

## Section 4

# Concurrent Objects - Correctness, Progress and Efficiency

## Concurrent Objects

- ❑ A **concurrent object** is a data object "shared" by concurrently executing processes.
- ❑ Each object has a type, which defines a set of possible values and a set of primitive operations that provide the only means to create and manipulate that object.
- ❑ Concurrent objects have been proposed as building blocks for the construction of more complex multi-processing systems.
  - Leads to a system that is simple and well-structured.

# Basic Concurrent Objects

## Multi-Writer (MW) Register

- All processes are allowed to execute update operations to the register

## Single-Writer (SW) Register

- Only one process is allowed to execute update operations to the register.

## Register size

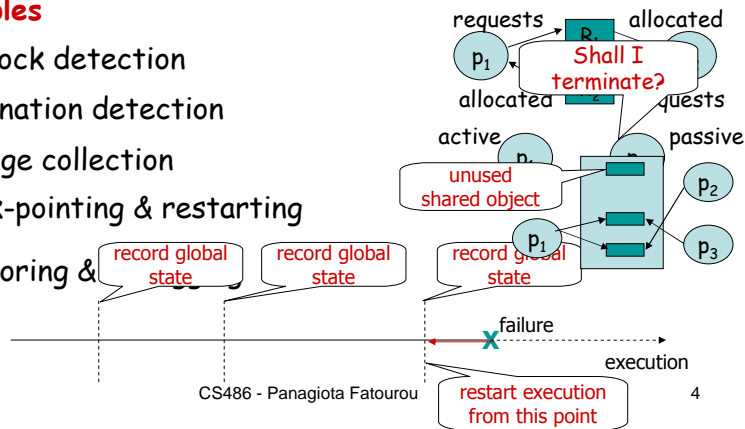
- A register is of bounded size if the set of values it may store is bounded. In the opposite case, we say that the register size is unbounded.
- Different machines support different sets of these operations in hardware.
- The hardware, then, guarantees that these operations are executed atomically.

# Global State Predicate Evaluation

In many problems of distributed computing, some action should take place only if some global predicate evaluates to TRUE.

## Examples

- deadlock detection
- termination detection
- garbage collection
- check-pointing & restarting
- monitoring &



## Debugging Distributed Algorithms

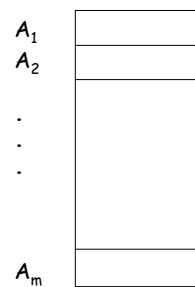
- ❑ The biggest difficulty in proving that a distributed algorithm is correct is the necessity to argue based on non-consistent versions of the shared variables.
- ❑ Calculating consistent views of the shared variables facilitates the verification of correctness of distributed algorithms.
- Calculating such consistent views is not however an easy problem. ☹

CS486 - Panagiota Fatourou

5

## Snapshot Object

- ❑ A **snapshot object** is a concurrent object that is composed by an array of  $m$  components,  $A_1, \dots, A_m$ , each capable to store a value from some set.
- ❑ The snapshot object supports two operations:
  - **UPDATE**( $i,v$ ): writes  $v$  to component  $A_i$ .
  - **SCAN**: returns a vector of consistent values, one for each component.
- ❑ Snapshot objects provide "consistent views" of a set of shared variables given that **UPDATE** operations may concurrently change the values of these variables.



Snapshot Object A

CS486 - Panagiota Fatourou

6

# Snapshot Objects

## Single-Writer Snapshot

- Only process  $p_i$  is allowed to execute UPDATES on component  $A_i$ .

## Multi-Writer Snapshot

- All processes are allowed to execute UPDATES on every component.

- Snapshots simplify the task of designing distributed algorithms! 😊
- They are too complicated to be provided by the hardware. ☹️

# The Problem

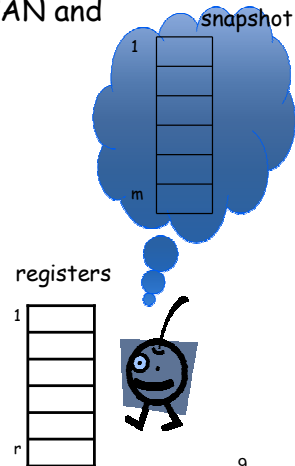
- Can we implement snapshot objects in systems that provide only r/w registers?
- If yes, how efficient can their implementation be in terms of:
  - Time complexity?
  - Space complexity?

## Implementations of Snapshots from r/w registers

- Use the registers to store the values of the snapshot components.
- Provide algorithms to implement SCAN and UPDATE.

### Efficiency

- **Step Complexity (SCAN or UPDATE)**
  - Maximum number of steps executed by any process in any execution in order to perform an operation.
- **Space Complexity**
  - Number (and size) of registers needed.



CS486 - Panagiota Fatourou

9

## Correctness and Termination

### Wait-Freedom

- Each process finishes the execution of its operation within a bounded number of its own steps.
- ➔ Wait-free algorithms are highly fault-tolerant!

### Linearizability Intuitively

- In each execution  $a$ , each SCAN and UPDATE should have the same response as if it has executed serially (or atomically) at some point in its execution interval. This point is called **linearization point** of the operation.

### and slightly more formally:

- For each (parallel) execution  $a$  produced by the implementation, there is a serial execution  $\sigma$  of the operations (SCAN and UPDATE) executed in  $a$ , s.t.:
  - each operation has the same response in  $\sigma$  as in  $a$ , and
  - $\sigma$  respects the partial order determined by the execution intervals of the operations.

CS486 - Panagiota Fatourou

10

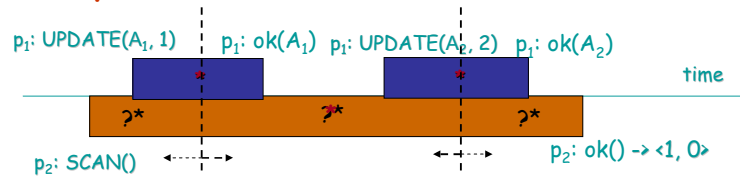
## Linearizability - Even more Formally

We say that an execution  $a$  is linearizable if it is possible to do all of the following:

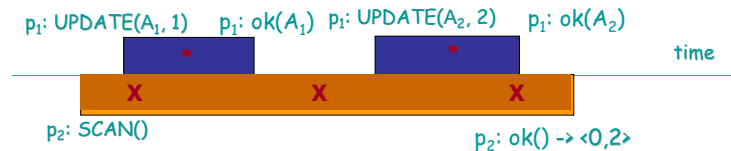
1. For each completed operation  $\pi$ , to insert a linearization point  $*\pi$  somewhere between  $\pi$ 's invocation and response in  $a$ .
2. To select a subset  $\Phi$  of the incomplete operations, and for each operation  $\rho$  in  $\Phi$ :
  - to select a response, and
  - to insert a linearization point  $*\rho$  somewhere after  $\rho$ 's invocation in  $a$ .
3. These operations and responses should be selected and these linearization points should be inserted, so that, in the sequential execution constructed by serially executing each operation at the point that its linearization point has been inserted, the response of each operation is the same as that in  $a$ .

## Linearizability - Examples

### Example of linearizable execution

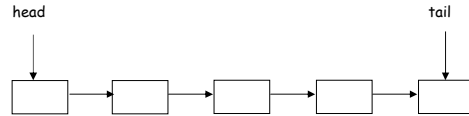


### Example of non-linearizable execution



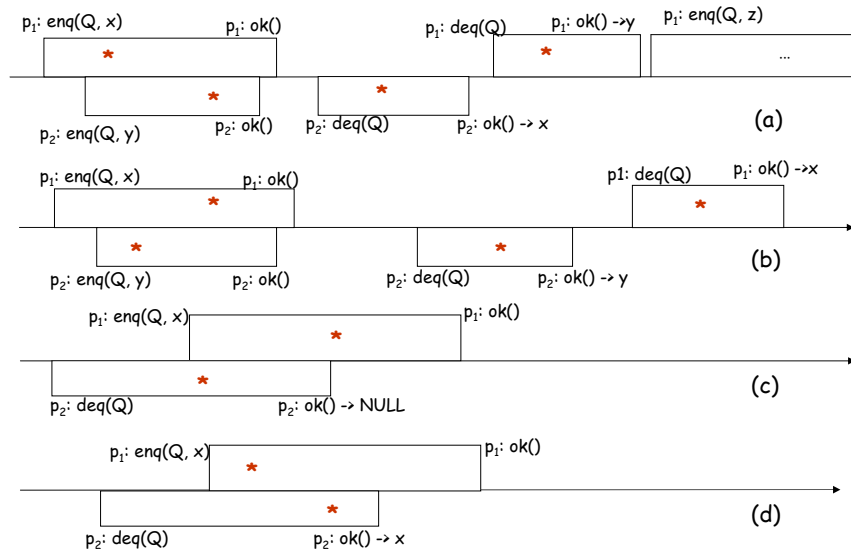
## Linearizability - Examples

- ❑ A **FIFO queue**  $Q$  supports the following two operations:
  - $\text{enq}(Q, v)$ : add an element with value  $v$  as the last element of  $Q$
  - $\text{deq}(Q)$ : deletes and returns the first element of  $Q$
- ❑ In a concurrent FIFO queue, several processes try to apply  $\text{enq}()$  and  $\text{deq}()$  operations concurrently.

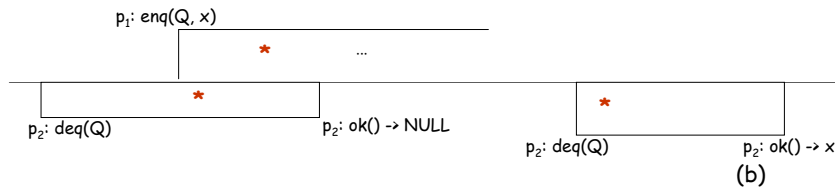
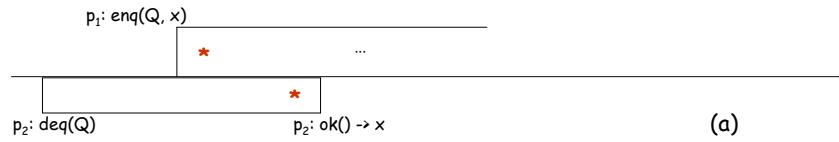


- Concurrent FIFO queues can be implemented by simpler concurrent objects, such as Test-And-Set and LL/SC registers.
- The same is true for other concurrent objects, such as stacks, lists, skip lists, graphs, etc.

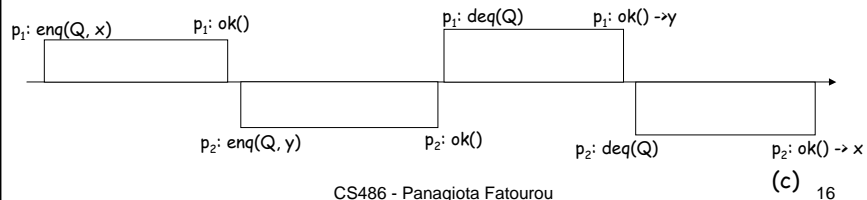
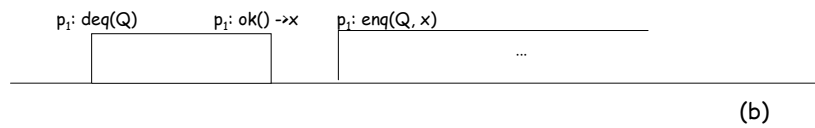
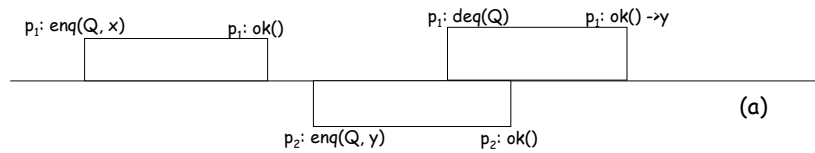
## Linearizable Executions



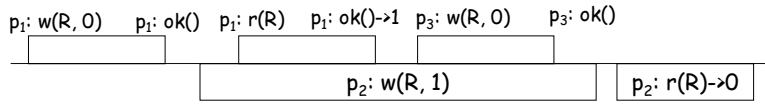
# Linearizable Executions



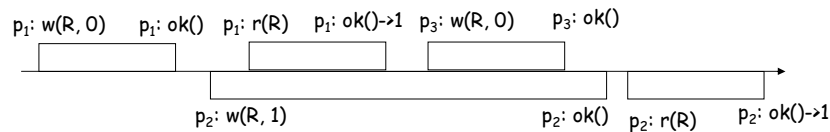
# Non-Linearizable Executions



## More Examples - Register Executions



(a) linearizable



(b) non-linearizable

## Linearizable implementations

- ❑ An implementation is **linearizable** if all the executions (histories) it produces are linearizable.
- ❑ In this part of the course, an implementation of a concurrent object is called **correct**, if it is linearizable.

## The Problem

- Can we implement snapshot objects in systems that provide only r/w registers?

CS486 - Panagiota Fatourou

19

## The trivial solution does not work!

- Assign a register  $R_i$  to each component  $A_i$ .
  - **UPDATE**( $i, v$ ): write( $R_i, v$ );
  - **SCAN**( $Q$ ): Read all registers and return a vector consisting of the values you read.
- Are there wait-free linearizable implementations of atomic snapshot objects?

**This algorithm is not linearizable!**

Processes	Register Values
$P_1$	$R_1$ $R_2$
read( $R_1$ )	0   0
write( $R_1, 1$ )	0   0
UPDATE( $A_1, 1$ )	0   0
UPDATE( $A_2, 2$ )	0   0
te( $R_2, 2$ )	0   0
read( $R_2$ )	0   2
SCAN() →	$\langle 0, 2 \rangle$

time

## A simple implementation using unbounded-size registers

- Each component,  $A_j$ ,  $1 \leq j \leq n$ , can be updated only by process  $p_j$  (single-writer snapshot).
- The implementation uses  $n$  registers  $R_1, \dots, R_n$ , one for each component. A register  $R_j$  has been assigned to each component  $A_j$ . UPDATE operations on  $A_j$  write only into  $R_j$ . Each register  $R_j$  is written only by  $p_j$  but it can be read by all processes (single-writer registers).
- Each register  $R_j$  is big enough to store the following information:
  - $val_j$ : the current value of  $A_j$
  - $tag_j$ : a timestamp used by  $p_j$  to differentiate the UPDATE operations it initiates
  - $view_j$ : a vector of  $n$  values, one for each component.
- This assumption is not realistic (since the registers are too big to be provided by the hardware).
- However, for the time being, we are only interested to design a simple such algorithm.

CS486 - Panagiota Fatourou

21

## The UnboundedSnapshot Algorithm

**UPDATE( $v$ ), code for process  $p_i$ :**

```
view := SCAN();
tag := increment  $p_i$ 's tag by 1;
write( $R_i$ ,  $\langle v, tag, view \rangle$ );
```

**SCAN, code for process  $p_i$ :**

repeatedly read  $R_1, \dots, R_n$  until:

1. Two **sets of reads** return the same  $R_j.tag$  for every  $j$ ; then, return the vector of values  $R_j.val$ ,  $1 \leq j \leq n$ , returned by the second set of reads, or
2. For some  $j$ , three distinct values of  $R_j.tag$  have been seen; return the vector of values in  $R_j.view$  associated with the last of the three values read in  $R_j.tag$ .

**Step Complexity =  $O(n^2)$  for SCAN and UPDATE.**  
**The implementation uses  $n$  SW registers of unbounded size.**

CS486 - Panagiota Fatourou

22

## Linearizability

**Definition:** A SCAN operation  $S$  returns a *consistent vector* if, for each component  $A_i$ ,  $S$  returns for  $A_i$  the value written by the last UPDATE on  $A_i$  for which the linearization point precedes the linearization point of  $S$  (or the initial value if such an UPDATE does not exist).

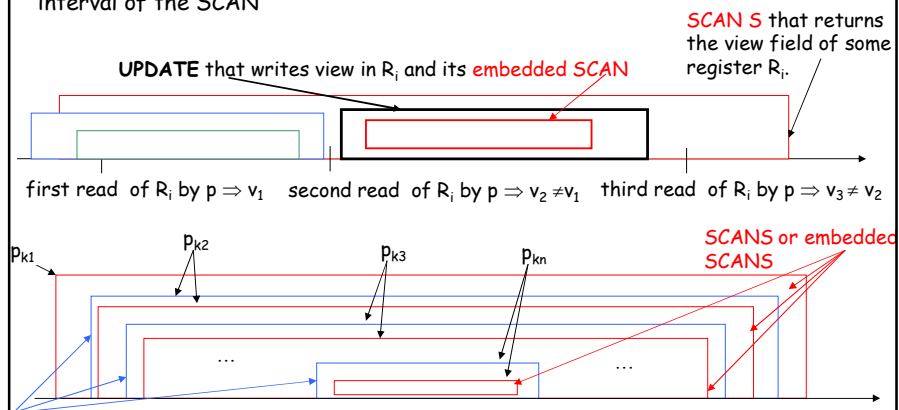
- To prove that the implementation is linearizable, we have to do the following:
  - assign linearization points to the operations of any execution
  - prove that the linearization point of each operation is within its execution interval, and
  - prove that each SCAN returns a consistent vector

CS486 - Panagiota Fatourou

23

## The UnboundedSnapshot Algorithm

➤ If process  $j$ , while performing repeated sets of reads on behalf of a SCAN, ever sees the same register  $R_i$  with three different tags, then it knows that some UPDATE $_i$  is completely contained within the execution interval of the SCAN



All  $p_{k1}, \dots, p_{kn}$  are simultaneously active. Since we have  $n$  processes in the system, the embedded SCAN executed by  $p_{kn}$  terminates by evaluating condition (1) of the SCAN code to TRUE.

## The UnboundedSnapshot Algorithm

### Assignment of Linearization Points

#### UPDATE Operations

- We insert the linearization point of an UPDATE at the point where its write occurs.

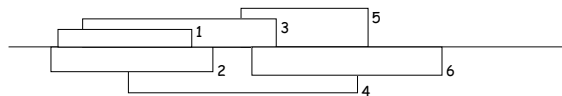
#### SCAN operations

- We assign linearization points not just to SCANS but also to embedded SCANS.
- We partition SCANS in two categories:
  - Direct SCANS: those that terminate by evaluating condition (1) to TRUE
  - Indirect SCAN: those that terminate by evaluating condition (2) to TRUE
- The linearization point of a direct SCAN can be placed anywhere within the end of the first of its last two sets of reads and the beginning of its second.

## The UnboundedSnapshot Algorithm

### Linearization points of indirect SCANS

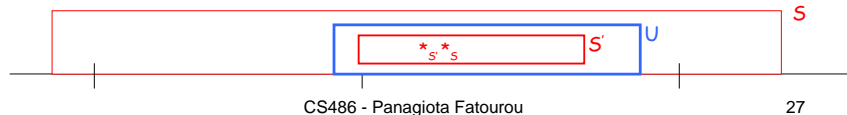
By induction on their response events (i.e., we linearize these SCANS one by one, in the order of their response events).



## The UnboundedSnapshot Algorithm

### Induction Base: Linearization point of indirect SCAN S which responds first.

- S returns a vector of values written by an UPDATE U (i.e., this vector has been calculated by the embedded SCAN S' of U).
- The execution interval of U and therefore also of S' (which is executed by U) is included in the execution interval of S.
- S' is a direct SCAN (since otherwise, the indirect SCAN that has the first response event would be S' and not S).
- Thus, a linearization point for S' has already been assigned.
- We linearize S at the same place as S'.



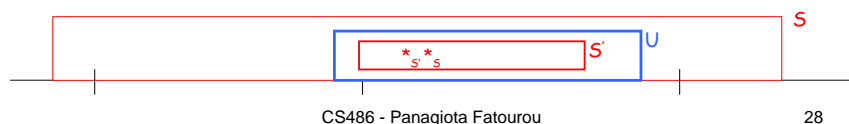
## The UnboundedSnapshot Algorithm

### Induction Hypothesis

- Let S be the indirect SCAN which has the k-th respond event. Assume that we have assigned linearization points to all SCAN operations for which it holds that their response events precede that of S.

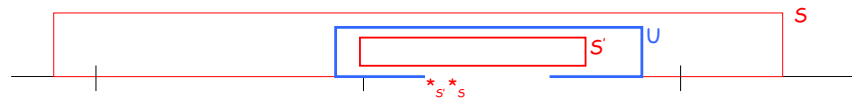
### Induction Step: We assign a linearization point to S.

- S returns a vector of values written by an UPDATE U (i.e., this vector has been calculated by the embedded SCAN S' of U).
- The execution interval of U and therefore also of S' (which is executed by U) is included in the execution interval of S.
- If S' is an indirect SCAN, by induction hypothesis, a linearization point has been assigned to S'. The same is TRUE for all direct SCANS.
- We linearize S at the same place as S'.



## The UnboundedSnapshot Algorithm

- **Lemma:** The linearization point of each SCAN or UPDATE is within its execution interval.
- **Proof:** For UPDATES and direct SCANS this is obvious.
- We prove the claim for indirect SCANS by induction on the order of their response events.
- **Induction Base**
  - Recall that the indirect SCAN  $S$  with the first response event is linearized at the same point as the direct SCAN  $S'$  from which it borrows its vector. Moreover, the execution interval of  $S'$  is included in the execution interval of  $S$ .
  - The linearization point of  $S'$  is in its execution interval  $\Rightarrow$  Thus, the linearization point of  $S$  is in its execution interval.

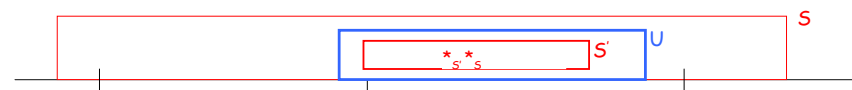


CS486 - Panagiota Fatourou

29

## The UnboundedSnapshot Algorithm

- **Induction Step:** We prove the claim for the indirect SCAN  $S$  with the  $k$ -th response event.
- $S$  returns a vector of values written by an UPDATE  $U$  (i.e., this vector has been calculated by the embedded SCAN  $S'$  of  $U$ ).
- The execution interval of  $U$  and therefore also of  $S'$  (which is executed by  $U$ ) is included in the execution interval of  $S$ .
- If  $S'$  is a direct SCAN, its linearization point is obviously within its execution interval. The same is TRUE, if  $S'$  is an indirect SCAN (by induction hypothesis since the execution interval of  $S'$  is included in the execution interval of  $S$ , and therefore the response event of  $S'$  precedes that of  $S$ ).
- The linearization point of  $S$  is placed at the same point as that of  $S' \Rightarrow$  thus, the linearization point of  $S$  is within its execution interval.



CS486 - Panagiota Fatourou

30

## The UnboundedSnapshot Algorithm

- **Lemma:** In every execution of UnboundedSnapshot, each SCAN operation returns a consistent vector of values.
- **Proof:** Obviously true for direct SCANS.
- For indirect SCANS, the claim will be proved by induction on the order of response events
- **Induction Base:** The indirect SCAN  $S$  which responds first, is linearized at the same point as the direct SCAN  $S'$  from which it borrows its vector. Since  $S'$  returns a consistent vector, the same holds for  $S$ .
- **Induction Step:** We prove the claim for the indirect SCAN  $S$  that has the  $k$ th response event.  $S$  returns the same vector of values as an embedded SCAN  $S'$  and is linearized at the same point as  $S'$ . Moreover, the execution interval of  $S'$  is within the execution interval of  $S$ .
- If  $S'$  is direct, it obviously returns a consistent vector. If  $S'$  is an indirect SCAN, by the induction hypothesis (since the response event of  $S'$  precedes that of  $S$ ), it follows that it returns a consistent vector. Thus,  $S$  returns a consistent vector.

A1	A2	A3	Vals	Scans
U(1)			[1,0,0]	
	U(1)		[1,1,0]	
				S1
U(3)			[3,1,0]	
	U(4)		[3,4,0]	
		U(5)	[3,4,5]	
		U(6)	[3,4,6]	
				S2
				S3
				S4

time ↓

\*<sub>U(1,1)</sub> [1,0,0] \*<sub>U(2,1)</sub> [1,1,0] \*<sub>S1</sub> \*<sub>U(1,3)</sub> [3,1,0] \*<sub>U(2,4)</sub> [3,4,0] \*<sub>U(3,5)</sub> [3,4,5] \*<sub>U(3,6)</sub> [3,4,6] \*<sub>S2</sub> \*<sub>S3</sub> \*<sub>S4</sub>

CS486 - Panagiota Fatourou

31

## Bibliography

These slides are based on material that appears in the following books:

- Herlihy and Wing, "Linearizability: a correctness condition for concurrent objects", *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 12(3): 463-492, 1990.
- N. Lynch, *Distributed Algorithms*, Morgan Kaufmann, 1996 (Chapter 13, Section 3).

CS486 - Panagiota Fatourou

32