

Performance Analysis of Coordinated Check pointing Protocols

Krishnamurthy Ramasubramanian

Asistant Professor, Department of CSE
Samskruti College of Engineering and Technology,
Kondapur (V), Ghatkesar (Mdl), Hyderabad

Abstract

Check-pointing can be coordinated, uncoordinated, or communication-induced. Log-based protocols combine check pointing with logging of nondeterministic events, encoded in tuples called determinants. Depending on how determinants are logged, log-based protocols can be pessimistic, optimistic, or causal. Throughout the survey, we highlight the research issues that are at the core of rollback recovery and present the solutions that currently address them. We also compare the performance of different rollback-recovery protocols with respect to a series of desirable properties and discuss the issues that arise in the practical implementations of these protocols. This paper presents the protocols which have been appeared in the literature for checkpointing in distributed systems.

Keywords: Checkpoint, checkpointing protocols, Distributed systems, rollback recovery, fault tolerant computing. Message-logging

1. Introduction

The term Distributed Systems consists of several computers that do not share memory or a clock, each computer having its own memory and runs its own operating system and communicate with each other by exchanging messages over a communication network [22]. A mobile distributed system (MDS) is a distributed system where some of processes are running on mobile hosts (MHs). A mobile distributed system having fixed and mobile station interconnected through a communication network. The fixed station is located at the fixed location and the mobile station moves from one location to another in the network. Mobile Hosts (MHs) are becoming common in distributed systems due to their accessibility, cost, and mobile connectivity. The term "mobile" means able to move while retaining its network

connection. Checkpoint-based rollback-recovery techniques can be classified into three categories: uncoordinated checkpointing, coordinated checkpointing, and communication-induced checkpointing. A distributed system containing more several processes that execute on geographically dispersed computers and collaborate via message-passing with each other to achieve a common goal [19]. checkpoint is one of the most prominent techniques for providing fault-tolerance, and can also be used for debugging and migration in both uniprocessor and distributed systems [20,21]. particularly, checkpointing is the act of saving a program's state on stable storage, and restart is the act of restarting an application from its saved state. Especially, if an application takes periodic checkpoints, then in case of failure, it is possible to restart it from the latest checkpoint, thereby avoiding loss of all the computation that was carried before that checkpoint. Many distributed check pointing protocols produce control overhead [22]. control overhead is the overhead due to control information. During the past years a large number of check pointing protocols have been proposed for distributed systems [5]. Most of these protocols were never implemented or tested. The distributed mobile systems use check pointing for providing fault tolerance. In this case, when fault or failures of process occur, an application with mobile should rollback to a consistent global checkpoint as close as possible to the end of the computation. A local checkpoint is a recorded state of process. A global checkpoint is a set of local checkpoints one from each process in a distributed system [6]. A consistent global checkpoint is one in which every message that has been received is also shown to have been sent in the corresponding state of sender. A distributed mobile system contains of both static Mobile service

Stations and Mobile Hosts. A set of wireless communication links and dynamic links can be established between a mobile service station and mobile host, and a set of high-speed communication link is assumed between the mobile service stations. A mobile service station may communicate with a number of mobile hosts but a mobile host communicates with the rest of the system via the mobile service station it is connected to.

2. Check pointing Protocols

Check pointing is a standard method for the repair of faults in systems. The idea is to save the state of the system on a stable periodic to prevent breakdowns. That way when you restart after a power failure, the state saved newest restored and execution resumes its course before the crash. The overall status of a distributed system is defined by the union of local states of all processes belonging to the system. Taking checkpoints is the process of periodically saving the state of a running process to durable storage. Checkpointing allows a process that fails to be restarted from the point its state was last saved, or its checkpoint. If the host processor has not failed, temporal redundancy can be used to roll back and restart the process on the same platform. As in other systems, this method is widely used in grids [21, 22]. Otherwise, if the host has failed, the process may be migrated, or transferred, to a different execution environment where it can be restarted from a checkpoint (a technique also referred to as failover). This section begins by discussing checkpoint and process migration methods used in commercial and science grid systems that are based on methods used in high performance cluster computing. This is followed by discussion of new methods being developed or adapted for scaled grid environments, together with related issues that need to be resolved. Most notable is the issue of finding efficient methods for checkpointing many concurrent, intercommunicating processes, so that in the event of failure, they can resume from a common saved state [9]. Check pointing can be initiated either from within grid systems or within applications. There are two main classes of protocols: coordinated checkpointing and message logging.

A. Coordinated checkpointing protocols:

Coordinated checkpointing is an attractive approach for transparently adding fault tolerance to distributed applications without requiring additional programmer efforts. In this approach, the state of each process in the system is periodically saved on stable storage, which is called a checkpoint of the process. To recover

from a failure, the system restarts its execution from a previous error free, consistent global state recorded by the checkpoints of all processes. More specifically, the failed processes are restarted on any available machine and their address spaces are restored from their latest checkpoints on stable storage. Other processes may have to roll back to their checkpoints on stable storage in order to restore the entire system to a consistent state. Coordinated checkpointing simplifies failure recovery and eliminates domino effects in case of failures by preserving a consistent global checkpoint on stable storage. However, the approach suffers from high overhead associated with the checkpointing process. Two approaches are used to reduce the overhead: First is to minimize the number of synchronization messages and the number of checkpoints, the second is to make the checkpointing process non-blocking. The protocol requires processes coordinate their checkpoints to form a consistent global state. A global state is consistent if it does not include any orphan messages (i.e, a message received but not already sent). This approach simplifies the recovery and avoids the domino effect, since every process always restarts at the resume point later. Also, the protocol requires each process to maintain only one permanent checkpoint in stable storage, reducing the overhead due to storage and release of checkpoints (garbage collection) Its main drawback however is the large latency that require interaction with the outside world, in this case the solution is to perform a checkpoint after every input / output. To improve the performance of the backup coordinated, several techniques have been proposed. We have implemented as non-blocking coordinated checkpointing and Communication induced checkpointing

1) **Non-blocking coordinated checkpointing** a non-blocking checkpointing algorithm does not require any process to suspend its underlying computation. When processes do not suspend their computations, it is possible for a process to receive a computation message from another process which is already running in a new checkpoint interval. If this situation is not properly dealt with, it may result in an inconsistency. This algorithm uses markers to coordinate the backup, and operates under the assumption of FIFO channels. a comparison of protocols for coordinated checkpoint blocking and non-blocking has been made. Experiments have shown that the synchronization between nodes induced by the protocol blocking further penalize the

performance of the calculation with a non-blocking protocol. However, using frequencies of taken checkpoints usual performance of the blocking approach is better on a cluster to high-performance communications.

2) **Communication induced checkpointing** this protocol defines two types of checkpoints [19]: local checkpoints taken by processes independently, to avoid the synchronization of coordinated backup and forced checkpoints based on messages sent and received and dependency information carried 'piggyback' on these posts, so to avoid the domino effect of uncoordinated backup, ensuring the advancement of online collection. Unlike coordinated checkpoint protocols, the additional cost due to the medium access protocol disappears because the protocol does not require any message exchange to force a checkpoint: this information is inserted piggyback on the messages exchanged.

B. Message-Logging protocols: Message logging is a common technique used to build systems that can tolerate process crash failure. These protocols required that each process occurs. Indeed, during the process execution, the determinants of messages are stored in volatile memory, before being saved periodically on stable support. The storage stable memory is asynchronous: the protocol does not require the application to be blocked during the backup memory stable. Induced latency is then very low. However, a failure may occur before the messages are saved on stable storage. In this case, the information stored in volatile memory of the process down is lost and the messages sent by this process are orphaned. This can produce a domino effect of rollbacks, which increases the recovery time. Thus, message logging protocols implement an abstraction of a resilient process in which the crash of a process is translated into intermittent unavailability of that process. All message logging protocols require that the state of a recovered process be consistent with the states of the other processes. This consistency requirement is usually expressed in terms of orphan processes, which are surviving processes whose states are inconsistent with the recovered state of crashed process. Thus, in the terminology of message logging, message logging protocols must guarantee that there are no orphan processes, either through careful logging of through a somewhat complex recovery protocol. The logging mechanism uses the fact that a process can be modeled as a sequence of deterministic state intervals, each event begins with a non-deterministic. An event may be receiving a

message, or issued or other event in the process. It is deterministic if from a given initial state, it always happens at the same final state. [19] The principle of Logging is to record on a reliable storage any occurrences of non-deterministic events to be able to replay them in recovering from a failure. During execution, each process performs periodic backups of their states, and recorded in a log information about messages exchanged between processes. There are three message-logging categories: pessimistic, optimistic, and causal.

i) **Pessimistic message-logging**

This protocol was designed under the assumption that a failure may occur after any nondeterministic event (i.e. message reception). Then, each message is saved on a stable storage before to be delivering to the application. These protocols are often made reference to the synchronized because when logging process logs an event of nondeterministic stable memory, it waits for an acknowledgment to continue its execution. In a pessimistic logging system, the status of each process can be recovered independently. This property has four advantages:

- Process can send messages to the outside without using a special protocol
- The process restarted at the most recent checkpoint.
- Recovery is simple because the effects of a failure are limited only on the fail process
- The garbage collector is simple

The main drawback is the high latency of communications, which results in degradation of the applications response time. Several approaches have been developed to minimize synchronizations:

- The use of semiconductor memories such as nonvolatile stable support
- The sender based message logging (SBML) [14] which preserves the determinant or the message in the volatile memory of the transmitter, instead of a remote memory

ii) **Optimistic message-logging**

This protocol uses the assumption that the logging of a message on reliable support will be complete before a failure, the determinants of messages are stored in volatile memory, before being saved periodically on stable support. The storage stable memory is asynchronous: the protocol does not require the application to be blocked during the backup memory stable. Induced latency is then very low. However, a failure may occur before the messages are saved on stable storage. In this case, the information stored in volatile memory of the process down is lost and the

messages sent by this process are orphaned. This can produce a domino effect of rollbacks, which increases the recovery time.

iii) Causal message-logging

This protocol combines the advantages of both previous methods. As optimistic logging, it avoids the synchronized access to stable, except during the input / output. As pessimistic logging, it allows the process to make interactions with the outside world independently, and does not create process orphan. Causal logging protocols piggyback determinants of messages previously received on outgoing messages so that they are stored by their receivers.

3. Checkpointing Protocols in Comparison

Many checkpointing protocols were incepted at a time where the communication overhead far exceeded the overhead of accessing stable storage. Furthermore, the memory available to run processes tended to be small. These tradeoffs naturally favored uncoordinated checkpointing schemes over coordinated checkpointing schemes. Current technological trends however have reversed this tradeoff. In modern systems, the overhead of coordinating checkpoints is negligible compared to the overhead of saving the states [10]. Using concurrent and incremental checkpointing, the overhead of either coordinated or uncoordinated checkpointing is essentially the same. Therefore, uncoordinated checkpointing is not likely to be an attractive technique in practice given the negligible performance gains. These gains do not justify the complexities of finding a consistent recovery line after the failure, the susceptibility to the domino effect, the high storage overhead of saving multiple checkpoints of each process, and the overhead of garbage collection. It follows that coordinated checkpointing is superior to uncoordinated checkpointing when all aspects are considered on the balance. A recent study has also shed some light on the behavior of communication-induced checkpointing [20]. It presents an analysis of these protocols based on a prototype implementation and validated simulations, showing that communication-induced checkpointing does not scale well as the number of processes increases. The occurrence of forced checkpoints at random points within the execution due to communication messages makes it very difficult to predict the required amount of stable storage for a particular application run. Also, this unpredictability affects the policy for placing local checkpoints and makes CIC

protocols cumbersome to use in practice.

Furthermore, the study shows that the benefit of autonomy in allowing processes to take local checkpoints at their convenience does not seem to hold. In all experiments, a process takes at least twice as many forced checkpoints as local, autonomous ones.

Check pointing	Advantages	Disadvantages
A	Process coordintes thecheckpointing	Large delay in computing the output
B	Lower run time overhead during execution	Recovery from the failure is slow
C	Eliminate useless checkpoint	Processes are forced to take additional checkpoint to advance the global recovery line
D	Improve efficient	Incorrect replay of messages can cause orphan

Comparison between Checkpointing protocols here A-Coordinated checkpointing; B-Un Coordinated checkpointing; C-Communicaion induced checkpointing; D-Message Logging based checkpointing

4. Performance Analysis of Distributed Checkpointing Protocols

Sync-and-Stop (SaS) is a coordinated checkpointing protocol [1]. It was shown in [2] that $SaS \in 1$ -rollback. In this protocol there are no forced checkpoints, therefore, $F(SaS) = 0$. Regarding the control overhead, in each phase of SaS, the coordinator broadcasts three messages and the other $n - 1$ processes send two reply messages. Notice that the protocol needs an 8-bit control messages. Therefore, $M(SaS) = 5(n - 1)(wm + 8 \cdot wb)$. Chandy-Lamport (C-L) [7] is a coordinated checkpointing protocol in which there is no need to block the application execution. C-L belongs to 1-rollback and since there are no forced checkpoints, $F(C-L) = 0$. In a fully connected network with n nodes, C-L generates $2n(n - 1)$ messages per checkpoint [1] and the marker since is 8-bit, where it should distinguish between different runs of C-L. Therefore, $M(C-L) = 2n(n - 1)(wm + 8 \cdot wb)$. Fixed-Dependency-Interval (FDI) was suggested in [19]. Wang showed that FDI is Z-path free (ZPF) [3]. By [1], $ZPF \subseteq 1$ -rollback. Therefore, $FDI \in 1$ -rollback. Also, the dependency vector is piggybacked on each message. Thus,

$M(\text{FDI}) = n \cdot \text{MR}(\text{E})$ for an execution E. However, the number of forced checkpoints clearly depends on the number of processes 50 60 70 communication induced checkpointing protocol ensuring ZCF by preventing potential Z-cycles from being created. By [2], BQC \in n-rollback. Moreover, Alvisi et al [4] showed that BQC is worse than BCS but $F(\text{BQC}) = 2$. Lastly, the protocol propagates n^2 32-bit values on each application message to help processes detect suspected Z-cycles. Therefore, we have that $M(\text{BQC}) = \text{MR}(\text{E})(32 \cdot n^2 \cdot w_b + \varphi)$, where φ is the delay for intercepting every data message. d-Bounded Cycles (d-BC) is a communication induced checkpointing protocol that allows bounded cycles to be formed [2]. By [2], d-BC belongs to $(n - 1)d$ -rollback. Upon a new checkpoint $C_{p,i}$, process p broadcasts a cut of size no more than $d \cdot n$, therefore, $M(\text{d-BC}) = n \cdot w_m + d \cdot n^2 \cdot w_b$. Moreover, d-BC forces checkpoints by calling C-L only if a cycle of size d is generated. Since a Z-cycle is a special case of a cycle, then the conditions of generating cycles and Z-cycles are almost equivalent. Also since ZCF = 1-BC, then by [4] we have that $F(1\text{-BC}) = 2$.

5. Related Work

There has been much work on checkpointing performance analysis [11, 12, 15, 17]. Most of these works do not take into account the rollback propagation. Ours is the first to incorporate all parameters that affect the performance in distributed environments into an analytical measure. Mishra and Wang [11] evaluated several checkpointing protocols by implementing and running them with test applications. Ziv and Bruck [17] compared four checkpointing protocols by using the Markov Reward Model [13]. Our approach differs from [17] in that we provide a technique for comparing any checkpointing protocol based on rollback propagation. Ziv and Bruck presented in [18] a checkpoint scheme for duplex systems, and conducted a performance analysis for their scheme in the duplex system. However, it is not a general system for distributed executions. Vaidya defined the overhead ratio for uniprocessor systems as a function of the checkpoint overhead and latency [14], and proved that the optimum checkpoint interval depends on α . Additionally, he claimed that the overhead ratio can be computed for distributed systems as in uniprocessor systems by taking the values of parameters either to be the maximum or the average over all processes. In [16], Vaidya computed the overhead ratio for the two-level recovery

approach. This approach tolerates single failures with a low overhead and multiple failures with a higher overhead. Plank and Thomason [14] presented a method for estimating the overhead ratio for coordinated checkpointing. By assuming coordinated checkpointing, they do not care about rollback propagation. Moreover, they do not address the control overhead incurred by control information.

6. Conclusion

We have reviewed some fundamental concepts of checkpointing protocols in distributed systems. This paper presents a comprehensive model of rollback recovery protocols that encompasses a wide range of checkpoint/restart protocols. Included coordinated checkpoint and uncoordinated checkpoint protocols. This model provides the first tool for a quantitative assessment of all these protocols. Hence the concept of checkpoint is introduced before planned disconnection so that checkpointing can be completed without any delay resulting enhanced fault tolerance in the proposed scheme.

References

- 1) J. S. Plank. Efficient Checkpointing on MIMD Architectures. PhD thesis, Princeton University, January 1993.
- 2) A. Agbaria, H. Attiya, R. Friedman, and R. Vitenberg. Quantifying Rollback Propagation in Distributed Checkpointing. In 20th Symposium on Reliable Distributed Systems, pages 36–45, New Orleans, October 2001.
- 3) Y. M. Wang. Consistent Global Checkpoints that Contain a Given Set of Checkpoints. IEEE Transactions on Computers, 42(4):456–486, April 1997.
- 4) L. Alvisi, E. Elnozahy, S. Rao, S. A. Husain, and A. D. Mel. An Analysis of Communication Induced Checkpointing. In Proceedings of the 29th Fault-Tolerance Computing Symposium, pages 242–249, Madison, Wisconsin, June 1999.
- 5) G.H. Forman and J. Zahorjan, The changes of Mobile computing, computer pp 38-47, Apr-1994
- 6) Ms. Pooja Sharma and Dr. Ajay khuntala " A survey of checkpointing Algorithm in Mobile Ad Hoc Network" Globl Journal of Computer Science and Technology 2012. [7] Sarmistha Neogy, Anupam siha, pradip k Das ,CCMUL: A Checkpointing protocol for distributed system processes, IEEE, 2004.
- 7) B. bhargava, S.R. Lian "Independent checkpointing and concurrent rollback for recovery in distributed systems- An Optimistic approach". proc 7th IEEE Symp. Reliable Distributed syst. pp 3-12 1988 oct.
- 8) L. Alvisi, E.N. Elnozahy, S. Rao, S. A. Husain and A. Del Mel. "An analysis of communication-induced checkpointing." In Proceedings of the Twenty Ninth International Symposium on Fault-Tolerant Computing, Jun. 1999.
- 9) D.B. Johnson. "Distributed system fault tolerance using message logging and checkpointing." Rice University, Dec. 1989.

- 10) S. Mishra and D. Wang. Choosing an Appropriate Checkpointing and Rollback Recovery Algorithm for LongRunning Parallel and Distributed Applications. In 11th ISCA International Conference on Computers and their Applications, San Francisco, CA, March 1996
- 11) J. S. Plank and M. G. Thomason. Processor allocation and checkpoint interval selection in cluster computing systems. *Journal of Parallel and Distributed Computing*, 61(11):1570–1590, November 2001.
- 12) K. S. Trivedi. *Probability and Statistics with Reliability, Queuing, and Computer Science Applications*. Prentice-Hall, USA, 1982.
- 13) Nitin Vaidya. On Checkpoint Latency. In Pacific Rim International Symposium on Fault-Tolerant Systems, Newport Beach, December 1995.
- 14) Nitin H. Vaidya. Another Two-Level Failure Recovery Scheme: Performance Impact of Checkpoint Placement and Checkpoint Latency. Technical Report TR94-068, Dept. of Computer Science, Texas A&M University, 1994.
- 15) Y. M. Wang. Consistent Global Checkpoints that Contain a Given Set of Checkpoints. *IEEE Transactions on Computers*, 42(4):456–486, April 1997.
- 16) Ziv and J. Bruck. Analysis of Checkpointing Schemes for Multiprocessor Systems. In Proceeding of the 13th Symposium on Reliable Distributed Systems, pages 52–61, 1994.
- 17) Ziv and J. Bruck. Efficient checkpointing over local area network. In IEEE Workshop on Fault-Tolerant Parallel and Distributed Systems, June 1994.
- 18) Ch.D.V.Subba Rao and MM Naidu : A new efficient coordinated checkpointing protocol combined with selective sender based message logging , IEEE,2008.
- 19) Acharya and B.R.Badrinath ,checkpointing distributed Applications on Mobil computers,proc.3rd Int'l conf.parallel and distributed Information systems, sept.1994.
- 20) R.Prakash and M.Singhal, "Low-cost checkpointing and failure recovery in mobile computing systems," IEEE Trans.parallel and distributed systems pp.1035-1048,oct 1996
- 21) Lalit kumar p.kumar "A synchronous checkpointing protocol for mobile distributed systems: probabilistic approach" *Int.Journal of information and computer society* 2007.

Author Profile



Krishnamurthy Ramasubramanian
Asst.Prof., Department of
CSE,Samskruti College of Engineering
and Technology,Kondapur (V),
Ghatkesar (Mdl), Hyderabad