



# CS G623: Advanced Operating Systems

## Lecture 10

**BITS Pilani**

Pilani Campus

Amit Dua

August 25, 2018

# Topics to be discussed

---



- DME
- Goals
- Approaches
  - Centralized
  - Distributed
  - Token based
- Leader-election algorithm
- Centralized approach
- Lamport's algorithm
- Richart-Agarwala's algorithm
- Maekwa's algorithm

# Distributed Mutual Exclusion

## What is mutual exclusion?

It is exclusive access to a shared resource or to the critical region.

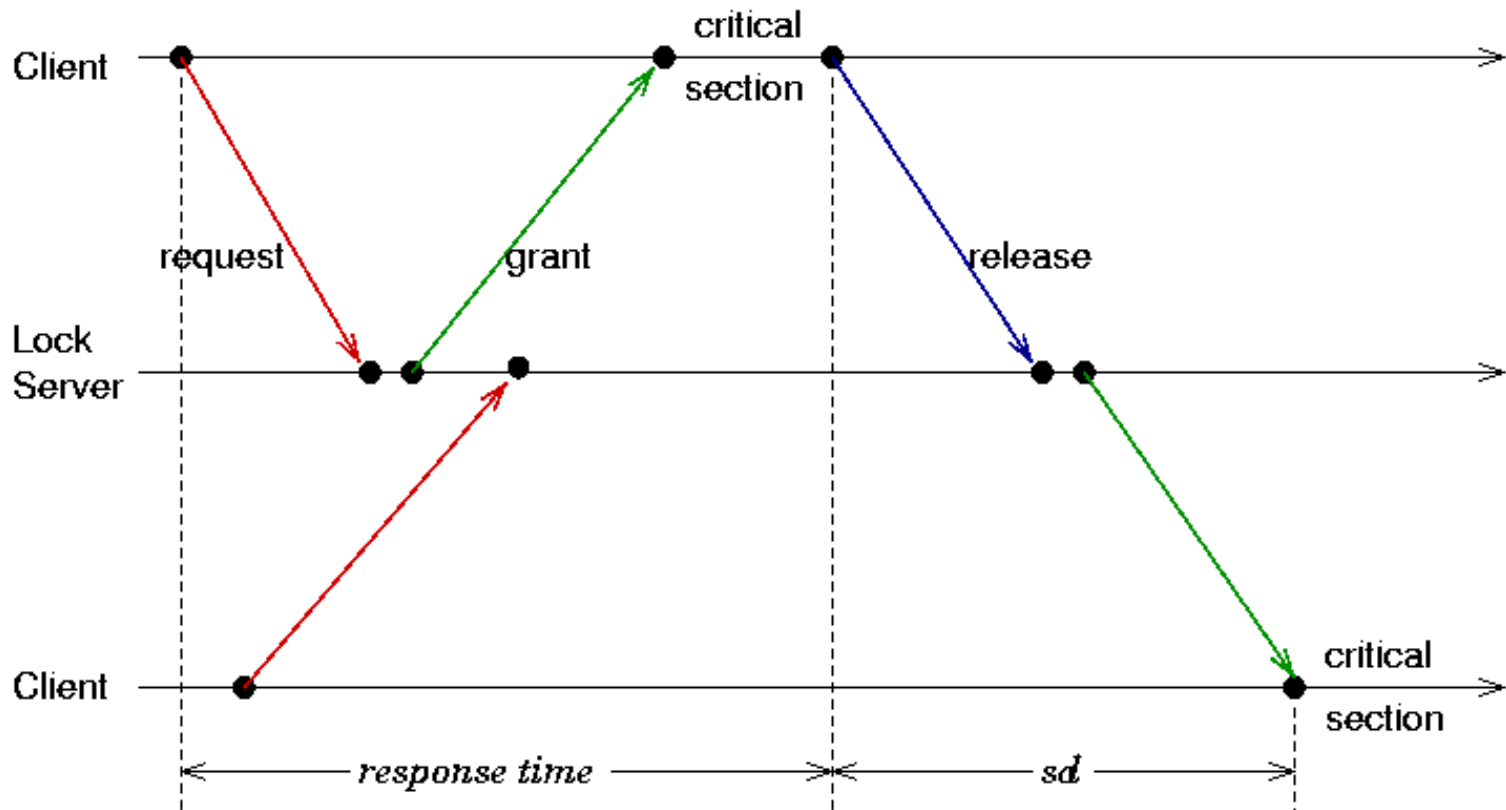
# DME Continued.....

- An algorithm used for implementing mutual exclusion must satisfy:
  - Mutual exclusion
  - No starvation,
  - Freedom from deadlock
  - Fault tolerant
- To handle mutual exclusion in a distributed system,
  - Centralized Approach
  - Distributed Approach
  - Token-Passing Approach
- All use message passing approach rather than shared variable

# Performance of DME Algorithms

- Performance of each algorithm is measured in terms of
  - no. of messages required for CS invocation
  - synchronization delay (leaving & entering)
  - response time (arrival, sending out, Entry & Exit)
- system throughput =  $1/(sd+E)$

# Performance Continued...



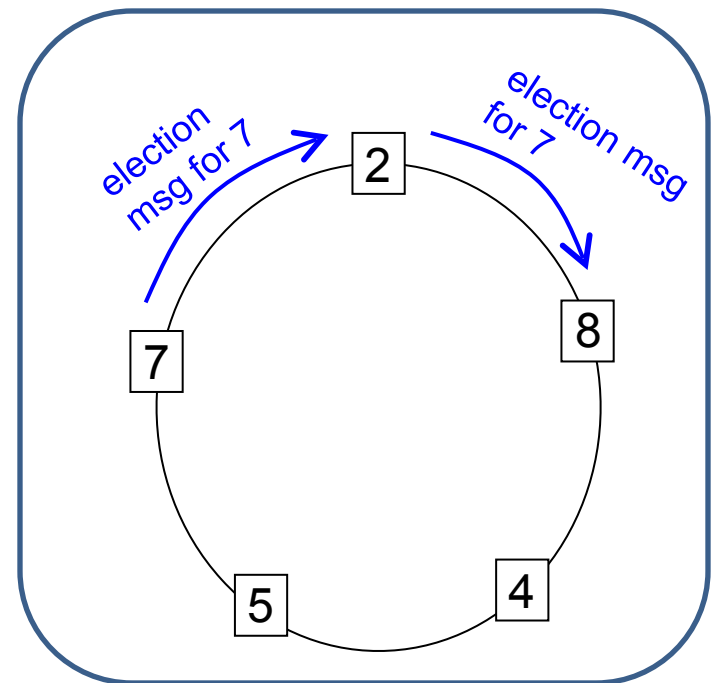
$$\text{System Throughputs} = \frac{1}{Sd + E}$$

# DME: A Centralized Algorithm

- One process is elected as the **coordinator/leader** (e.g., the one running on the machine with the highest network address) which resolves conflicts.

# An example leader election algorithm

- Each node has a **unique identifier**
- Nodes only send messages clockwise
- Each node acts on its own
- Protocol:
  - A node send election message with its own id clockwise
  - Election message is forwarded if id in message larger than own message
  - Otherwise message discarded
- A node becomes leader if it sees its own election message



Lelang-Chang Algorithm



# Central Approach: Adv.s & Disadv.s

- Mutual exclusion can be achieved.
  - No process waits forever.
  - Easy to implement.
  - It can be used for general resource allocation rather than just managing mutual exclusion.
- 
- Single point of failure
  - If process normally blocks after making a request, difficult to distinguish a dead coordinator from “access denied”.

# Lamport's DME

## Requesting the critical section.

1. When a site  $S_i$  wants to enter into the CS, it sends a **REQUEST**( $T=ts_i, i$ ) message to all the sites in its request set  $R_i$  and places the request on `request_queuei`.
2. When a site  $S_j$  receives the **REQUEST**( $ts_i, i$ ) message from site  $S_i$ , it returns a **timestamped REPLY** message to  $S_i$  and places site  $S_i$ 's request on `request_queuej`.

## Executing the critical section.

Site  $S_i$  enters the CS when the **two following conditions** hold:

1.  $S_i$  has received a message with timestamp larger than ( $ts_i, i$ ) from all other sites.
2.  $S_i$ 's request is at the top of `request_queuei`.

# Continued...

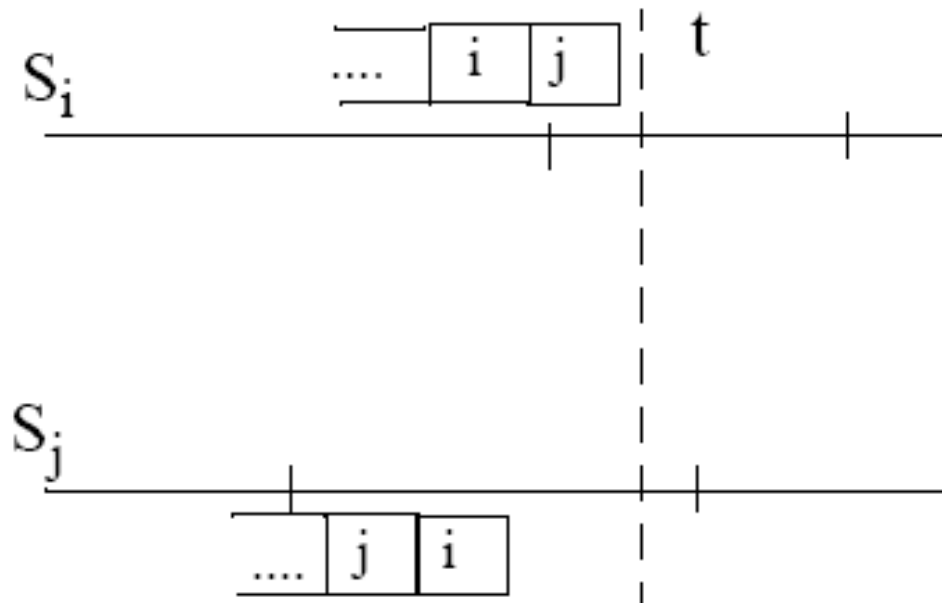
## Releasing the critical section.

1. Site  $S_i$ , upon exiting the CS, removes its' request from the top of its' request queue and sends a timestamped **RELEASE message** to all the sites in its request set.
2. When a site  $S_j$  receives a RELEASE message from site  $S_i$ , it removes  $S_i$ 's request from its request queue.
3. When a site removes a request from its request queue, its' own request may become the top of the queue, enabling it to enter into the CS. The algorithm executes CS requests in the increasing order of timestamps.

# Example Lamport's DME

# Correctness

- Suppose that both  $S_i$  and  $S_j$  were in CS at the same time ( $t$ ).
- Then we have:



No. of messages,  
Synch delay,  
Optimization if  
any?)

# Ricart-Agrawala DME

## Requesting the CS

1. A site  $S_i$  wanting to enter into CS sends a timestamped request message to all sites in its request set.
2. Upon reception of a request message from  $S_i$ , site  $S_j$  immediately sends a timestamped *reply* message if and only if:
  1.  $S_j$  is not currently interested in the critical section OR
  2.  $S_j$ 's request has a lower priority (*usually this means having a later timestamp*)
3. Otherwise,  $S_j$  will defer the reply message. This means that a reply will be sent only after  $S_j$  has finished using the critical section itself.

## Executing in the Critical Section

1.  $S_i$  enters into CS after receiving reply from all other sites

## Releasing the critical section

4. Upon exiting the critical section, the site  $S_i$  sends reply messages to all deferred requests.

# Maekawa's DME Algorithm

- A site requests permission only from a subset of sites.
- Request set of sites  $S_i$  &  $S_j$ :  $R_i, R_j$  such that  $R_i$  and  $R_j$  will have at least one common site ( $S_k$ ).  $S_k$  mediates conflicts between  $R_i$  and  $R_j$ .
- A site can send only one REPLY message at a time, i.e., a site can send a REPLY message only after receiving a RELEASE message for the previous REPLY message.

# Request Subsets

- Rules for generating request sets:
  - Sets  $R_i$  and  $R_j$  have at least one common site.
  - $S_i$  is always in  $R_i$ .
  - Cardinality of  $R_i$ , i.e., the number of sites in  $R_i$  is  $K$ .
  - Any site  $S_i$  is in  $K$  number of  $R_i$ 's.  $N = K(K - 1) + 1$
- What would be the request set for  $N = 3$ ?



# Request Subsets continued...

- What would be the request set for  $N = 7$ ?

- Example: Finite Projective Planes

$$S(0) = \{0,5,6\} \checkmark$$

$$S(1) = \{1,3,6\}$$

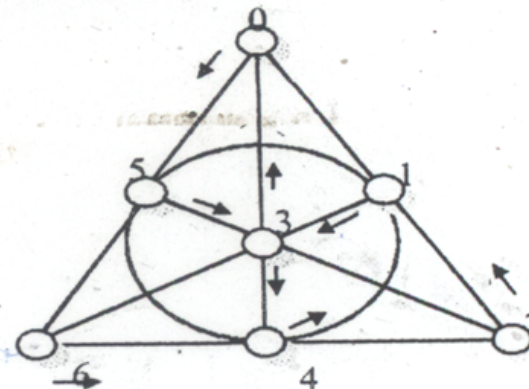
$$S(2) = \{2,1,0\} \checkmark$$

$$S(3) = \{3,0,4\} \checkmark$$

$$S(4) = \{4,1,5\}$$

$$S(5) = \{5,3,2\}$$

$$S(6) = \{6,4,2\}$$



# Request Set with N=13

$$R_1 = \{ 1, 2, 3, 4 \}$$

$$R_2 = \{ 2, 5, 8, 11 \}$$

$$R_3 = \{ 3, 6, 8, 13 \}$$

$$R_4 = \{ 4, 6, 10, 11 \}$$

$$R_5 = \{ 1, 5, 6, 7 \}$$

$$R_6 = \{ 2, 6, 9, 12 \}$$

$$R_7 = \{ 2, 7, 10, 13 \}$$

$$R_8 = \{ 1, 8, 9, 10 \}$$

$$R_9 = \{ 3, 7, 9, 11 \}$$

$$R_{10} = \{ 3, 5, 10, 12 \}$$

$$R_{11} = \{ 1, 11, 12, 13 \}$$

$$R_{12} = \{ 4, 7, 8, 12 \}$$

# Maekawa's DME Algorithm

## Requesting the critical section

1. A site  $S_i$  requests access to the CS by sending REQUEST( $i$ ) messages to all the sites in its request set  $R_i$ .
2. When a site  $S_j$  receives the REQUEST( $i$ ) message, it sends a
  - REPLY( $j$ ) message to  $S_i$  provided it hasn't sent a REPLY message to a site from the time it received the last RELEASE message.
  - Otherwise, it queues up the REQUEST for later consideration.

## Executing the critical section

1. Site  $S_i$  accesses the CS only after receiving REPLY messages from all the sites in  $R_i$ .

## Releasing the critical section

1. After the execution of the CS is over, site  $S_i$  sends RELEASE( $i$ ) message to all the sites in  $R_i$ .
2. When a site  $S_j$  receives a RELEASE( $i$ ) message from site  $S_i$ , it sends a REPLY message to the next site waiting in the queue and deletes that entry from the queue. If the queue is empty, then the site updates its state to reflect that the site has not sent out any REPLY message.

# Maekawa's Example

# Maekawa's DME Correctness

- **Proof of correctness (by contradiction)**

Suppose  $S_i$  and  $S_j$  are in the CS at the same time  $R_i \wedge R_j = \{ S_k \}$ .

Then,  $S_k$  must have sent REPLY to both  $S_i$  and  $S_j$ , which is not allowed.

- **Performance**

messages/CS = ?

synchronization delay = ?

Do you foresee any problem?

# Solution to Deadlock

- **FAILED**

A FAILED message from site  $S_i$  to site  $S_j$  indicates that  $S_i$  cannot grant  $S_j$ 's request.

- **INQUIRE**

An INQUIRE message from  $S_i$  to  $S_j$  indicates that  $S_i$  would like to find out from  $S_j$ ...

- **YIELD**

A YIELD message from site  $S_i$  to  $S_j$  indicates that  $S_i$  is returning the permission to  $S_j$ .

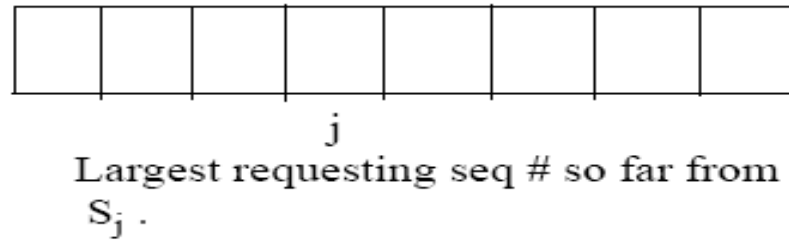
# Token-based DME Algorithms

- A site enters CS if it possesses the **token** (only one token for the System).
- The major difference is the way the token is **searched**.
- Use **sequence numbers** instead of timestamps
  - Used to distinguish requests from same site
  - Keep advancing independently at each site
- The proof of mutual exclusion is **trivial**

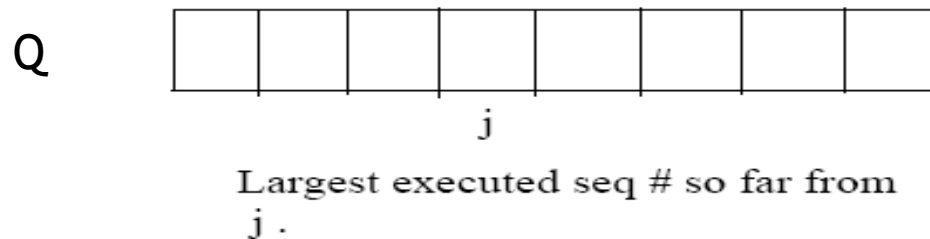


# Suzuki-Kasami's Broadcast Algorithm

A site has RN:



Token: A queue & LN



- $LN[j]$  is the sequence no. of the request that site  $S_j$  executed most recently.

# The Algorithm

- **Requesting the critical section.**

1. If the requesting site  $S_i$  does not have the token, then it **increments** its sequence number,  $RN_i[i]$ , and sends a **REQUEST( $i, sn$ )** message to all other sites. ( $sn$  is the updated value of  $RN_i[i]$ .)
2. When a site  $S_j$  receives this message, it sets  **$RN_j[i]$  to  $\max(RN_j[i], sn)$** . If  $S_j$  has the idle token, it sends the token to  $S_i$  if  **$RN_j[i] = LN[j] + 1$** .

- **Executing the critical section.**

3. Site  $S_i$  executes the CS when it has received the token.

- **Releasing the critical section.**

Having finished the execution of the CS, site  $S_i$  takes the following actions:

4. It sets  $LN[i]$  element of the token array equal to  $RN_i[i]$ .
5. For every site  $S_j$  whose ID is not in the token queue, it appends its ID to the token queue if  **$RN_i[j] = LN[j] + 1$** .
6. If token queue is nonempty after the above update, then it deletes the top site ID from the queue and sends the token to the site indicated by the ID.

# Example1

# Example2

# Analysis

- **Correctness**

Mutex is trivial.

- Theorem:

A requesting site enters the CS in finite time.

- Proof:

A request enters the token queue in finite time. The queue is in FIFO order, and there can be a maximum  $N-1$  sites ahead of the request.

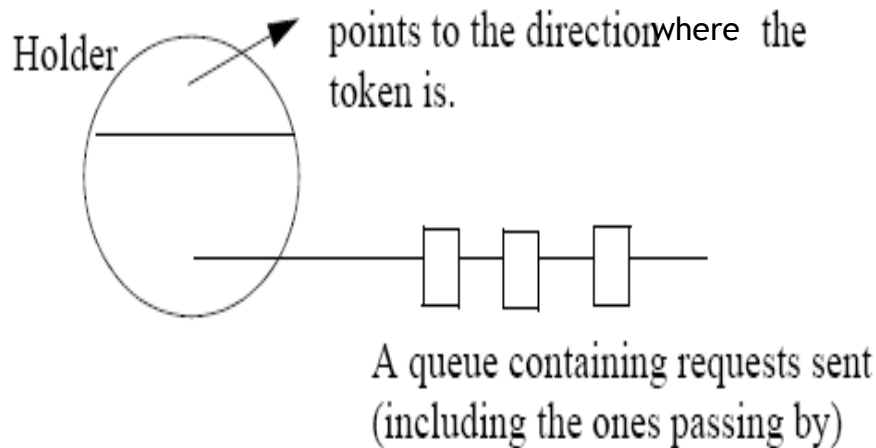
- **Performance**

0 or  $N$  messages per CS invocation.

Synchronous delay: 0 or  $T$  (only 1 message)

# Raymond's Tree-based DME Algorithm

A node:



# Raymonds' DME Continued...

Let us see a tree...

# The Algorithm

- **Requesting the critical section.**

1. When a site wants to enter the CS, it sends a REQUEST message to the node along the directed path to the root, provided it does not hold the token and its **request\_q** is empty. It then adds to its request\_q.

2. When a site on the path receives this message, it places the REQUEST in its **request\_q** and sends a REQUEST message along the directed path to the root provided it has **not** sent out a REQUEST message on its outgoing edge.

3. When the root site receives a REQUEST message, it sends the token to the site from which it received the REQUEST message and **sets its holder variable** to point at that site.

4. When site receives the token, it **deletes** the top entry from its **request\_q**, sends the token to the site indicated in this entry, and sets its *holder* variable to point at that site. If the **request\_q** is **nonempty** at this point, then the site sends a REQUEST message to the site which is pointed at by *holder* variable.



# The Algorithm Continued....

- **Executing the critical section**

5. A site enters the CS when it receives the token and its **own entry** is at the top of its **request\_q**. In this case, the site deletes the top entry from its request\_q and enters the CS.

- **Releasing the critical section**

6. If its **request\_q** is nonempty, then it deletes the top entry from its request\_q, sends the token to that site, and sets its *holder* variable to point at that site.

7. If the **request\_q** is nonempty at this point, then the site sends a REQUEST message to the site which is pointed at by the *holder* variable.

# Example 1

# Example2

# Analysis

- **Proof of Correctness**

Mutex is trivial.

Finite waiting: All the requests in the system form a FIFO queue and the token is passed in that order.

- **Performance**

$O(\log N)$  messages per CS invocation.

Sync. delay:  $(T \log N) / 2$

The average distance between two sites is  $\log N / 2$ .

# Singhal's Heuristic Algorithm

- Data Structures:
  - $S_i$  maintains  $SV_i[1..M]$  and  $SN_i[1..M]$  for storing information on other sites: **state** and **highest sequence number**.
  - Token contains 2 arrays:  $TSV[1..M]$  and  $TSN[1..M]$ .
  - States of a site
    - R : requesting CS
    - E : executing CS
    - H : Holding token, idle
    - N : None of the above
  - Initialization:
    - $SV_i[j] := N$ , for  $j = M .. i$ ;  $SV_i[j] := R$ , for  $j = i-1 .. 1$ ;  
 $SN_i[j] := 0$ ,  $j = 1..M$
    - Site  $S_1$  is in state H.
    - Token:  $TSV[j] := N$  &  $TSN[j] := 0$ ,  $j = 1 .. M$ .

# Algorithm Continued...

- **Requesting CS**

- If  $S_i$  has no token and requests CS:
  - $SV_i[i] := R$ .  $SN_i[i] := SN_i[i] + 1$ .
  - Send REQUEST( $i, sn$ ) to sites  $S_j$  for which  $SV_i[j] = R$ . (sn: sequence number, updated value of  $SN_i[i]$ ).
- Receiving REQUEST( $i, sn$ ): if  $sn \leq SN_j[i]$ , ignore. Otherwise, update  $SN_j[i]$  and do:
  - $SV_j[j] = N \rightarrow SV_j[i] := R$ .
  - $SV_j[j] = R \rightarrow$  If  $SV_j[i] \neq R$ , set it to  $R$  & send REQUEST( $j, SN_j[j]$ ) to  $S_i$ . Else do nothing.
  - $SV_j[j] = E \rightarrow SV_j[i] := R$ .
  - $SV_j[j] = H \rightarrow SV_j[i] := R$ ,  $TSV[i] := R$ ,  $TSN[i] := sn$ ,  $SV_j[j] = N$ . Send token to  $S_i$ .

- **Executing CS**

after getting token. Set  $SV_i[i] := E$ .

# Algorithm Continued...

- **Releasing CS**
  - $SVi[i] := N$ ,  $TSV[i] := N$ . Then, do:
    - For other  $Sj$ : if  $(SNI[j] > TSN[j])$ , then  
 $\{TSV[j] := SVi[j]; TSN[j] := SNI[j]\}$  //update token info from local info
    - else  $\{SVi[j] := TSV[j]; SNI[j] := TSN[j]\}$  // otherwise
  - If  $SVi[j] = N$ , for all  $j$ , then set  $SVi[i] := H$ . Else send token to a site  $Sj$  provided  $SVi[j] = R$ .
- Fairness of algorithm will depend on choice of  $Si$ , since no queue is maintained in token.
- Arbitration rules to ensure fairness used.
- Performance
  - Low to moderate loads: average of  $N/2$  messages.
  - High loads:  $N$  messages (all sites request CS).
  - Synchronization delay:  $T$ .

