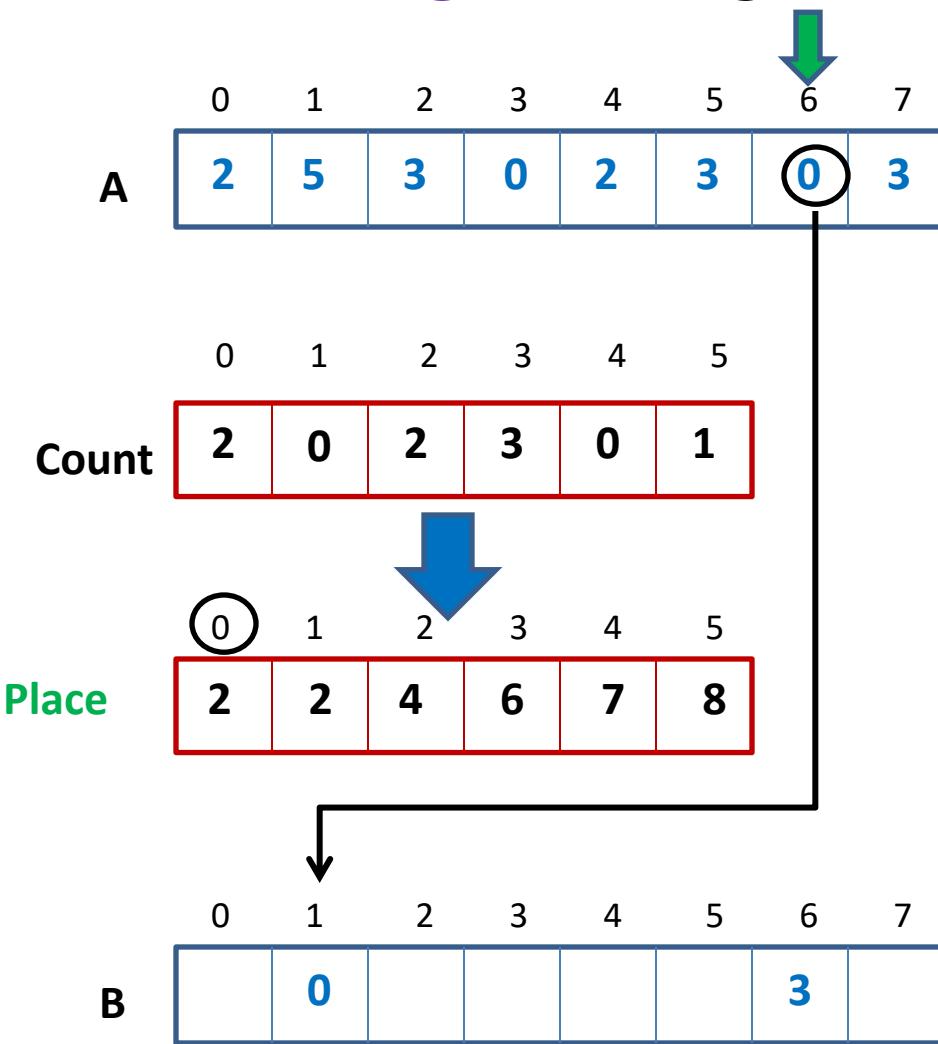
# Counting sort: algorithm for sorting integers



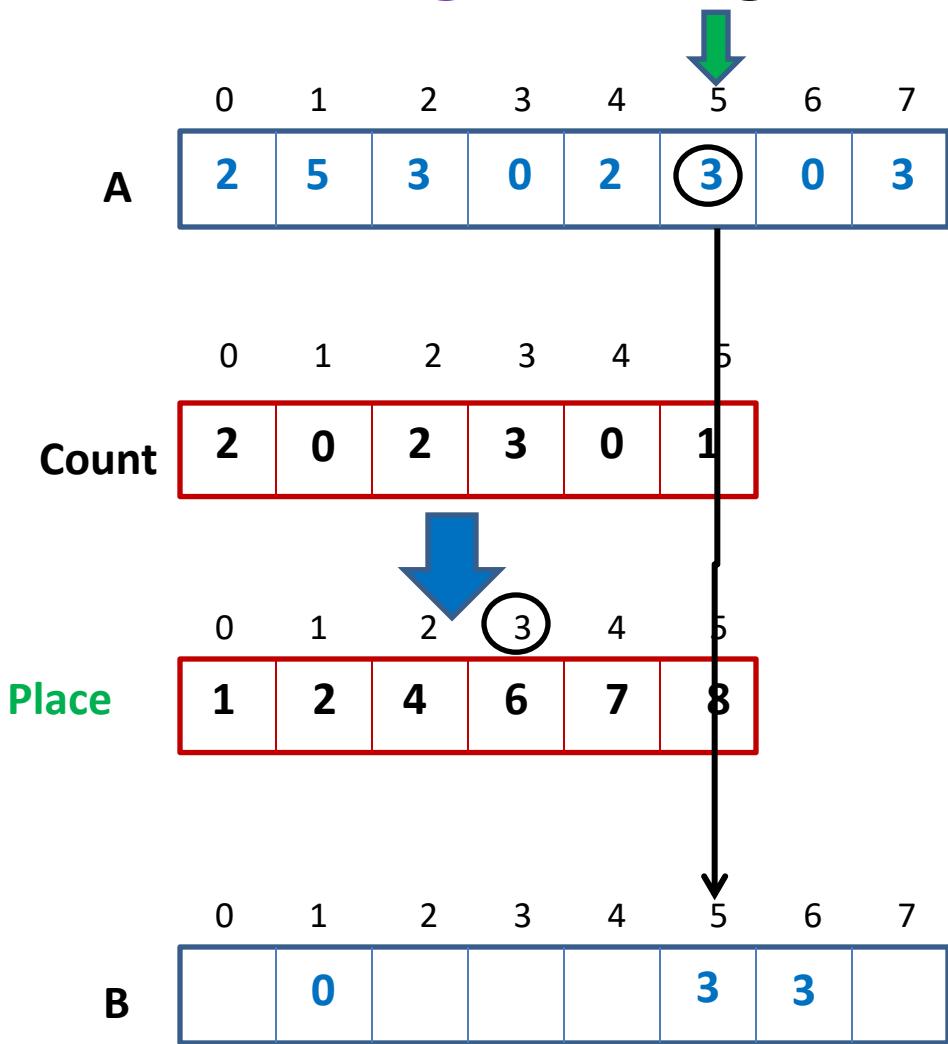
If  $A[i]=j$ ,  
where should  $A[i]$  be  
placed in B ?

Certainly after all those elements in A  
which are smaller than  $j$

# Counting sort: algorithm for sorting integers



# Counting sort: algorithm for sorting integers



# Types of sorting algorithms

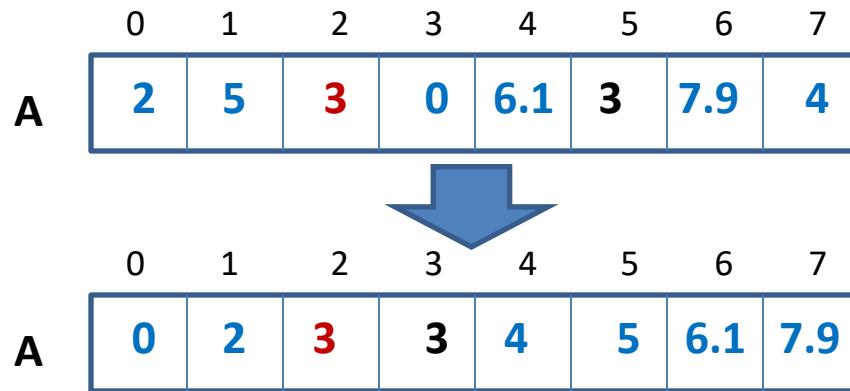
**In Place** Sorting algorithm:

A sorting algorithm which uses

**Example:** Heap sort, Quick sort.

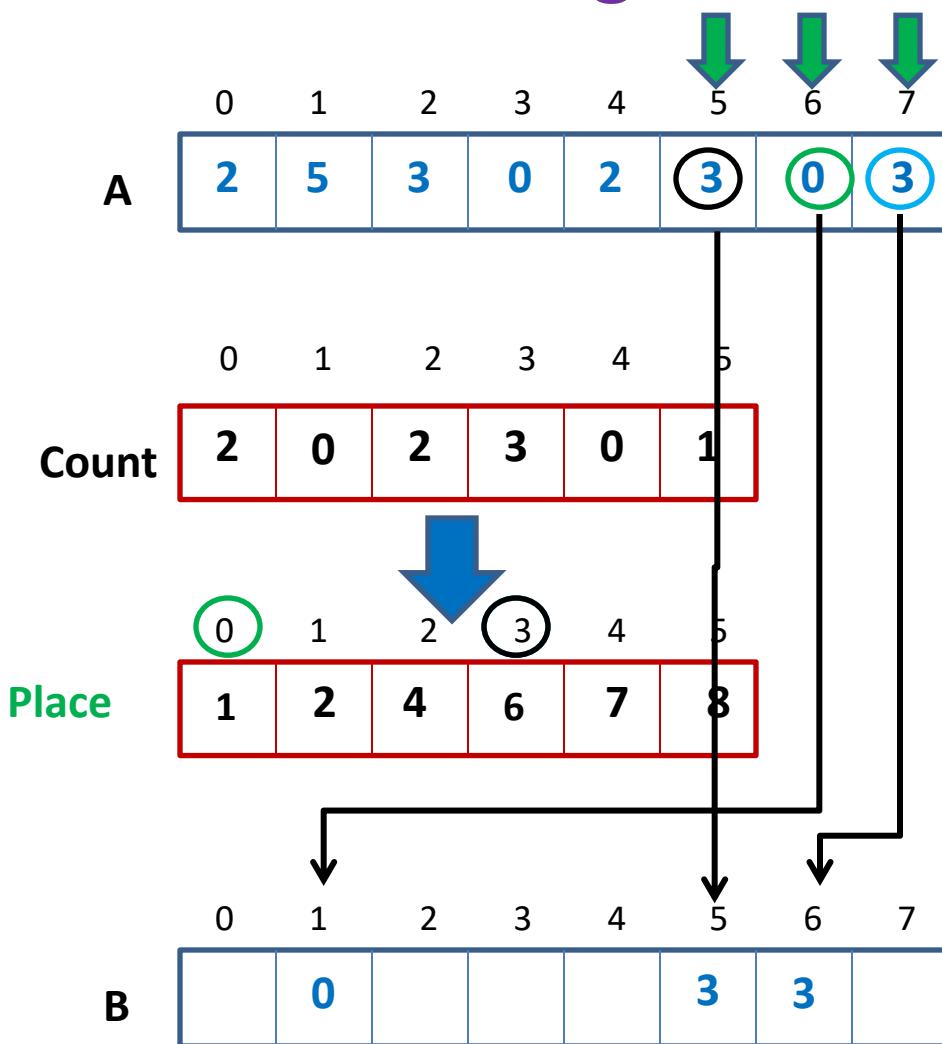
**Stable** Sorting algorithm:

A sorting algorithm which preserves



**Example:** Merge sort.

# Counting sort: a visual description



Why did we scan elements of **A** in reverse order (from index  $n - 1$  to **0**) while placing them in the final sorted array **B**?

Answer:

- To ensure that Counting sort is **stable**.
- The reason why stability is required will become clear soon 😊

# Counting sort: algorithm for sorting integers

Algorithm ( $A[0 \dots n - 1], k$ )

For  $j=0$  to  $k - 1$  do  $\text{Count}[j] \leftarrow 0$ ;

For  $i=0$  to  $n - 1$  do  $\text{Count}[A[i]] \leftarrow \text{Count}[A[i]] + 1$ ;

$\text{Place}[0] \leftarrow \text{Count}[0]$ ;

For  $j=1$  to  $k - 1$  do  $\text{Place}[j] \leftarrow \text{Place}[j - 1] + \text{Count}[j]$  ;

For  $i=n - 1$  to  $0$  do

{      $B[\text{Place}[A[i]] - 1] \leftarrow A[i]$ ;

$\text{Place}[A[i]] \leftarrow \text{Place}[A[i]] - 1$ ;

}

return  $B$ ;

Each arithmetic operations

involves  $O(\log n + \log k)$  bits

# Counting sort: algorithm for sorting integers

## Key points of Counting sort:

- It performs arithmetic operations involving  $O(\log n + \log k)$  bits  
( $O(1)$  time in word RAM).
- It is a **stable** sorting algorithm.

**Theorem:** An array storing  $n$  integers in the range  $[0..k - 1]$   
can be sorted in  $O(n+k)$  time and  
using total  $O(n+k)$  space in word RAM model.

- For  $k \leq n$ ,  
→ For  $k = n^t$ ,  
(too bad for  $t > 1$ . ☹)

## Question:

How to sort  $n$  integers in the range  $[0..n^t]$  in

# Radix Sort

# Digits of an integer

507266

No. of digits = 6

value of digit  $\in \{0, \dots, 9\}$

The binary representation of the decimal number 507266 is shown as a sequence of bits: 1011000101011111. Four horizontal brackets above the sequence group the bits into four groups of three, indicating they represent nibbles or bytes. This visualizes how integers can have different digit counts depending on the base used.

1011000101011111

No. of digits = 4

value of digit  $\in \{0, \dots, 15\}$

It is up to us how we define digit ?

# Radix Sort

**Input:** An array  $\mathbf{A}$  storing  $n$  integers, where

- (i) each integer has exactly  $d$  digits.
- (ii) each digit has value  $< k$
- (iii)  $k < n$ .

**Output:** Sorted array  $\mathbf{A}$ .

**Running time:**

$O(dn)$  in word RAM model of computation.

**Extra space:**

$O(n + k)$

**Important points:**

- makes use of a count sort.
- Heavily relies on the fact that count sort is a stable sort algorithm.

# Demonstration of Radix Sort through example

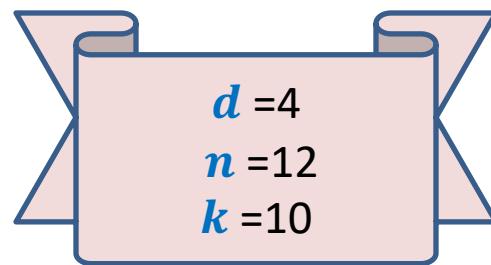
A



2	0	1	2
1	3	8	5
4	9	6	1
5	8	1	0
2	3	7	3
6	2	3	9
9	6	2	4
8	2	9	9
3	4	6	5
7	0	9	8
5	5	0	1
9	2	5	8



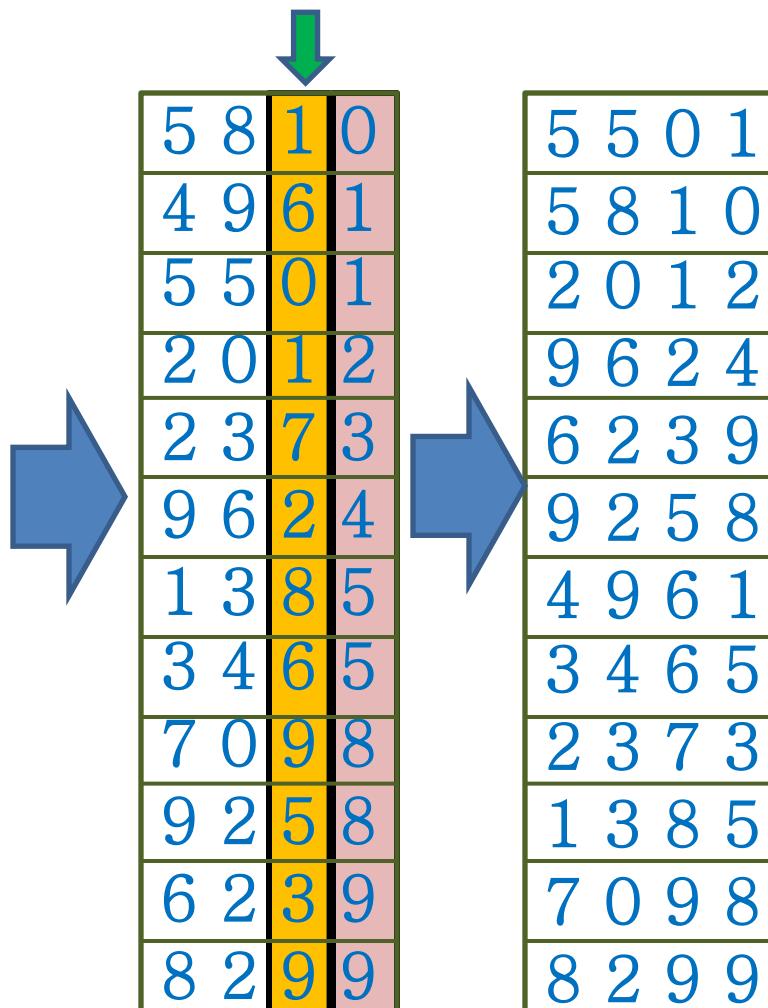
5	8	1	0
4	9	6	1
5	5	0	1
2	0	1	2
2	3	7	3
9	6	2	4
1	3	8	5
3	4	6	5
7	0	9	8
9	2	5	8
6	2	3	9
8	2	9	9



# Demonstration of Radix Sort through example

A

2	0	1	2
1	3	8	5
4	9	6	1
5	8	1	0
2	3	7	3
6	2	3	9
9	6	2	4
8	2	9	9
3	4	6	5
7	0	9	8
5	5	0	1
9	2	5	8



5	5	0	1
5	8	1	0
2	0	1	2
9	6	2	4
6	2	3	9
9	2	5	8
4	9	6	1
3	4	6	5
2	3	7	3
1	3	8	5
7	0	9	8
8	2	9	9

# Demonstration of Radix Sort through example

A

2 0 1 2
1 3 8 5
4 9 6 1
5 8 1 0
2 3 7 3
6 2 3 9
9 6 2 4
8 2 9 9
3 4 6 5
7 0 9 8
9 2 5 8
5 5 0 1
9 2 5 8

5 8 1 0
4 9 6 1
5 5 0 1
2 0 1 2
2 3 7 3
9 6 2 4
1 3 8 5
3 4 6 5
7 0 9 8
9 2 5 8
6 2 3 9
8 2 9 9
4 9 6 1
3 4 6 5
5 5 0 1
9 2 5 8
8 2 9 9

↓

5	5	0	1
5	8	1	0
2	0	1	2
9	6	2	4
6	2	3	9
9	2	5	8
4	9	6	1
3	4	6	5
2	3	7	3
1	3	8	5
7	0	9	8
8	2	9	9

2 0 1 2
7 0 9 8
6 2 3 9
9 2 5 8
2 3 7 3
1 3 8 5
3 4 6 5
5 5 0 1
9 6 2 4
5 8 1 0
4 9 6 1

# Demonstration of Radix Sort through example

A

2 0 1 2	5 8 1 0
1 3 8 5	4 9 6 1
4 9 6 1	5 5 0 1
5 8 1 0	2 0 1 2
2 3 7 3	9 6 2 4
6 2 3 9	6 2 3 9
9 6 2 4	9 2 5 8
8 2 9 9	1 3 8 5
3 4 6 5	4 9 6 1
7 0 9 8	3 4 6 5
5 5 0 1	9 2 5 8
9 2 5 8	6 2 3 9

5 5 0 1	5 8 1 0
5 8 1 0	2 0 1 2
2 0 1 2	9 6 2 4
9 6 2 4	6 2 3 9
6 2 3 9	9 2 5 8
9 2 5 8	1 3 8 5
1 3 8 5	4 9 6 1
4 9 6 1	3 4 6 5
3 4 6 5	7 0 9 8
7 0 9 8	2 3 7 3

↓

2 0 1 2	7 0 9 8	6 2 3 9	9 2 5 8	8 2 9 9	2 3 7 3	1 3 8 5	3 4 6 5	4 9 6 1	5 5 0 1	5 8 1 0	6 2 3 9	7 0 9 8	8 2 9 9	9 2 5 8	9 6 2 4
2 0 1 2	7 0 9 8	6 2 3 9	9 2 5 8	8 2 9 9	2 3 7 3	1 3 8 5	3 4 6 5	4 9 6 1	5 5 0 1	5 8 1 0	6 2 3 9	7 0 9 8	8 2 9 9	9 2 5 8	9 6 2 4
2 0 1 2	7 0 9 8	6 2 3 9	9 2 5 8	8 2 9 9	2 3 7 3	1 3 8 5	3 4 6 5	4 9 6 1	5 5 0 1	5 8 1 0	6 2 3 9	7 0 9 8	8 2 9 9	9 2 5 8	9 6 2 4
2 0 1 2	7 0 9 8	6 2 3 9	9 2 5 8	8 2 9 9	2 3 7 3	1 3 8 5	3 4 6 5	4 9 6 1	5 5 0 1	5 8 1 0	6 2 3 9	7 0 9 8	8 2 9 9	9 2 5 8	9 6 2 4
2 0 1 2	7 0 9 8	6 2 3 9	9 2 5 8	8 2 9 9	2 3 7 3	1 3 8 5	3 4 6 5	4 9 6 1	5 5 0 1	5 8 1 0	6 2 3 9	7 0 9 8	8 2 9 9	9 2 5 8	9 6 2 4

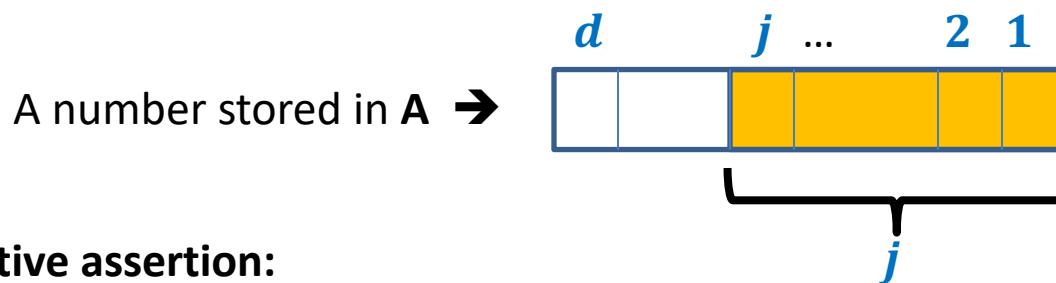
Can you see where we are exploiting the fact that  
Countsort is a **stable** sorting algorithm ?

# Radix Sort

**RadixSort(A[0... $n - 1$ ],  $d$ ,  $k$ )**

```
{  For  $j=1$  to  $d$  do
    Execute CountSort(A,  $k$ ) with  $j$ th digit as the key;
    return A;
}
```

**Correctness:**



**Inductive assertion:**

At the end of  $j$ th iteration, array **A** is sorted according to the last  $j$  digits.

During the induction step, you will have to use the fact that **Countsort** is a **stable** sorting algorithm.

# Radix Sort

**RadixSort(A[0... $n - 1$ ],  $d$ ,  $k$ )**

```
{  For  $j=1$  to  $d$  do
    Execute CountSort(A, $k$ ) with  $j$ th digit as the key;
    return A;
}
```

**Time complexity:**

- A single execution of **CountSort(A, $k$ )** runs in  $O(n + k)$  time and  $O(n + k)$  space.
- For  $k < n$ ,
  - a single execution of **CountSort(A, $k$ )** runs in  $O(n)$  time.
  - Time complexity of radix sort =  $O(dn)$ .
- → Extra space used =  $O(n)$

**Question:** How to use Radix sort to sort  $n$  integers in range  $[0..n^t]$  in  $O(tn)$  time and  $O(n)$  space ?

**Answer:**

The diagram illustrates the conversion of bit widths. On the left, a blue box labeled "1 bit" has a green arrow pointing to a white box labeled "log  $n$  bits". Below this, a larger blue box labeled " $t \log n$  bits" has a green arrow pointing to a white box labeled " $t$  bits".

$d$	$k$	Time complexity
$t \log n$	2	$O(tn \log n)$ 😕
$t$	$n$	$O(tn)$ 😊



# Power of the word RAM model

- **Very fast** algorithms for **sorting integers**:  
**Example:**  $n$  integers in range  $[0..n^{10}]$  in  $O(n)$  time and  $O(n)$  space ?
- **Lesson:**  
**Do not** always go after **Merge sort** and **Quick sort** when input is integers.
- **Interesting programming exercise:**  
**Compare** **Quick sort** with **Radix sort** for sorting **long** integers.

# **Data Structures and Algorithms**

**(ESO207)**

## **Lecture 40**

- **Search data structure for integers : Hashing**

# Data structures for searching

in **O(1)** time

# Motivating Example

**Input:** a given set  $S$  of 1009 positive integers

**Aim:** Data structure for searching

## Example

```
{  
123, 579236, 1072664, 770832456778, 61784523, 100004503210023, 19,  
762354723763099, 579, 72664, 977083245677001238, 84, 100004503210023,  
...  
}
```

**Data structure :** Array storing  $S$  in sorted order

**Searching :** Binary search

$\mathbf{O}(\log |S|)$  time

Can we perform  
search in  $\mathbf{O}(1)$  time ?

# Problem Description

$\mathcal{U}$  :

$\mathcal{S} \subseteq \mathcal{U}$ ,

$n = |\mathcal{S}|$ ,

$n \ll m$

A search query:

**Aim:** A data structure for a given set  $\mathcal{S}$

that can facilitate search in  $O(1)$  time

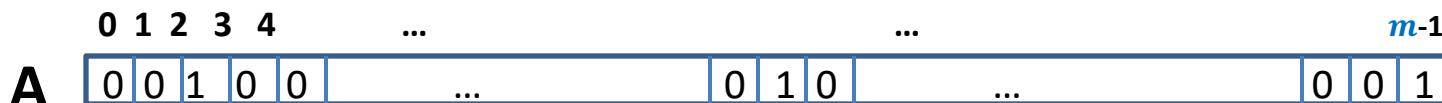
# A trivial data structure for $O(1)$ search time

Build a 0-1 array  $A$  of size  $m$  such that

$A[i] = 1$  if  $i \in S$ .

$A[i] = 0$  if  $i \notin S$ .

Time complexity for searching an element in set  $S$  :  $O(1)$ .



This is a totally Impractical data structure because  $n \ll m$  !

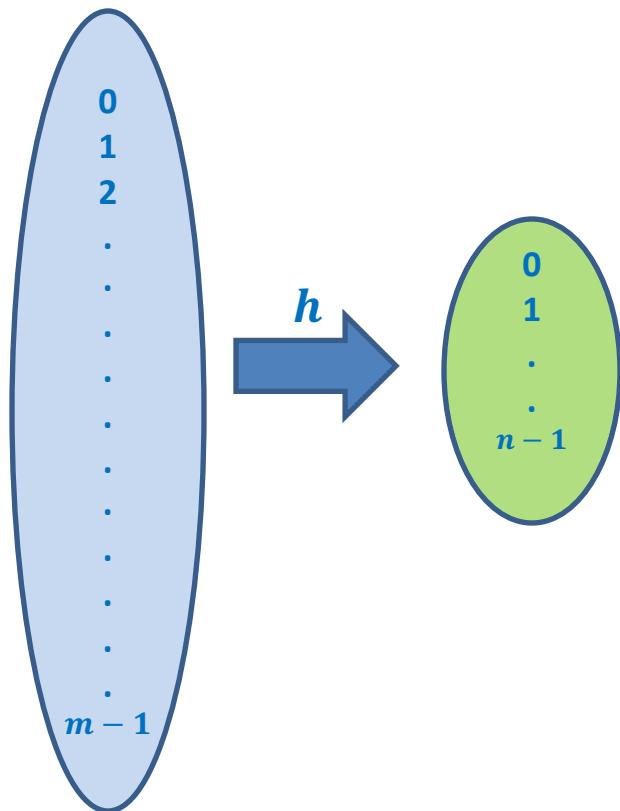
Example:  $n$  = few thousands,  $m$  = few trillions.

**Question:**

Can we have a data structure of  $O(n)$  size that can answer a search query in  $O(1)$  time ?

**Answer:** Hashing

# Hash function



## Hash function:

$h$  is a mapping from  $U$  to  $\{0, 1, \dots, n - 1\}$  with the following characteristics.

- Space required for  $h$
- $h(i)$  computable in  $O(1)$

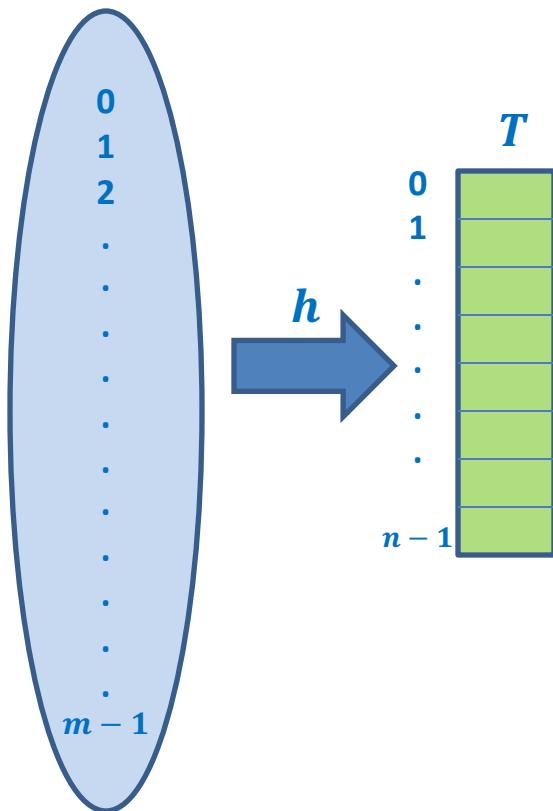
## Example:

## Hash value:

For a given hash function  $h$ , and  $i \in U$ .

$h(i)$  is called hash value of  $i$

# Hash function, hash value



## Hash function:

$h$  is a mapping from  $U$  to  $\{0, 1, \dots, n - 1\}$  with the following characteristics.

- Space required for  $h$ : a few words.
- $h(i)$  computable in  $O(1)$  time in word RAM.

Example:  $h(i) = i \bmod n$

## Hash value:

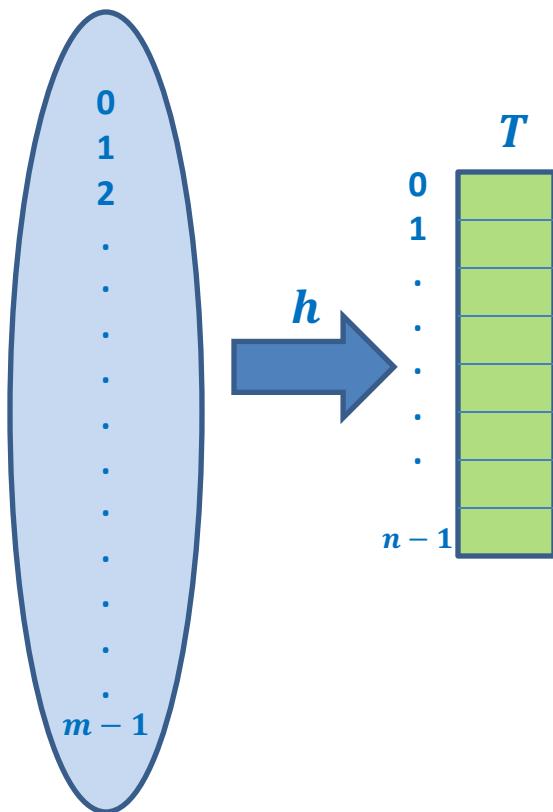
For a given hash function  $h$ , and  $i \in U$ .

$h(i)$  is called hash value of  $i$

## Hash Table:

An array  $T[0 \dots n - 1]$

# Hash function, hash value, hash table



## Hash function:

$h$  is a mapping from  $U$  to  $\{0, 1, \dots, n - 1\}$  with the following characteristics.

- **Space** required for  $h$  : a few **words**.
- $h(i)$  computable in **O(1)** time in **word RAM**.

Example:  $h(i) = i \bmod n$

## Hash value:

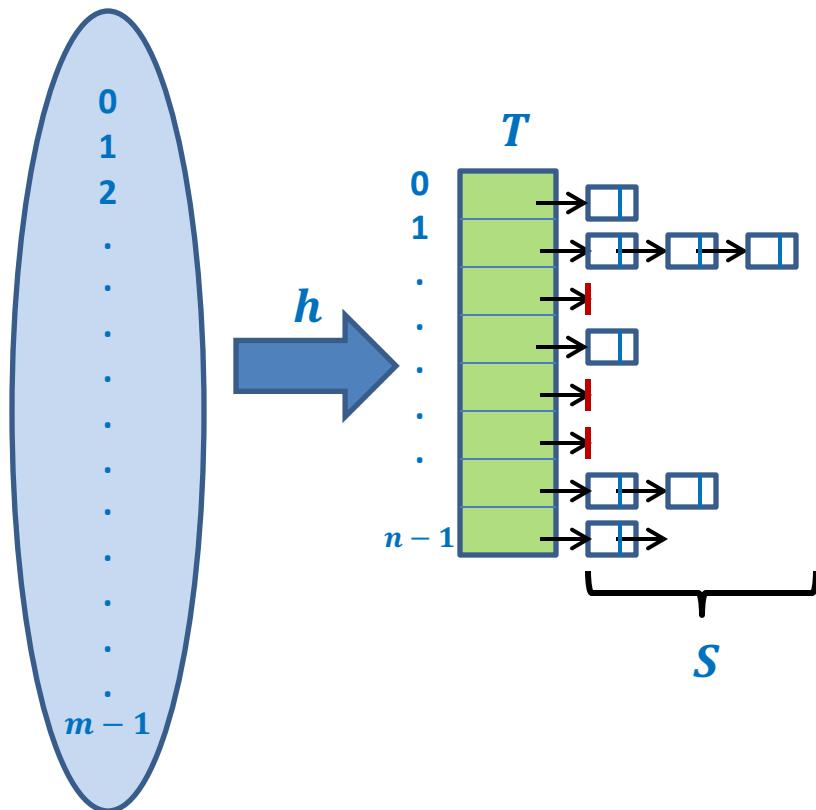
For a given hash function  $h$ , and  $i \in U$ .

$h(i)$  is called hash value of  $i$

## Hash Table:

An array  $T[0 \dots n - 1]$

# Hash function, hash value, hash table



## Hash function:

**$h$**  is a mapping from  $\mathcal{U}$  to  $\{0, 1, \dots, n - 1\}$  with the following characteristics.

- **Space** required for  $h$  : a few **words**.
  - $h(i)$  computable in **O(1)** time in **word RAM**.

**Example:**  $h(i) = i \bmod n$

## Hash value:

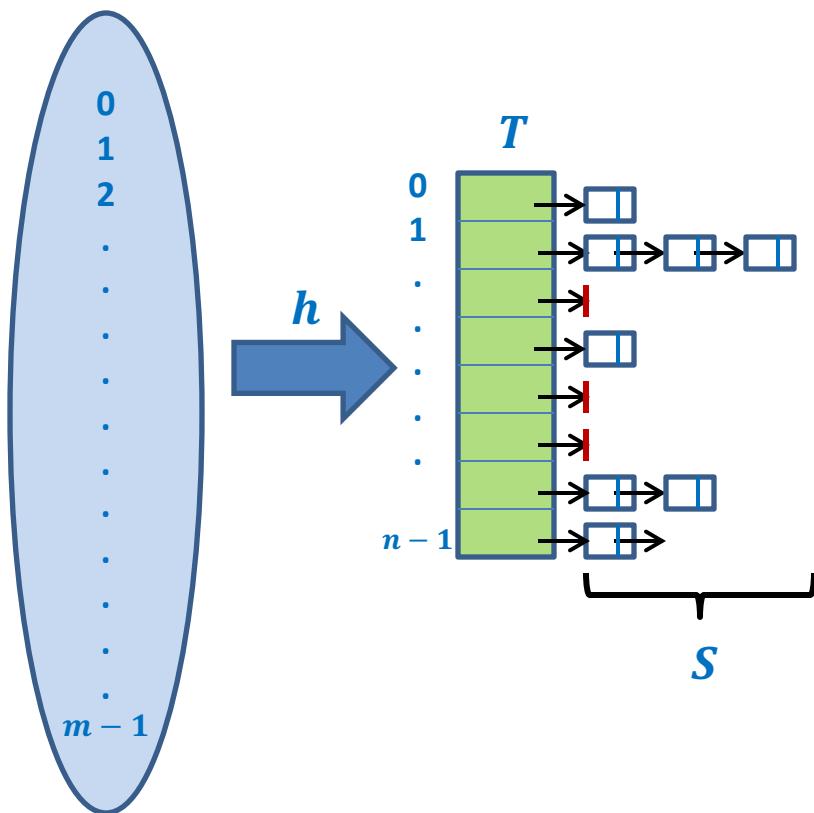
For a given hash function  $\textcolor{blue}{h}$ , and  $i \in \textcolor{blue}{U}$ .

**$h(i)$**  is called hash value of  **$i$**

## Hash Table:

An array  $T[0 \dots n - 1]$  of pointers storing  $S$ .

# Hash function, hash value, hash table



**Question:**

How to use  $(h, T)$  for searching an element  $i \in U$ ?

**Answer:**

$$k \leftarrow h(i);$$

Search element  $i$  in the list  $T[k]$ .

**Time complexity for searching:**

**O**(length of the **longest** list in  $T$ ).

# Efficiency of Hashing depends upon hash function

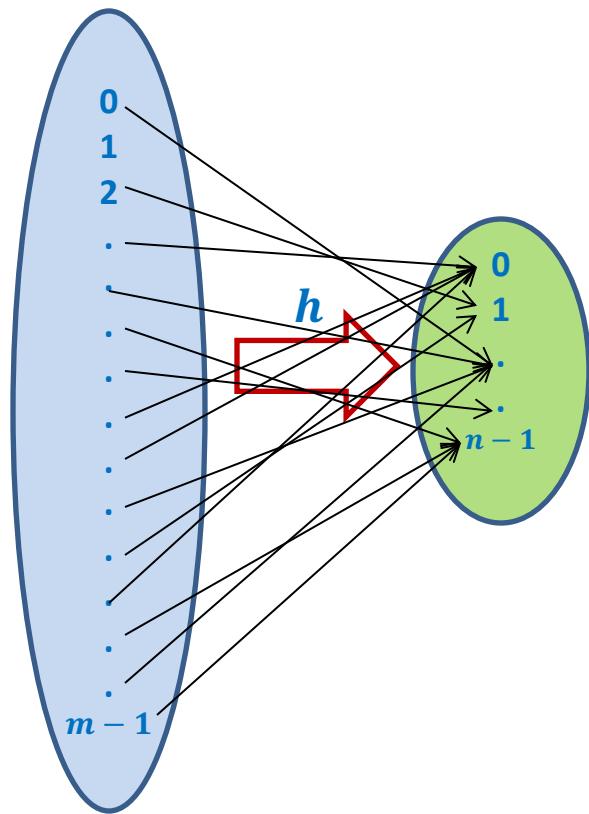
A hash function  $h$  is good if it can **evenly** distributes  $S$ .

**Aim:** To search for a good hash function for a given set  $S$ .



There **can not be** any hash function  $h$  which is good for every  $S$ .

# Hash function, hash value, hash table



For every  $h$ , there exists a subset of  $\left\lceil \frac{m}{n} \right\rceil$  elements from  $U$  which are hashed to same value under  $h$ .

So we can always construct a subset  $S$  for which all elements have same hash value

- All elements of this set  $S$  are present in a single list of the hash table  $T$  associated with  $h$ .
- $O(n)$  worst case search time.

# Why does hashing work **so well** in Practice ?

$$h(i) = i \bmod n$$

Because the set **S** is usually a **uniformly random** subset of **U**.

Let us do a **theoretical analysis**  
**to prove** this fact.

# Why does hashing work so well in Practice ?

$$h(x) = x \bmod n$$

Let  $y_1, y_2, \dots, y_n$  denote  $n$  elements selected randomly uniformly from  $U$  to form  $S$ .

**Question:**

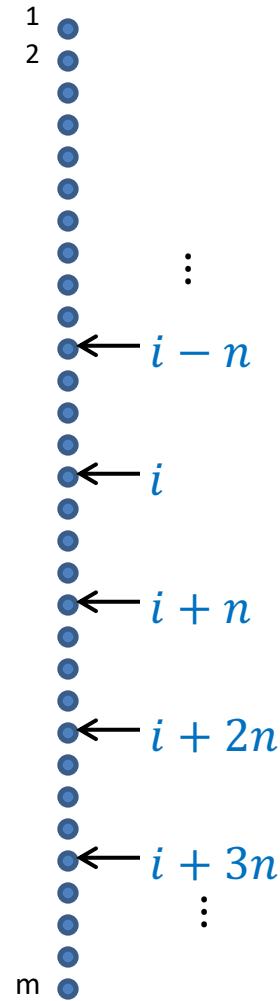
What is expected number of elements of  $S$  colliding with  $y_1$ ?

**Answer:** Let  $y_1$  takes value  $i$ .

$P(y_j \text{ collides with } y_1) = ??$

How many possible values can  $y_j$  take ? m - 1

How many possible values can collide with  $i$  ?



# Why does hashing work so well in Practice ?

$$h(x) = x \bmod n$$

Let  $y_1, y_2, \dots, y_n$  denote  $n$  elements selected randomly uniformly from  $U$  to form  $S$ .

**Question:**

What is expected number of elements of  $S$  colliding with  $y_1$ ?

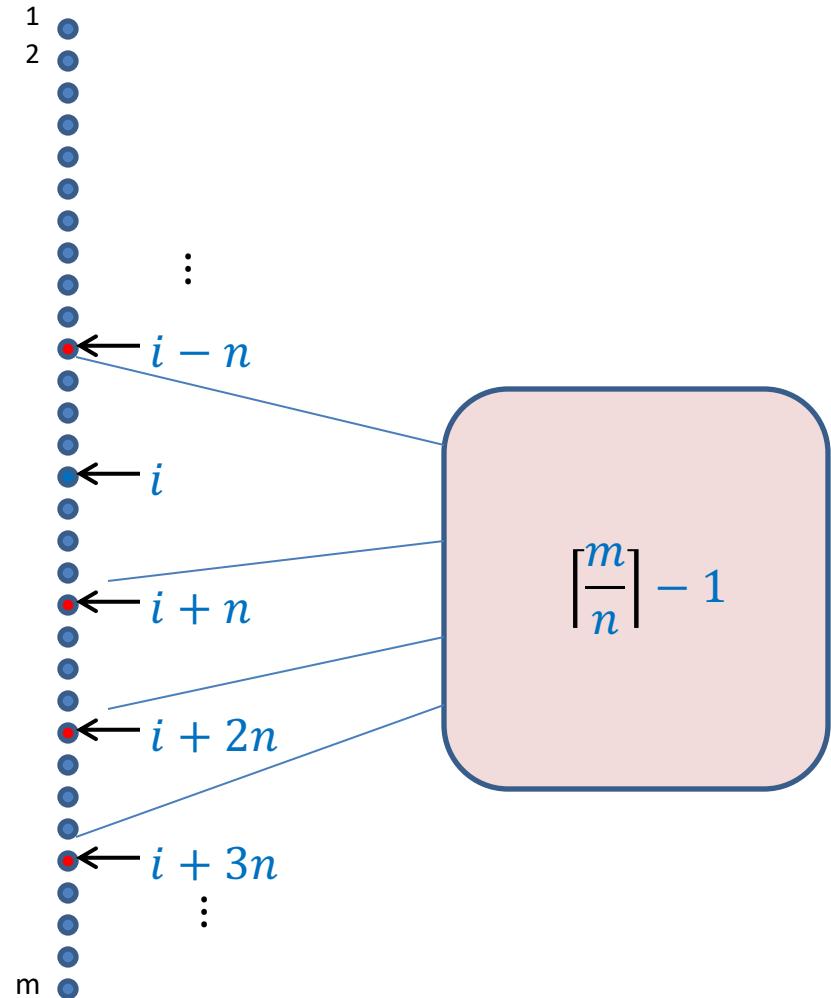
**Answer:** Let  $y_1$  takes value  $i$ .

$P(y_j \text{ collides with } y_1) =$

$$\frac{\left[\frac{m}{n}\right] - 1}{m-1}$$

Expected number of elements of  $S$  colliding with  $y_1$  =

$$\begin{aligned} &= \frac{\left[\frac{m}{n}\right] - 1}{m-1} (n-1) \\ &= O(1) \end{aligned}$$



# Why does hashing work **so well** in Practice ?

## Conclusion

1.  $h(i) = i \bmod n$  works so well because  
for a uniformly random subset of  $U$ ,  
the **expected** number of collision at an index of  $T$  is  $O(1)$ .

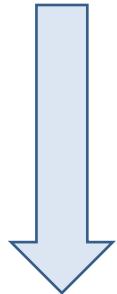
It is easy to fool this hash function such that it achieves  $O(s)$  search time.  
(do it as a simple exercise).

This makes us think:

“How can we achieve worst case  $O(1)$  search time for a given set  $S$ .”

# Hashing: theory

1953



1984

$U : \{0, 1, \dots, m - 1\}$

$S \subseteq U,$

$n = |S|,$

**Theorem [FKS, 1984]:**

A hash table and hash function can be computed in average  $O(n)$  time for a given  $S$  s.t.

Space :  $O(n)$

Query time: worst case  $O(1)$

**Ingredients :**

- elementary knowledge of **prime numbers**.
- The algorithms use **simple randomization**.

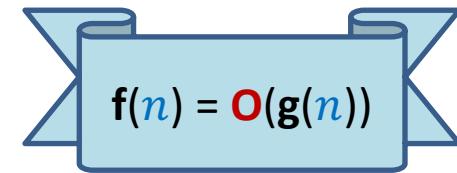
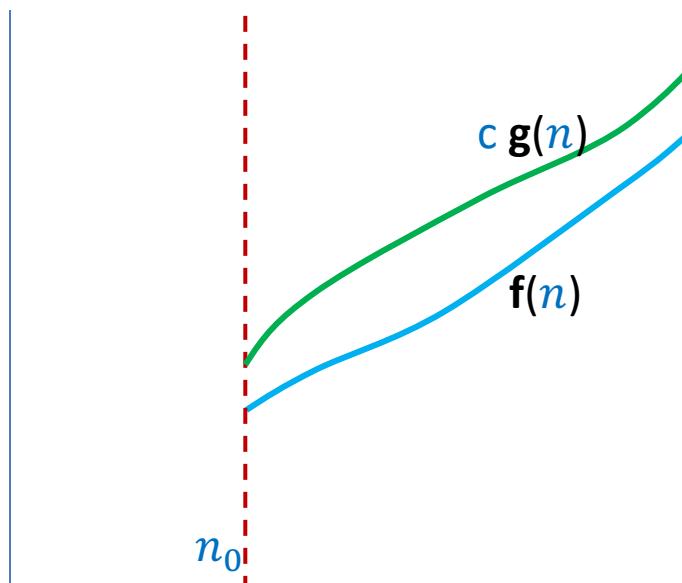
# Order notation

**Definition:** Let  $f(n)$  and  $g(n)$  be any two increasing functions of  $n$ .

$f(n)$  is said to be

if there exist constants  $c$  and  $n_0$  such that

$$f(n) \leq c g(n) \quad \text{for all } n > n_0$$



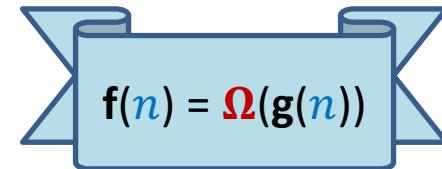
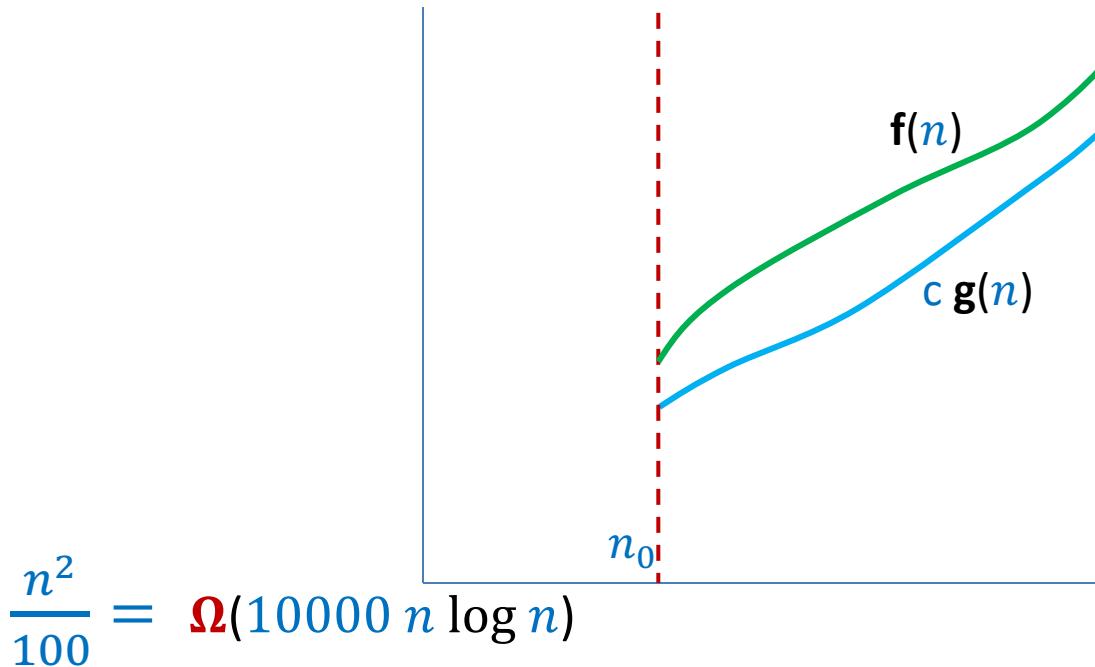
# Order notation extended

**Definition:** Let  $f(n)$  and  $g(n)$  be any two increasing functions of  $n$ .

$f(n)$  is said to be

if there exist constants  $c$  and  $n_0$  such that

$$f(n) \geq c g(n) \quad \text{for all } n > n_0$$



# Order notation extended

## Observations:

- $f(n) = O(g(n))$  if and only if  $g(n) = \Omega(f(n))$

## One more Notation:

If  $f(n) = O(g(n))$  and  $g(n) = O(f(n))$ , then

$$g(n) = \Theta(f(n))$$

## Examples:

- $\frac{n^2}{100} = \Theta(10000 n^2)$

- Time complexity of Quick Sort is  $\Omega(n \log n)$
- Time complexity of Merge sort is  $\Theta(n \log n)$