

Hệ thống phân tán



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN

Các thuật toán đảm bảo truy cập tuần tự

- Không sử dụng token:
 - A site/process có thể vào miền găng khi một điều kiện nào đó thỏa.
 - Thuật toán đảm bảo điều kiện chỉ thỏa cho chỉ một site/process.
- Sử dụng token:
 - Một token duy nhất đc chia sẻ giữa các sites/processes.
 - Site giữ token đc phép vào miền găng.
 - Phải xử lý tình huống mất token, site giữ token bị tắt, khả năng nhiều token, v.v...

Mô hình hệ thống

- ☐ Giả thiết có nhiều site muốn vào miền găng, các yêu cầu đc đưa vào hàng đợi.
- ☐ Site States: Requesting CS, executing CS, idle.
- ☐ Requesting CS: blocked until granted access, cannot make additional requests for CS.
- ☐ Executing CS: using the CS.
- ☐ Idle: action is outside the site. In token-based approaches, *idle* site can have the token.

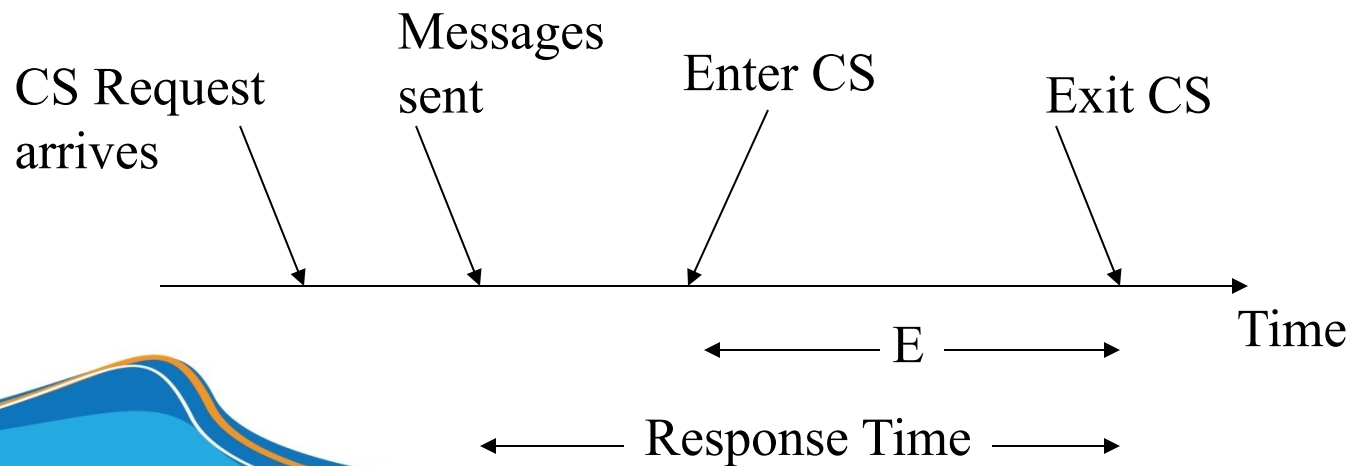
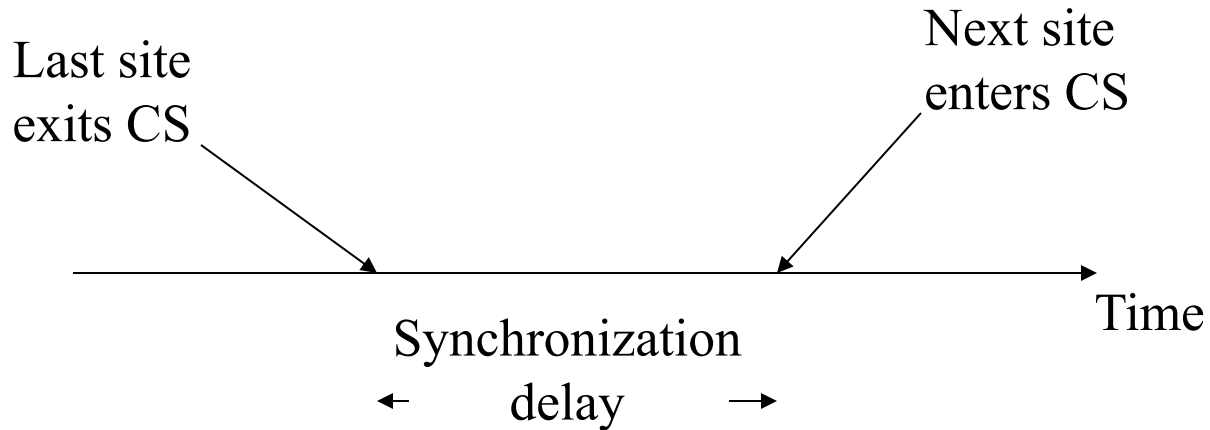
Truy cập luân phiên: Yêu cầu

- ☐ Không bị deadlock: two or more sites should not endlessly wait on conditions/messages that never become true/arrive.
- ☐ Không bị starvation: No indefinite waiting.
- ☐ Công bằng: Order of execution of CS follows the order of the requests for CS. (equal priority).
- ☐ Chịu lỗi: recognize “faults”, reorganize, continue. (e.g., loss of token).

Hiệu suất

- ☐ Số thông điệp trao đổi để vào miền găng: should be minimized.
- ☐ Synchronization delay,: tgian từ khi một site rời miền găng và site kế tiếp vào miền găng.
- ☐ Response time: khoảng tgian từ khi gửi yêu cầu vào miền găng đến khi rời khỏi miền găng.
- ☐ System throughput, : số lượng site thực hiện xong miền găng trong một đơn vị tgian.
- ☐ If sd is synchronization delay, E trung bình tgian thực thi trong miền găng: system throughput = $1 / (sd + E)$.

Hiệu suất



Performance ...

- Low and High Load:
 - ▣ Low load: No more than one request at a given point in time.
 - ▣ High load: Always a pending mutual exclusion request at a site.
- Best and Worst Case:
 - ▣ Best Case (low loads): Round-trip message delay + Execution time. $2T + E$.
 - ▣ Worst case (high loads).
- Message traffic: low at low loads, high at high loads.
- Average performance: when load conditions fluctuate widely.

Giải pháp đơn giản

- ☐ Control site: cấp quyền vào miền găng.
- ☐ A site sends REQUEST message to control site.
- ☐ Controller cấp quyền từng site một.
- ☐ Synchronization delay: $2T$ \rightarrow site rời miền găng gửi thông báo đến controller, controller gửi quyền đến site mới.
- ☐ System throughput: $1/(2T + E)$. If synchronization delay giảm còn T , throughput tăng gấp đôi.
- ☐ Controller bị bottleneck, có thể dẫn đến tắt nghẽn.

Thuật toán không sử dụng token

□ Notations:

- S_i : Site i
- R_i : Request set, containing the ids of all S_i s from which permission must be received before accessing CS.
- Non-token based approaches use time stamps to order requests for CS.
- Smaller time stamps get priority over larger ones.

□ Lamport's Algorithm

- $R_i = \{S_1, S_2, \dots, S_n\}$, i.e., all sites.
- Request queue: maintained at each S_i . Ordered by time stamps.
- Assumption: message delivered in FIFO.

Thuật toán Lamport

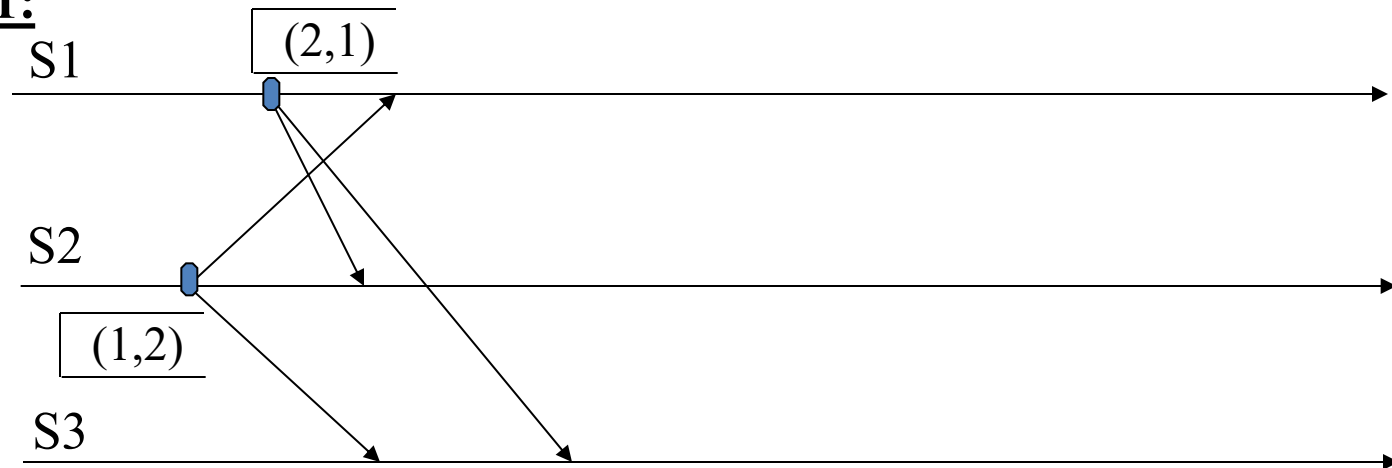
- Requesting CS:
 - Send REQUEST(t_{si}, i). (t_{si}, i): Request time stamp. Place REQUEST in *request_queue_i*.
 - On receiving the message; s_j sends time-stamped REPLY message to s_i . S_i 's request placed in *request_queue_j*.
- Executing CS:
 - S_i has received a message with time stamp larger than (t_{si}, i) from all other sites.
 - S_i 's request is the top most one in *request_queue_i*.
- Releasing CS:
 - Exiting CS: send a time stamped RELEASE message to all sites in its request set.
 - Receiving RELEASE message: S_j removes S_i 's request from its queue.

Thuật toán Lamport

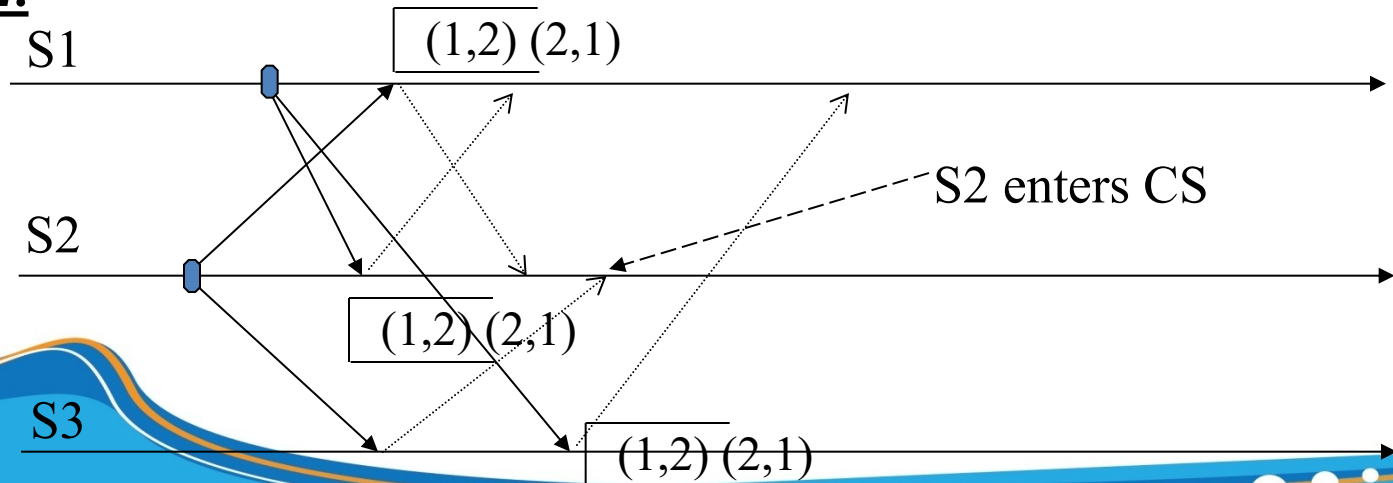
- Performance.
 - $3(N-1)$ messages per CS invocation. $(N - 1)$ REQUEST, $(N - 1)$ REPLY, $(N - 1)$ RELEASE messages.
 - Synchronization delay: T
- Optimization
 - Giữ lại reply msg. (e.g.,) S_j nhận REQUEST từ S_i sau khi nó gửi REQUEST message với time stamp lớn hơn của S_i 's. KHÔNG gửi REPLY message.
 - Messages giảm còn khoảng $2(N-1)$ tới $3(N-1)$.

Ví dụ thuật toán Lamport

Step 1:

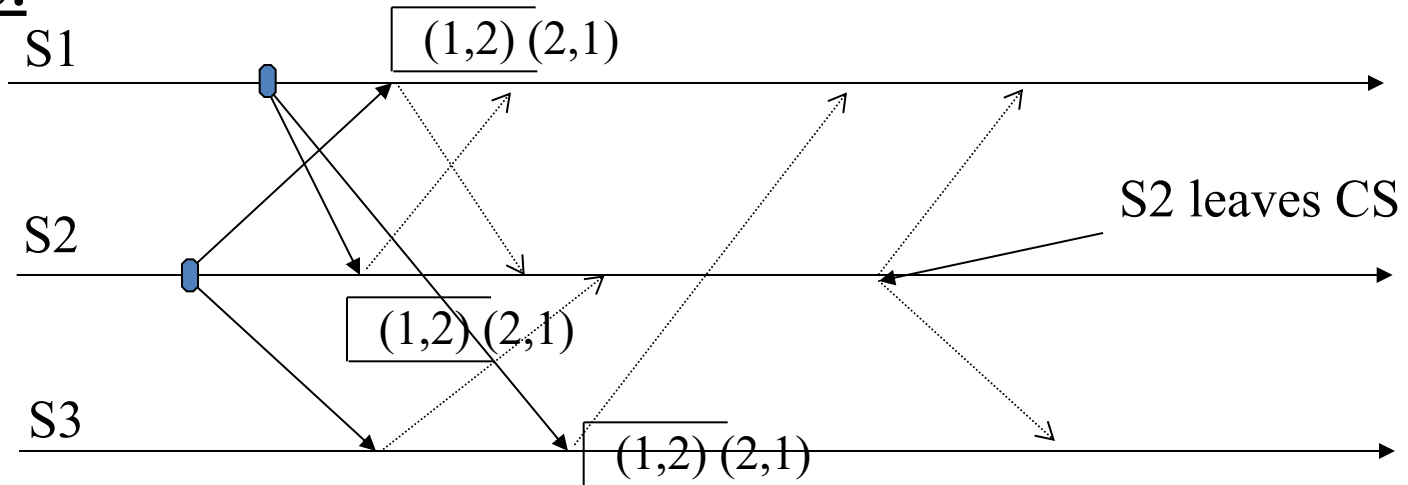


Step 2:

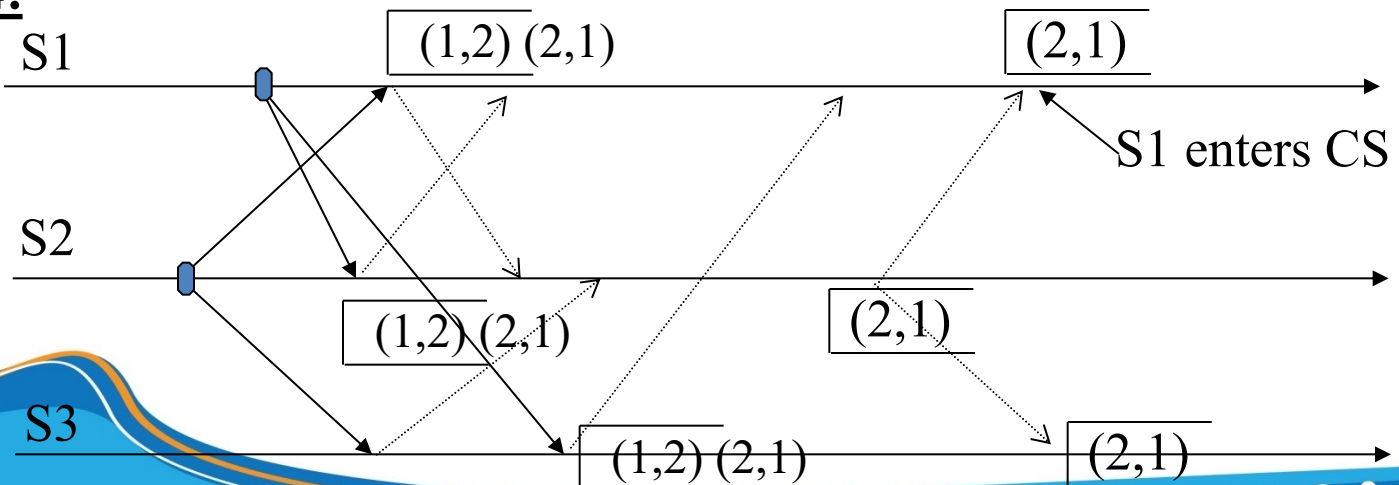


Ví dụ Lamport's:

Step 3:



Step 4:



Thuật toán Ricart-Agrawala

- Requesting critical section
 - S_i sends time stamped REQUEST message
 - S_j sends REPLY to S_i , if
 - S_j is not requesting nor executing CS
 - If S_j is requesting CS and S_i 's time stamp is smaller than its own request.
 - Request is deferred otherwise.
- Executing CS: after it has received REPLY from all sites in its request set.
- Releasing CS: Send REPLY to all deferred requests. i.e., a site's REPLY messages are blocked only by sites with smaller time stamps

Ricart-Agrawala: Performance

□ Performance:

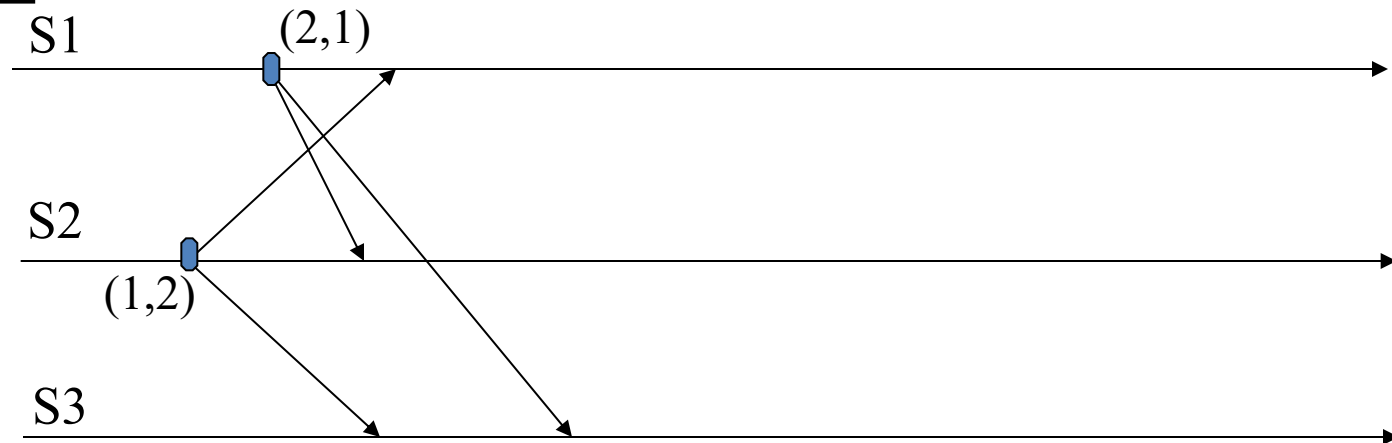
- $2(N-1)$ messages per CS execution. $(N-1)$ REQUEST + $(N-1)$ REPLY.
- Synchronization delay: T .

□ Optimization:

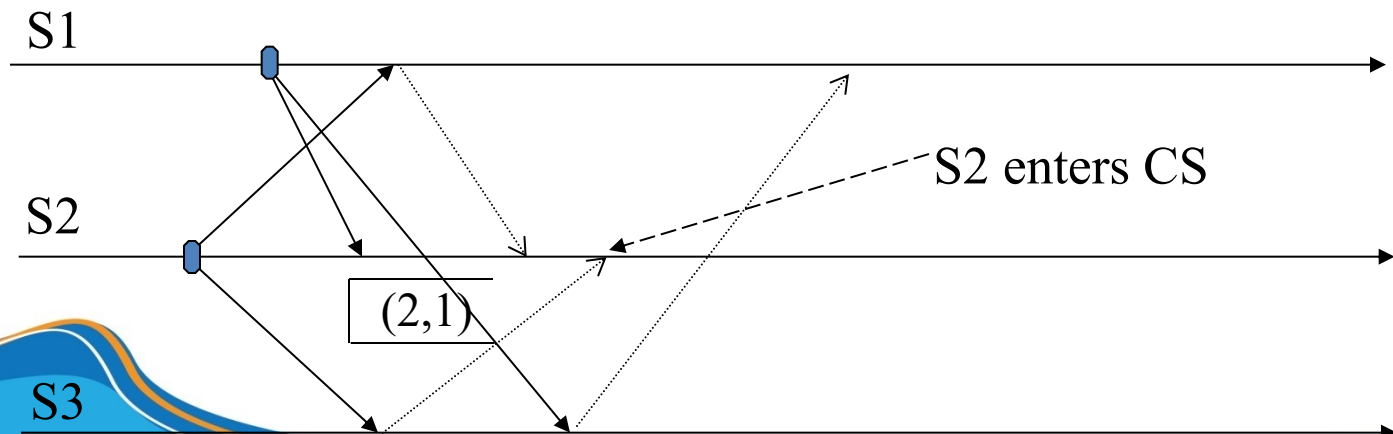
- Khi Si nhận REPLY message từ S_j -> có quyền vào miền găng cho đến khi
 - S_j gửi REQUEST message và Si gửi REPLY message.
 - Truy cập miền găng nhiều lần
 - A site requests permission from dynamically varying set of sites: 0 to $2(N-1)$ messages.

Ví dụ Ricart-Agrawala

Step 1:

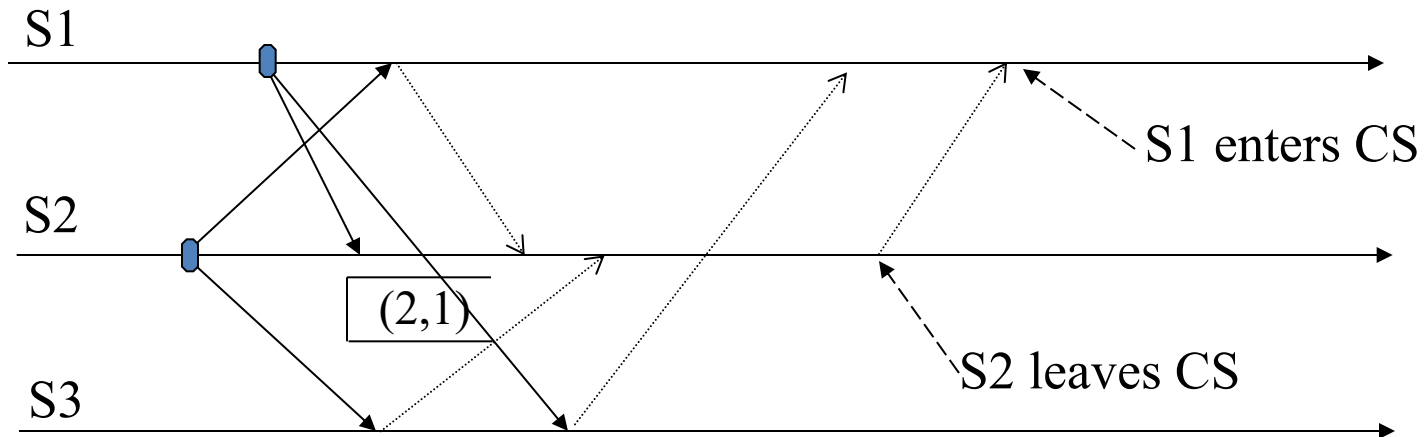


Step 2:



Ví dụ Ricart-Agrawala

Step 3:



Thuật toán Maekawa

- 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 Sets Rules:
 - 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. $N = K(K - 1) + 1 \rightarrow K = \text{square root of } N$.

Thuật toán Maekawa

- Requesting CS
 - S_i sends REQUEST(i) to sites in R_i .
 - S_j sends REPLY to S_i if
 - S_j has NOT sent a REPLY message to any site after it received the last RELEASE message.
 - Otherwise, queue up S_i 's request.
- Executing CS: after getting REPLY from all sites in R_i .
- Releasing CS
 - send RELEASE(i) to all sites in R_i
 - Any S_j after receiving RELEASE message, send REPLY message to the next request in queue.
 - If queue empty, update status indicating receipt of RELEASE.

Các tập con

- Example $k = 2$; ($N = 3$).
 - ▣ $R1 = \{1, 2\}$; $R3 = \{1, 3\}$; $R2 = \{2, 3\}$
- Example $k = 3$; $N = 7$.
 - ▣ $R1 = \{1, 2, 3\}$; $R4 = \{1, 4, 5\}$; $R6 = \{1, 6, 7\}$;
 - ▣ $R2 = \{2, 4, 6\}$; $R5 = \{2, 5, 7\}$; $R7 = \{3, 4, 7\}$;
 - ▣ $R3 = \{3, 5, 6\}$
- Algorithm in Maekawa's paper (uploaded in Lecture Notes web page).

Thuật toán Maekawa

- Performance
 - Synchronization delay: $2T$
 - Messages: 3 times square root of N (one each for REQUEST, REPLY, RELEASE messages)
- Deadlocks
 - Message deliveries are not ordered.
 - Assume S_i, S_j, S_k concurrently request CS
 - R_i intersection $R_j = \{S_{ij}\}$, $R_j R_k = \{S_{jk}\}$, $R_k R_i = \{S_{ki}\}$
 - Possible that:
 - S_{ij} is locked by S_i (forcing S_j to wait at S_{ij})
 - S_{jk} by S_j (forcing S_k to wait at S_{jk})
 - S_{ki} by S_k (forcing S_i to wait at S_{ki})
 - \rightarrow deadlocks among S_i, S_j , and S_k .

Xử lý Deadlocks

- Si *yields* to a request if that has a smaller time stamp.
- A site suspects a deadlock when it is locked by a request with a higher time stamp (lower priority).
- Deadlock handling messages:
 - ▣ FAILED: from Si to Sj -> Si has granted permission to higher priority request.
 - ▣ INQUIRE: from Si to Sj -> Si would like to know Sj has succeeded in locking all sites in Sj's request set.
 - ▣ YIELD: from Si to Sj -> Si is returning permission to Sj so that Sj can yield to a higher priority request.

Xử lý Deadlocks

- REQUEST(t_{si}, i) to S_j :
 - S_j is locked by $S_k \rightarrow S_j$ sends FAILED to S_i , if S_i 's request has higher time stamp.
 - Otherwise, S_j sends INQUIRE(j) to S_k .
- INQUIRE(j) to S_k :
 - S_k sends a YIELD (k) to S_j , if S_k has received a FAILED message from a site in S_k 's set. (or) if S_k sent a YIELD and has not received a new REPLY.
- YIELD(k) to S_j :
 - S_j assumes it has been released by S_k , places S_k 's request in its queue appropriately, sends a REPLY(j) to the top request in its queue.
- Sites may exchange these messages even if there is no real deadlock. Maximum number of messages per CS request: 5 times square root of N .

Thuật toán dựa trên token

- Unique token circulates among the participating sites.
- A site can enter CS if it has the token.
- Token-based approaches use sequence numbers instead of time stamps.
 - ▣ Request for a token contains a sequence number.
 - ▣ Sequence number of sites advance independently.
- Correctness issue is trivial since only one token is present -> only one site can enter CS.
- Deadlock and starvation issues to be addressed.

Thuật toán Suzuki-Kasami

- If a site without a token needs to enter a CS, broadcast a REQUEST for token message to all other sites.
- Token: (a) Queue of request sites (b) Array $LN[1..N]$, the sequence number of the most recent execution by a site j .
- Token holder sends token to requestor, if it is not inside CS. Otherwise, sends after exiting CS.
- Token holder can make multiple CS accesses.
- Design issues:
 - ▣ Distinguishing outdated REQUEST messages.
 - Format: $REQUEST(j,n) \rightarrow$ j th site making n th request.
 - Each site has $RNi[1..N] \rightarrow RNi[j]$ is the largest sequence number of request from j .
 - ▣ Determining which site has an outstanding token request.
 - If $LN[j] = RNi[j] - 1$, then S_j has an outstanding request.

Thuật toán Suzuki-Kasami

□ Passing the token

- After finishing CS
- (assuming S_i has token), $LN[i] := RN[i]$
- Token consists of Q and LN . Q is a queue of requesting sites.
- Token holder checks if $RN[j] = LN[j] + 1$. If so, place j in Q .
- Send token to the site at head of Q .

□ Performance

- 0 to N messages per CS invocation.
- Synchronization delay is 0 (if the token holder repeats CS) or T .

Ví dụ Suzuki-Kasami

Step 1: S1 has token, S3 is in queue

Site	Seq. Vector RN	Token Vect. LN	Token Queue
S1	10, 15, 9	10, 15, 8	3
S2	10, 16, 9		
S3	10, 15, 9		

Step 2: S3 gets token, S2 in queue

Site	Seq. Vector RN	Token Vect. LN	Token Queue
S1	10, 16, 9		
S2	10, 16, 9		
S3	10, 16, 9	10, 15, 9	2

Step 3: S2 gets token, queue empty

Site	Seq. Vector RN	Token Vect. LN	Token Queue
S1	10, 16, 9		
S2	10, 16, 9	10, 16, 9	<empty>
S3	10, 16, 9		

Thuật toán Singhal's Heuristic

- Instead of broadcast: each site maintains information on other sites, guess the sites likely to have the token.
- 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 $j = i-1 .. 1$; $SN_i[j] := 0$, $j = 1..N$. S1 (Site 1) is in state H.
 - Token: $TSV[j] := N$ & $TSN[j] := 0$, $j = 1 .. N$.

Thuật toán Singhal's Heuristic ...

- 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 $SV_i[i] := E$.

Singhal's Heuristic ...

- Releasing CS
 - $SV_i[i] := N, TSV[i] := N$. Then, do:
 - For other S_j : if $(SN_i[j] > TSN[j])$, then $\{TSV[j] := SV_i[j]; TSN[j] := SN_i[j]\}$
 - else $\{SV_i[j] := TSV[j]; SN_i[j] := TSN[j]\}$
 - If $SV_i[j] = N$, for all j , then set $SV_i[i] := H$. Else send token to a site S_j provided $SV_i[j] = R$.
 - Fairness of algorithm will depend on choice of S_i , 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 .

Ví dụ Singhal

Sn	R	R	R	R	N
S4	R	R	R	N	N
S3	R	R	N	N	N
S2	R	N	N	N	N
S1	H	N	N	N	N
	1	2	3	4	n

(a) Initial Pattern

Each row in the matrix has increasing number of Rs.

Sn	R	R	R	R	N
S4	R	R	R	N	N
S2	R	N	R	N	N
S1	N	N	R	N	N
S3	N	N	H	N	N
	1	2	3	4	n

(b) Pattern after S3 gets the token from S1.

Stair case is pattern can be identified by noting that S1 has 1 R and S2 has 2 Rs and so on. Order of occurrence of R in a row does not matter.

Ví dụ Singhal

- Assume there are 3 sites in the system. Initially:
 - Site 1: $SV1[1] = H$, $SV1[2] = N$, $SV1[3] = N$. $SN1[1]$, $SN1[2]$, $SN1[3]$ are 0.
 - Site 2: $SV2[1] = R$, $SV2[2] = N$, $SV2[3] = N$. SNs are 0.
 - Site 3: $SV3[1] = R$, $SV3[2] = R$, $SV3[3] = N$. SNs are 0.
 - Token: TSVs are N. TSNs are 0.
- Assume site 2 is requesting token.
 - S2 sets $SV2[2] = R$, $SN2[2] = 1$.
 - S2 sends REQUEST(2,1) to S1 (since only S1 is set to R in SV[2])
- S1 receives the REQUEST. Accepts the REQUEST since $SN1[2]$ is smaller than the message sequence number.
 - Since $SV1[1]$ is H: $SV1[2] = R$, $TSV[2] = R$, $TSN[2] = 1$, $SV1[1] = N$.
 - Send token to S2
- S2 receives the token. $SV2[2] = E$. After exiting the CS, $SV2[2] = TSV[2] = N$. Updates SN, SV, TSN, TSV. Since nobody is REQUESTing, $SV2[2] = H$.
- Assume S3 makes a REQUEST now. It will be sent to both S1 and S2. Only S2 responds since only $SV2[2]$ is H ($SV1[1]$ is N now).

Thuật toán Raymond

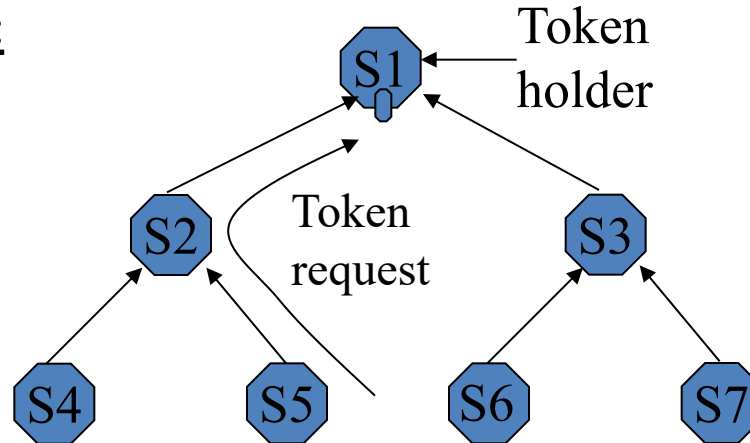
- Sites are arranged in a logical directed tree. Root: token holder. Edges: directed towards root.
- Every site has a variable *holder* that points to an immediate neighbor node, on the directed path towards root. (Root's holder point to itself).
- Requesting CS
 - If S_i does not hold token and request CS, sends REQUEST *upwards* provided its *request_q* is empty. It then adds its request to *request_q*.
 - Site on path to root receiving REQUEST -> propagate it up, if its *request_q* is empty. Add request to *request_q*.
 - Root on receiving REQUEST -> send token to the site that forwarded the message. Set *holder* to that forwarding site.
 - Any S_i receiving token -> delete top entry from *request_q*, send token to that site, set *holder* to point to it. If *request_q* is non-empty now, send REQUEST message to the *holder* site.

Thuật toán Raymond...

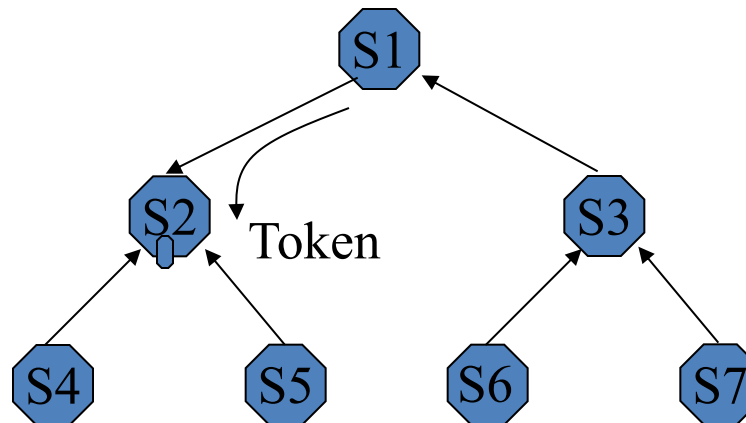
- Executing CS: getting token with the site at the top of *request_q*. Delete top of *request_q*, enter CS.
- Releasing CS
 - ▣ If *request_q* is non-empty, delete top entry from *q*, send token to that site, set *holder* to that site.
 - ▣ If *request_q* is non-empty now, send REQUEST message to the *holder* site.
- Performance
 - ▣ Average messages: $O(\log N)$ as average distance between 2 nodes in the tree is $O(\log N)$.
 - ▣ Synchronization delay: $(T \log N) / 2$, as average distance between 2 sites to successively execute CS is $(\log N) / 2$.
 - ▣ Greedy approach: Intermediate site getting the token may enter CS instead of forwarding it down. Affects fairness, may cause starvation.

Ví dụ Raymond

Step 1:

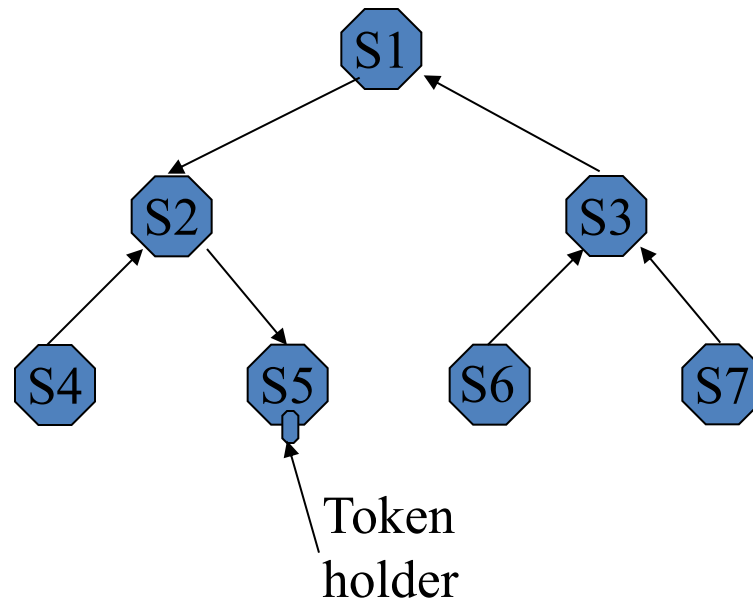


Step 2:



Ví dụ Raymond...

Step 3:



So sánh

Non-Token	Resp. Time(l)	Sync. Delay	Messages(l)	Messages(h)
Lamport	$2T+E$	T	$3(N-1)$	$3(N-1)$
Ricart-Agrawala	$2T+E$	T	$2(N-1)$	$2(N-1)$
Maekawa	$2T+E$	$2T$	$3*sq.rt(N)$	$5*sq.rt(N)$
Token	Resp. Time(l)	Sync. Delay	Messages(l)	Messages(h)
Suzuki-Kasami	$2T+E$	T	N	N
Singhal	$2T+E$	T	$N/2$	N
Raymond	$T(\log N)+E$	$T\log(N)/2$	$\log(N)$	4