

AN EFFICIENT TOKEN BASED ALGORITHM FOR MUTUAL EXCLUSION IN DISTRIBUTED SYSTEM

Fahim Kawsar* Shariful Hasan Shaikot**Md. Shahariar Saikat⁰
Md. Abdur Razzaque⁰⁰ M. A. Mottalib⁰⁰

ABSTARCT

Many distributed computations involving the sharing of resources among various processes require that a resource be allocated to a single process at a time. Therefore, mutual exclusion is a fundamental problem in any distributed computing system. This problem must be solved to synchronize the access to shared resources in order to maintain their consistency and integrity. The major goal of this paper is to get the reader acquainted with a new approach towards the ring based technique for mutual exclusion in a distributed system. An algorithm is proposed based on the idea of generating token by the competing processes to enter the critical section and thus eliminating idle time message passing and reducing communication overhead.

Keywords: Mutual exclusion, critical section, token, ring structure.

1. INTRODUCTION

A distributed computing system is a collection of autonomous computing sites that do not share a global or common memory and communicate solely by exchanging messages over a communication facility. In a distributed computing system any given site (also referred to as "node") has only a partial or incomplete view of the total system and a system-wide common clock does not exist. Processes must share common hardware or software resources, cooperating in such a way that they can work in parallel and independently of each other. The access to a shared resource must be synchronized to ensure that only one process is making use of the resource at a given time. The problem of coordinating the execution of critical sections by each process is solved by providing mutually exclusive access in time to the critical section (CS). Each process must request permission to enter its critical section and must release it after it has completed its execution. A mutual exclusion algorithm must satisfy the following requirements [1, 2]:

- i. At most one process can execute its critical section at a given time.
- ii. If no process is in its critical section, any process requesting to enter its critical section must be allowed to do so at finite time.
- iii. When competing processes concurrently request to enter their respective critical sections, the selection cannot be postponed indefinitely.
- iv. A requesting process cannot be prevented by another one to enter its critical section within a finite delay.

To simplify, an algorithm must provide mutually exclusive access to the source, ensure deadlock freedom, ensure starvation freedom, and must provide some fairness in the order that requests are granted. The algorithm presented in this paper is based on the token ring approach and satisfies the mentioned requirements in a way that minimize the communication overhead and ensure deadlock freedom, ensure starvation freedom. The

* Nakajima Lab, Waseda University, Japan.

** Dept. of Computer Science, Asian University of Bangladesh, Dhaka, Bangladesh.

⁰ Grameen Software Limited, Dhaka, Bangladesh.

⁰⁰ CIT Dept. Islamic University of Technology, Gazipur-1704, Bangladesh.

competing nodes generate a token for the permission to enter CS. The token traverses the logical ring structure. A node can enter CS if and only if it receives back its generated token.

The performance of the algorithm presented here will be evaluated using the total number of messages required for a node to enter the critical section as a criterion. Message traffic should be minimized in order to decrease the overhead in the communications network.

The remainder of the paper is organized as follows. Section 2 compares our work with existing research in distributed scheduling. The algorithm is presented in Section 3. Finally, Section 4 concludes the paper.

2. PREVIOUS WORKS

Two approaches have been used to implement a mutual exclusion mechanism in a distributed computing system. In a *centralized* approach, one of the nodes functions as a central coordinator. The central coordinator is fully responsible for having all the information of the system and for granting permission to make use of a shared resource. In a *distributed* approach, the decision-making is distributed across the entire system. This paper only considers the distributed approach. Distributed mutual exclusion algorithms are designed based on two basic principles: the existence of token in the system or the collection of permission from nodes in the system.

2.1 Permission-based Algorithm

All the permission-based algorithms are introduced so far basically work in the same way. The node that wants to enter the critical section sends messages to other processors. Also, associated with each request there is a timestamp. When there are competitions for the critical section, the one with the lowest timestamp should enter first.

In Lamport's event ordering mutual exclusion algorithm [3], a node that wants to enter the critical section, broadcasts a message to all nodes in the system. The node that made the request enters the critical section if it received responses from all other nodes. After the node finished with the critical section, it again broadcasts a message to all other nodes. For a N node system a total of $3(N-1)$ messages are required to handle one request.

Ricart-Agarwala's algorithm [4] is very similar to Lamport's algorithm. The difference is that in Ricart-Agrawala's algorithm, the response message is deferred. Similar to Lamport's algorithm, Ricart-Agrawala's algorithm requires totally ordered events and all nodes being alive. About the number of messages, it does not need the release message, so it requires $2(N-1)$ messages for handling one request.

Maekawa's algorithm [5] associates a set of nodes is with each node, and this set has a nonempty intersection with every set associated with each other node. A node i must obtain permission from all other nodes in its home set S_i before it can enter its critical section (CS). The number of messages required to handle a request is 3 times the size of the request set. For a system with N nodes, the size of each request set is roughly square root of N. So total $3\sqrt{N}$ messages are required to handle a request.

2.2 Token-based Algorithm

The simplest of token-based algorithm is the Token Ring algorithm [6]. In this algorithm, the nodes in the system form a logical ring. A token is passed around the ring.

A node can enter the critical section if it holds the token. In average $N/2$ messages are required to handle one request in a N node system.

In Suzuki-Kasami's broadcast algorithm [7], when a node wants to enter the critical section; it broadcasts a message to all other nodes. Whoever holds the token sends the token directly to the node that wants the token. The algorithm requires N messages for handling each request.

In Raymond's tree-based algorithm [8] the token is always kept at the root node. When a node wants to enter the critical section, it sends a request to its parent. The parent sends a request to its parent, recursively, this request will reach the root node. The root node, upon receiving the request sends the token down to the child that requested the token and is on top of the request queue. Once the node gets the token, it can enter the critical section. In this algorithm, it requires an average of $2\log N$ messages for handling each request.

3. THE PROPOSED ALGORITHM

Our proposed algorithm is based on the token ring algorithm. The following assumptions and conditions for the distributed environment are considered while designing the algorithm:

- i. All nodes in the system are assigned unique identification numbers from 1 to N .
- ii. There is only one requesting process executing at each node. Mutual exclusion is implemented at the node level.
- iii. Processes are competing for a single resource.
- iv. At any time, each process initiates at most one outstanding request for mutual exclusion.
- v. All the nodes in the system are fully connected.

The following aspects about the reliability of the underlying communications network should be considered.

- Message delivery guaranteed.
- Message-order preservation.
- Message transfer delays are finite, but unpredictable.
- The topology of the network is known.

The network may be of any topology with no inherent ordering of the processes. In software a logical ring is constructed in which each process is assigned a position in the ring as shown in Figure 1. The ring positions may be allocated in numerical order of network addresses. All the matter is that each process knows who is next in line after itself.

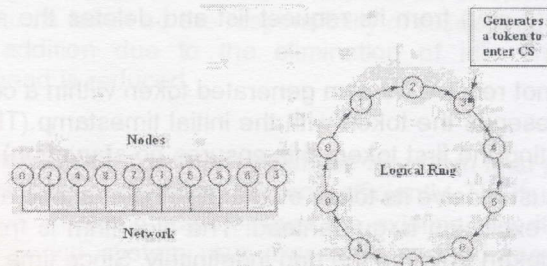


Figure 1: Logical ring of unordered node on the network.

The basic idea is that whenever a process needs to enter a critical section it generates a token and puts it in the network. On receiving a token a node reacts in the following manner

- a) A process P wants to enter the critical section and generates a token, makes a copy of the token in it's node & passes it to the next node. The token structure is

```
token
{
  node no;
  timestamp;
};
```

- b) A process Q receives a token. It reacts to the token in the following ways:

- If Q has no intention to enter into the critical section it simply passes the token to the next node.
- If Q is in the critical section it puts the token in its request list. The request list structure is

```
request list
{
  token
  next list;
};
```

When Q exits the critical section it sends the tokens to the next node sequentially (if any) from its request list.

- If Q has already generated a token but not yet received that back it compares the incoming token's timestamp with it's generated token's timestamp. If it's token's timestamp is higher then it passes the token to the next node otherwise it puts the token in it's request list.
 - If the timestamps in both the tokens are equal then Q looks for the node_no. If Q's node_no is lower (lowest number highest priority) then it adds the token in its request list otherwise passes it to the next node.
- c) If the process P receives a token, it compares the token with the stored token copy. If it matches then (that is no node wants to enter into the critical section) it enters into the critical section. When P exits the critical section it sends the tokens to the next node sequentially (if any) from its request list and deletes the associated copy and original token.
- d) If process P does not receive it's own generated token within a certain timeout period (token is lost), P resends the token with the initial timestamp. (The timestamp that is used while generating the first token, this ensures no starvation)

Since each node must receive its token eventually and at a time at most one node can have its token mutual exclusion is guaranteed. The algorithm is free from deadlock as there is no way that a token traverse the ring indefinitely. Since time stamping is used so it ensures starvation freedom.

The timeout period for regenerating token must be designed in such a manner that unnecessarily no token should be resend. If multiple copies of the same token arrive at a originator node after it performs its operation in the CS or it is in the CS, it absorbs the token. (Since there is no copy of any token to be matched)

Another problem is that if a node dies there is no way to detect it so token can be lost. To avoid this acknowledgement of receiving a token may be used. Thus node failure can easily be detected. At that point the dead node can be removed from the group and the token holder can throw the token over the head of the dead process to the next node down the line. Of course, doing so requires that everyone maintains the current ring configuration.

From the working principle of this algorithm it is visible that for a N node system we need N messages to be transferred for handling one request of entering critical section.

3.1 Performance Gain

The most attractive feature of this newly proposed method is that there is no need for passing the token around the ring when no node requires it that is idle period token passing is eliminated. It reduces communication overhead to a great extent. The Performance gain of the proposed method is shown in Table 1.

Table 1 : Performance gain of proposed method.

| Algorithm | Messages/Request |
|------------------------------------|--------------------|
| Lamport's Algorithm | $3(N-1)$ |
| Ricart-Agarwals Algorithm | $2(N-1)$ |
| Maekawas Algorithm | $3vN$ or $5vN$ |
| Token Ring Algorithm | Avg $N/2$ |
| Suzuki-Kasamis Broadcast Algorithm | N |
| Raymonds Tree-based Algorithm | Avg $2\text{Log}N$ |
| Ours Algorithm | N |

$N = \text{Number of Node}$, $\text{Avg} = \text{Average}$

So from the Table 1 we find that this newly proposed algorithm will provide similar performance as some the existing techniques provide. But the attractive potentiality of this algorithm lies in error handling techniques. As in all token ring algorithm if the token is lost then detection and regeneration of the token is a big problem. Who will be responsible for regenerating the token and how to detect the loss of the token?

Our proposed method handles the loss of token in an efficient manner and as each node is responsible for its token so the responsibility of resending tokens completely lies on each node. In addition due to the elimination of idle time message passing communication overhead is reduced.

4. CONCLUSION

The motivation towards the development of this algorithm is to present a method that guarantee mutual exclusion and works fairly and efficiently. The fault tolerance capability of this algorithm clearly makes it superior over the existing algorithms. Considering the rapid growth of distributed system, our presented method may provide a lucrative approach towards the solution of mutual exclusion problem.

REFERENCES

- [1] Maekawa, M.; Oledchoeft, A.E.; and Oledchoeft, R.R., "Operating Systems, Advanced Concepts," Benjamin-Cummings, 1987.
- [2] Silberschatz, A. and Peterson, J.L., "Operating System Concepts," Addison-Wesley, Alternate edition, 1988.
- [3] Lamport, L., "Time, clocks, and the ordering of events in a distributed system," Communications of the ACM, vol. 21, no. 7, July 1978, pp. 558-565.
- [4] Ricart, G. and Agrawala, A., "An optimal algorithm for mutual exclusion in computer networks," Communications of the ACM, vol. 24, no. 1, Jan. 1981, pp. 9-17.
- [5] Maekawa, M., "A sqrt(n) algorithm for mutual exclusion in decentralized systems," ACM Transactions on Computer Systems, vol 3, no. 2, May 1985, pp. 145-159.
- [6] Agrawal, D., Elabbai, A., "An efficient and fault-tolerant solution for distributed mutual exclusion," ACM Transactions on Computer Systems, vol. 9, no. 1, Feb. 1991, pp. 1-20.
- [7] Suzuki, I. and Kasami, T., "A distributed mutual exclusion algorithm," ACM Transactions on Computer Systems, vol. 3, no. 4, Nov. 1985, pp. 344-349.
- [8] Raymond, K., "A tree-based algorithm for distributed Mutual Exclusion," ACM Transactions on Computer Systems, vol. 7, no. 1, Feb. 1989, pp. 61-77.

REFERENCES

- [1] Maekawa, M.; Oledchoeff, A.E.; and Oledchoeff, R.R., "Operating Systems, Advanced Concepts," Benjamin-Cummings, 1987.
- [2] Silberschatz, A. and Peterson, J.L., "Operating System Concepts," Addison-Wesley, Alternate edition, 1988.
- [3] Lamport, L., "Time, clocks, and the ordering of events in a distributed system," Communications of the ACM, vol. 21, no. 7, July 1978, pp. 558-565.
- [4] Ricart, G. and Agrawala, A., "An optimal algorithm for mutual exclusion in computer networks," Communications of the ACM, vol. 24, no. 1, Jan. 1981, pp. 9-17.
- [5] Maekawa, M., "A sqrt(n) algorithm for mutual exclusion in decentralized systems," ACM Transactions on Computer Systems, vol 3, no. 2, May 1985, pp. 145-159.
- [6] Agrawal, D., Elabbai, A., "An efficient and fault-tolerant solution for distributed mutual exclusion," ACM Transactions on Computer Systems, vol. 9, no. 1, Feb. 1991, pp. 1-20.
- [7] Suzuki, I. and Kasami, T., "A distributed mutual exclusion algorithm," ACM Transactions on Computer Systems, vol. 3, no. 4, Nov. 1985, pp. 344-349.
- [8] Raymond, K., "A tree-based algorithm for distributed Mutual Exclusion," ACM Transactions on Computer Systems, vol. 7, no. 1, Feb. 1989, pp. 61-77.

8. SI units must be used in the manuscript. However, other units may be used in parenthesis.
9. Tables should be referred to in consecutive Arabic numerical. Each table must have a table caption.
10. Line drawings must be in a form suitable for reproduction e.g., laser print, drawn in indian ink on white paper or on tracing paper. Photographs should have a glossy finish. Each figure must have a number and a figure caption. Electronic mode is preferred.
11. References should be set out in alphabetical order of the author's last name in a list at the end of the article. They should be given in standard form as in the following examples:
 - (a) Journal
Bloomer G. and Write A., "Scheduling of Vehicles from Factory to Depot." *Operations Research*, Vol. 12, January, pp. 590-598, 1984.
 - (b) Book
Best, John., and Kahn, James V., *Research in Education*, Prentice-Hall, New Jersey, 1986.
 - (c) Monograph
Syedali, M. M. " Computer Controlled Car", Thesis, M.Sc. Engineering, Department of Mechanical Engineering, BUET, 1990
 - (d) Conference paper
Hasan, M. and Ullah, M.S. "Tourism development in the Chittagong Hill Tracts of Bangladesh after the peace agreement of 1997", a paper submitted to the Regional Conference on physical mobility and development in the mountains, at Tribuvan University, 15-17, March, 2000 Kathmandu, Nepal, pp. 12.
 - (e) Unpublished paper
Ahmadi, R and Tang : Production Allocation with Dual Provisioning, Working Paper, Anderson Graduate School of Management, UCLA (1991)
12. The University does not accept responsibility for loss or damage of manuscript while in mail.
13. The responsibility for opinions in the contributions rests entirely on their authors.
14. The author(s) must submit declaration that the paper was not published elsewhere.
15. In case of joint papers, communication will be made with the first author
16. The University will reserve the copyright of the paper once it is accepted for publication in the journal. The authors must obtain written permission from REASP. IUT for publication elsewhere.

Procedures for acceptance of papers and publications :

1. Papers submitted for publication will be referred to the listed reviewers for assessment. However, the Editorial Board will make initial screening of the papers.
2. After the assessment, the authors may be requested to modify/clarify certain points.
3. Accepted (modified/corrected) papers will be published in the next issue of the journal.