



# **CONCURRENT SYSTEMS LECTURE 3**

Prof. Daniele Gorla



## MUTEX with specialized HW primitives

---

Atomic R/W registers provide quite a basic computational model.

We can strengthen the model by adding specialized HW primitives, that essentially perform in an atomic way the combination of some atomic instructions.

Usually, every operating system provides at least one specialized HW primitive.

The most common ones are:

- **Test&set**: atomic read+write of a boolean register
- **Swap**: atomic read+write of a general register
- **Fetch&add**: atomic read+increase of an integer register
- **Compare&swap**: atomic comparison+write of a general register; returns a boolean (the result of the comparison)
- ...





## MUTEX with Test&set

Let X be a boolean register; the **Test&set** primitive is implemented as follows:

```
X.test&set() :=  
    tmp ← X  
    X ← 1  
    return tmp
```

} *atomic (by hardware means)*

By using this primitive, MUTEX can be ensured by this simple protocol:

Initialize X at 0

lock() :=

wait X.test&set() = 0

return

unlock() :=

X ← 0

return





## MUTEX with Swap

Let  $X$  be a general register; the **Swap** primitive is implemented as follows:

```
X.swap(v) :=  
    tmp ← X  
    X ← v  
    return tmp
```

} *atomic (by hardware means)*

By using this primitive, the previous protocol for MUTEX can be adapted to the swap primitive by noting that

```
X.test&set() = X.swap(1)
```





## MUTEX with Compare&swap

Let X be a boolean register; the **Compare&swap** primitive is implemented as follows:

```
X.compare&swap(old, new) :=  
    if X = old then X ← new  
        return true  
    return false
```

} *atomic*

By using this primitive, MUTEX can be obtained as follows:

Initialize X at 0

```
lock() :=  
    wait X.compare&swap(0,1)=true  
    return
```

```
unlock() :=  
    X ← 0  
    return
```





## MUTEX with Fetch&add

Up to now, all solutions enjoy deadlock freedom, but allow for starvation

→ use Round Robin to promote the liveness property

Let  $X$  be an integer register; the **Fetch&add** primitive is implemented as follows:

```
X.fetch&add(v) :=  
    tmp ← X  
    X ← X+v  
    return tmp
```

} *atomic*

By using this primitive, MUTEX can be obtained as follows:

Initialize TICKET and NEXT at 0

```
lock() :=  
    my_tick ← TICKET.fetch&add(1)  
    wait my_tick = NEXT  
    return
```

```
unlock() :=  
    NEXT ← NEXT+1  
    return
```

*Bounded bypass  
with bound  $n-1$*





## Safe Registers

---

Atomic R/W and specialized HW primitives provide some form of atomicity

→ is it possible to enforce MUTEX without atomicity?

A **MRSW Safe register** is a register that provides READ and WRITE such that:

1. Every READ that does not overlap with a WRITE returns the value stored in the register
2. A READ that overlaps with a WRITE returns any value (of the register domain)

A **MRMW Safe register** behaves like a MRSW safe register, when WRITE operations do not overlap; otherwise, in case of overlapping WRITES, the register can contain any value (of the register domain)

This is the weakest type of register that is useful in concurrency





## Bakery algorithm (Lamport 1974)

Idea:

- Every process gets a ticket
- Because we don't have atomicity, tickets may be not unique
- Tickets can be made unique by pairing them with the process ID
- The smallest ticket (seen as a pair) grants the access to the CS

Initialize FLAG[i] to down and MY\_TURN[i] to 0, for all i

lock(i) :=

FLAG[i] ← up

MY\_TURN[i] ← max{MY\_TURN[1], ..., MY\_TURN[n]}+1

FLAG[i] ← down

forall j ≠ i

wait FLAG[j] = down

wait (MY\_TURN[j] = 0 OR

⟨MY\_TURN[i], i⟩ < ⟨MY\_TURN[j], j⟩)

} doorway

} bakery

(including CS)

unlock(i) :=

MY\_TURN[i] ← 0





---

**Lemma 1:** Let  $p_i$  enter the bakery before  $p_j$  enters the doorway; then,

$$\text{MY\_TURN}[i] < \text{MY\_TURN}[j].$$

*Proof:*

- Let  $t$  be the value of  $\text{MY\_TURN}[i]$  after that  $p_i$  exits the doorway
- When  $p_j$  computes its ticket, it reads  $t$  from  $\text{MY\_TURN}[i]$  (there is no write overlapping with this read)
- Hence,  $\text{MY\_TURN}[j]$  is at least  $t+1$





**Lemma 2:** Let  $p_i$  be in the CS and  $p_j$  is in the doorway or in the bakery; then,

$$\langle \text{MY\_TURN}[i], i \rangle < \langle \text{MY\_TURN}[j], j \rangle.$$

*Proof:*

If  $p_i$  is in the CS, it has terminated its first wait for  $j$

→ let's consider the read of  $\text{FLAG}[j]$  done by  $p_i$  that terminates such wait

W.r.t. the execution of  $p_j$ , it can be that

- This read overlaps with  $\text{FLAG}[j] \leftarrow \text{up}$ : by Lemma1,  $\text{MY\_TURN}[i] < \text{MY\_TURN}[j]$  and  $\checkmark$
- This read is contained within the computation of  $\text{MY\_TURN}[j]$ 
  - this is not possible, since  $\text{MY\_TURN}$  is computed with the  $\text{FLAG}$  up
- This read overlaps with  $\text{FLAG}[j] \leftarrow \text{down}$  or this read happens when  $p_j$  is in the bakery:
  - $\text{MY\_TURN}[j]$  has been decided and no write will change it until  $p_j$  is in the bakery
  - $\text{MY\_TURN}[j] > 0$  (it has been obtained by summing 1 to some natural number)
  - When  $p_i$  has evaluated the second wait for  $j$ , it found  $\langle \text{MY\_TURN}[i], i \rangle < \langle \text{MY\_TURN}[j], j \rangle$  and  $\checkmark$





**MUTEX:**  $p_i$  and  $p_j$  cannot simultaneously be in the C.S.

*Proof:* By contradiction, by Lemma2 applied twice, we would have

$$\langle \text{MY\_TURN}[i], i \rangle < \langle \text{MY\_TURN}[j], j \rangle \text{ and } \langle \text{MY\_TURN}[j], j \rangle < \langle \text{MY\_TURN}[i], i \rangle$$

**Deadlock freedom:** by contradiction, assume that there is a lock but nobody enters its CS

- All processes in the bakery (call this set  $Q$ ) are blocked in their wait
- The first wait cannot block for ever
  - All  $p_i \in Q$  have their FLAG down
  - All  $p_i \notin Q$  have their FLAG down (if they're not in the doorway) or will eventually put their FLAG down (they cannot remain in the doorway for ever)
- The second wait cannot block all of them for ever
  - Tickets can be totally ordered (lexicographically)
  - Let  $\langle \text{MY\_TURN}[i], i \rangle$  be the minimum
  - The second wait evaluated by  $p_i$  eventually succeeds for all  $j$ 
    - If  $p_j$  is before the doorway  $\rightarrow \text{MY\_TURN}[j] = 0$
    - If  $p_j$  is in the doorway  $\rightarrow \text{MY\_TURN}[i] < \text{MY\_TURN}[j]$  (bec.of Lemma1)
    - If  $p_j$  is in the bakery, by assumption  $\langle \text{MY\_TURN}[i], i \rangle < \langle \text{MY\_TURN}[j], j \rangle$  since it is the minimum



---

## Bounded bypass (with bound $n-1$ ):

Let  $p_i$  and  $p_j$  competing for the CS and  $p_j$  wins

Then,  $p_j$  enters its CS, completes it, unlocks and then invokes lock again

- If  $p_i$  has entered the CS  $\rightarrow \checkmark$
- Otherwise, by Lemma 1,  $MY\_TURN[i] < MY\_TURN[j]$   
 $\rightarrow p_j$  cannot bypass  $p_i$  again!
- At worse,  $p_i$  has to wait all other processes before entering its CS  
(indeed, since there is no deadlock, when  $p_i$  is waiting somebody enters the CS)





## Aravind's algorithm (2011)

Problem with Lamport's alg.: registers must be unbounded (every invocation of lock potentially increases the counter by 1  $\rightarrow$  domain of the registers is all naturals!)

For all processes, we have a FLAG and a STAGE (both binary MRSW), and a DATE (a MRMW register that ranges from 1 to  $2n$ )

For all  $i$ , initialize

- FLAG[ $i$ ] to down
- STAGE[ $i$ ] to 0
- DATE[ $i$ ] to  $i$

```
lock(i) :=
  FLAG[i]  $\leftarrow$  up
  repeat
    STAGE[i]  $\leftarrow$  0
    wait ( $\forall j \neq i$ . FLAG[j] = down OR
          DATE[i] < DATE[j])
    STAGE[i]  $\leftarrow$  1
  until  $\forall j \neq i$ . STAGE[j] = 0
```

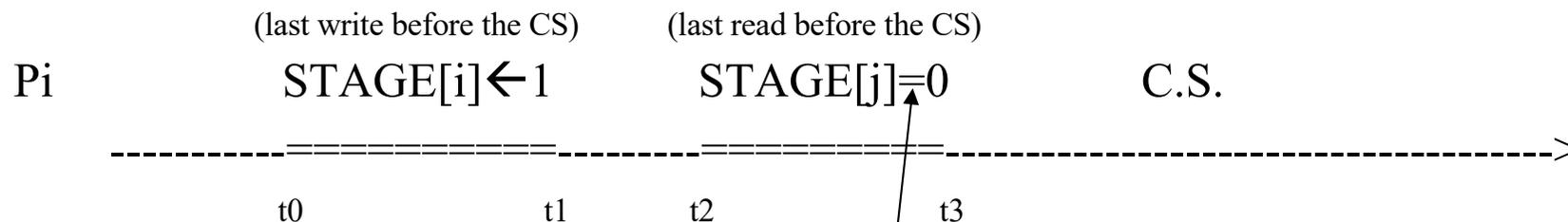
```
unlock(i) :=
  tmp  $\leftarrow$   $\max_j$ {DATE[j]}+1
  if tmp  $\geq$  2n
    then  $\forall j$ .DATE[j]  $\leftarrow$  j
    else DATE[i]  $\leftarrow$  tmp
  STAGE[i]  $\leftarrow$  0
  FLAG[i]  $\leftarrow$  down
```





**Thm.:** if  $p_i$  is in the CS, then  $p_j$  cannot simultaneously be in its CS.

*Proof:* By contradiction. Let us consider the execution of  $p_i$  leading to its CS:



Let's consider the last write of  $p_j$  before its CS (i.e.,  $STAGE[j] \leftarrow 1$ )

- It cannot complete before  $t_2$ , because of
- So, it must overlap with  $[t_2, t_3]$  or happen after  $t_3$ 
  - but then the last read of  $p_j$  from  $STAGE[i]$  happens after  $t_1$
  - $p_j$  finds  $STAGE[i] = 1$  and cannot enter its CS

**Cor.:** DATE is never written concurrently.





**Lemma 1:** exactly every  $n$  CSs there is a reset of DATE.

*Proof:* Because of the previous corollary

- The first CS leads  $\max_j \{DATE[j]\}$  to  $n+1$
- The second CS leads  $\max_j \{DATE[j]\}$  to  $n+2$
- ...
- The  $n$ -th CS leads  $\max_j \{DATE[j]\}$  to  $n+n = 2n \rightarrow$  RESET

**Lemma 2:** there can be at most one reset of DATE during an invocation of lock.

*Proof:*

Let  $p_i$  invoke lock. If no reset occurs  $\rightarrow \checkmark$

Otherwise, let us consider the moment in which a reset occurs.

If  $p_i$  is the next process that enters the CS  $\rightarrow \checkmark$

Otherwise, let  $p_j$  be the process that enters; its next date is  $n+1 > DATE[i]$

$\rightarrow p_j$  cannot surpass  $p_i$  again (before a RESET)

The worst case is when all proc's perform lock together and  $i=n$

$\rightarrow$  all  $p_1, \dots, p_{n-1}$  surpass  $p_n$

$\rightarrow$  then  $p_n$  enters and it resets the DATE in its unlock

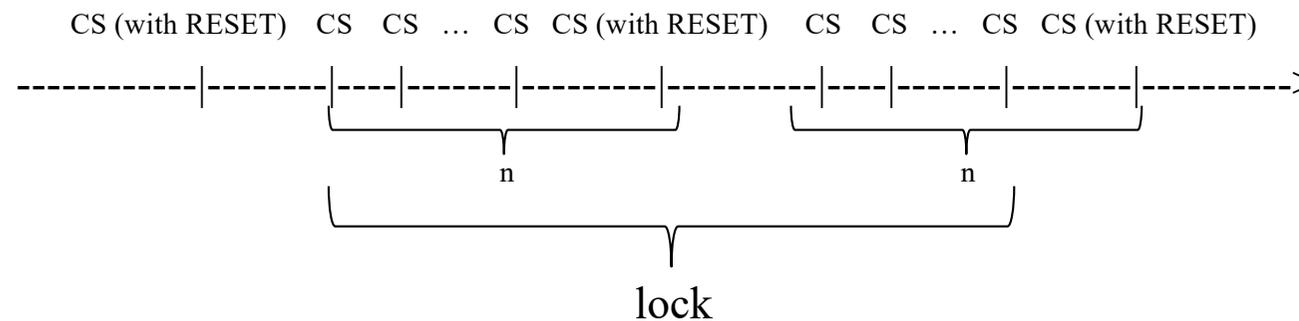




**Thm.:** The algorithm satisfies bounded bypass with bound  $2n-2$ .

*Proof:*

By Lemma1:



By Lemma2:

i.e., there is an upper bound of  $2n-1$

→ this bound is not reachable, whereas the bound reachable is  $2n-2$ :

- $p_n$  invokes lock alone, completes its CS (the first after the reset) and its new DATE is  $n+1$
- Then all processes invoke lock simultaneously
- $p_n$  has to wait all other processes to complete their CSs (after that  $p_i$  completes its CS it has  $DATE[i] \leftarrow n+i+1$ )
- When  $p_{n-1}$  completes its CS, its new DATE will be  $n+(n-1)+1 = 2n \rightarrow$  RESET
- Now all  $p_1, \dots, p_{n-1}$  invoke lock again and complete their CSs (after that  $p_i$  completes its CS, now it has  $DATE[i] \leftarrow n+i$ )
- So,  $p_n$  has to wait  $n-1$  CSs for the reset, and another  $n-1$  CSs before entering again





## Improvement of Aravind's algorithm

Consider the following revision of Aravind's UNLOCK:

```
unlock(i) :=  
   $\forall j \neq i. \text{if } \text{DATE}[j] > \text{DATE}[i] \text{ then } \text{DATE}[j] \leftarrow \text{DATE}[j] - 1$   
   $\text{DATE}[i] \leftarrow n$   
   $\text{STAGE}[i] \leftarrow 0$   
   $\text{FLAG}[i] \leftarrow \text{down}$ 
```

Since the LOCK is like before, the revised protocol satisfies MUTEX.

Furthermore, you can prove that it satisfies bounded bypass with bound  $n-1$

→ EXERCISE!

