# High Performance Mobile Checkpointing Algorithm

Surender Kumar*
Dept. of Information Technology (IT),
Haryana College of Tech. & Mgmt.(HCTM)
Ambala Road, Kaithal(HRY), INDIA
ssjangra20@rediffmail.com

R.K. Chauhan
Dept. of Computer Sc. & App.(DCA)
Kurukshetra University,
Kurukshetra(HRY), INDIA

Parveen Kumar
Deptt. of Computer Sc. & Engg.
Meerut Institute of Engg. & Tech.(MIET)
Meerut(U.P), INDIA

*Abstract:* A mobile computing system consists of mobile and stationary nodes, connected to each other by communication network. This system raises several requirement and restrictions, such as limited battery life, mobility, disconnection of hosts and lack of stable storage not taken into account by conventional distributed systems. Our checkpointing algorithm forces a minimum number of processes to take checkpoint and reduces the delay incurred in determining the consistent global state for mobile systems in the comparison of the traditional coordinated checkpointing algorithms. By reducing the checkpoint cost like searching, saving checkpoints, and transmitting data through wireless link we can deal with different checkpointing issues for mobile systems.

*Keywords:* Mobile Systems, Coordinated checkpointing, Consistent Checkpoints, Global Snapshot, , Non-blocking, orphan message.

## I. INTRODUCTION

A mobile computing system is a distributed system where some of the nodes are mobile computers (Mobile Hosts (MHs)) [9]. As time passes mobile computers location gets change. To communicate with MHs, mobile support stations (MSSs) are added. An MSS communicate with other MSS by wired networks and with MHs with wireless network. Each of saved state is called snapshot (checkpoint). All the processes in the system take their checkpoints periodically.

The checkpointing techniques do not require user interaction and can be classified into following categories: (a) Uncoordinated checkpointing (b) Coordinated checkpointing (c) Quasi-Synchronous (d) Message – Login based checkpointing [14]. In this paper we concentrate on coordinated checkpointing technique which maintains a consistent snapshot of system all the times. A consistent global snapshot indicates set of N local snapshots (checkpoints) one from each process forming a consistent system state which can be used to restart process execution upon a failure.

It is desirable to minimize the amount of lost work by restoring the system to most recent consistent global checkpoint. A good snapshot collection algorithm should be Non-Blocking i.e. which does not force the nodes in the system to stop their computations during snapshot collection. An efficient algorithm keeps minimum effort required for collecting a consistent snapshot to a minimum. The snapshot collection algorithm by Chandy and Lamport forces every node to take its local snapshots but the computation is allowed to continue while the global snapshot is being collected [1]. In Koo and Toueg's algorithm all the nodes are not forced to take their local snapshots [7]. However, the underlying computation is suspended during snapshot collection.

We propose a new coordinated checkpoint algorithm which is non- blocking and efficient that forces a minimal set of nodes to take their snapshot and underlying computation is not suspended during snapshot collection.

The rest of the paper is organized as follows: Section 2 presents the system model. Section 3 presents the basic idea, minimal dependency information at each node and the node mobility management of our algorithm. Section 4 presents a Non-Blocking Coordinated checkpoint algorithm with data structure used and an example describing our algorithm. Section 5 presents the related work on which our algorithm is based. Finally, Section 6 presents conclusion.

## II. SYSTEM MODEL

A message passing system consists of N fixed number of nodes that communicate each other only through messages. Some of the nodes may change their location with time. They are referred to as mobile hosts or MHs [9]. Static Nodes are referred to as mobile support stations or MSSs are connected to each other by a static network [Fig. 1]. An MH can be directly connected to at most one MSS at any given time and can communicate with other MHs and MSSs only through the MSS to which it is directly connected. The system does not have any shared memory or a global clock. Hence all the communication and synchronization takes place through messages.
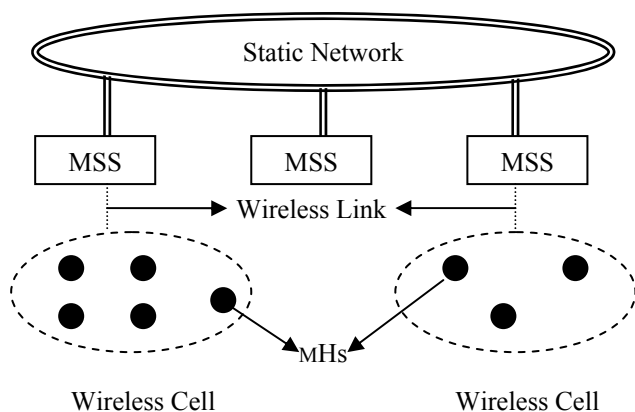
*Figure 1. Mobile Distributed System*

The messages generated by underlying distributed application will be referred to as computation messages. Messages generated by the nodes to advance checkpoints, handle failures and for recovery will be referred to as system messages. In this paper the horizontal lines extending towards right hand side represent the execution of each process (MH) and arrows between them represent the messages. Processes have access to a stable storage device that survives failures. The number of tolerated process failures may vary from 1 to N [14].

### III.   RELATED WORK

In Chandy-Lamport algorithm [1] control messages are sent to all the nodes for consistent global checkpoint. Hence message send overhead is increased along all the channels of network.

Acharya-Badrinath algorithm [9] proposed an uncoordinated checkpointing algorithm for mobile distributed system because they found the limitations of high cost to receive request messages along every channel in network and absence of local checkpoint of MH during disconnect interval in coordinated checkpoint algorithm.

In Koo-Toueg Algorithm [7] the underlying computation is blocked. There is direct dependency approach is used while global snapshot collection. Such algorithm is not suitable for concurrent initiation.

In Venkatessan and Juang's optimistic failure recovery algorithm [15] no dependency information is send with the computation messages. Hence while recovery process too many rollbacks occur.

In [3] Guohong Cao and Mukesh Singhal had proposed an efficient algorithm that neither forces all the processes to take checkpoints nor blocks the underlying computation during checkpointing and which significantly reduces the number of checkpoints. In this paper it is described that there does not exist a non-blocking algorithm that forces only a minimum number of processes to take checkpoints. Their algorithm requires minimum number of processes to take tentative checkpoints and thus minimizes the workload on stable storage server. Their algorithm has three kinds of checkpoints: tentative, permanent and forced. Tentative and permanent checkpoints are saved on stable storage. Forced checkpoints do not need to be saved on stable storage. They can be saved on any where even in the main memory.

In [4] Guohong Cao and Mukesh Singhal had introduced the concept of "Mutable Checkpoint" which is neither a tentative checkpoint nor a permanent checkpoint to design efficient checkpointing algorithms for mobile computing system. Mutable Checkpoint can be saved anywhere e.g. the main memory or local disk of MHs. Taking a mutable checkpoint avoids the overhead of transferring large amount of data to stable storage at MSSs over the wireless network.

### IV.   MANAGING CHECKPOINITNG ISSUES

Checkpointing cost in mobile distributed systems have three components: (1) Searching of MHs due to the mobile nature of MHs (2) To save the checkpoints due to the lack of stable storage on MH. (2) Transmit and receive the checkpointing from MH to MSS vice-versa, due to the wireless connectivity between both. Taking into consideration of different checkpoint issues and reducing the checkpoint cost, our new checkpoint algorithm mainly focuses on following:

#### A.   *Low Power Consumption at MHs Level:*

Our proposed checkpointing algorithm uses the following tricks to save the battery power of an MH:

- Due to the limited power and computation ability of MHs, and unlimited power and computation ability of an MSS, in our approach all the dependency information and other MHs related information are maintained in the corresponding MSS instead of itself.
- The checkpoint request are sent to the process in minimum set (which are directly or transitively depended on the initiator) to reduces the power consumption.

#### B.   *Efficiently use of Limited Available Memory of MHs*

As the most of the data are stored on the stable storage of the MSS and only needful information are stored on to the MHs.

#### C.   *Efficiently use of Limited Wireless Bandwidths*

To optimize the wireless bandwidth we use the concept of soft checkpoint, these soft checkpoint are stored on the local memory of MH and stored only if the local checkpoint time expires or any untimely event occurs.

#### D.   *Reducing Searching Cost:*

Corresponding MSS keeps the record message received, sent, location and checkpoint record. By maintaining the location information of an MH on MSS, we reduce the searching cost and shorten the checkpointing latency.

#### E.   *Handling Mobility*

As the mobile nature of an MH, an MHs may not stay forever in only one MSS and local information may become out of date without proper handling of handoffs. In our approach to handle the handoffs each MSS maintains a location list of MHs that join and leave that MSS. When MHi depart from a cell and joins with another cell, it first sends a DEPART () message to the old MSS, and then established a new connection by sending JOIN message by attaching the ID of own MH and old MSS. Once a MH join with the new MSS, then new MSS receive the old data related to the new joined MH by sending message OLDDATA (MHid, old MSSid). In such

way old MSS sends all the details to the new MSS and update the LOCATION information of the MH on its own level. New MSS store all the necessary on its own storage.

### F. *Handling Disconnection*

Disconnection of an MH is a voluntary operation and it may take arbitrary period of time. At the time of disconnection from $MSS_1$:

- MH takes its local checkpoint which is stored at $MSS_1$ as disconnect_checkpoint$_i$ which serves request messages for MH to take checkpoint
- Stores its dependency vector Di at MSS1.
- The Computation messages, for MH arriving at $MSS_1$ during disconnect interval are stored at $MSS_1$ until the end of the interval.

.
### Reconnection

At the time of reconnection to $MSS_2$: MH executes a reconnection algorithm. The reconnection algorithm sends a message through $MSS_2$ to $MSS_1$. On receiving the message $MSS_1$ executes the following steps:

- If $MSS_1$ had processed request message for MH then disconnect_checkpoint$_i$ and the buffered messages are sent to MH.
- If no checkpoint request for MH was received by $MSS_1$ during disconnect interval only buffered messages are sent.
- After that $MSS_1$ removes the buffered messages, disconnect_checkpoint$_i$ and MH's dependency vector.

When the data sent by $MSS_1$ arrives at MH, MH executes the following actions:

- If the received data contains disconnect_checkpoint$_i$, MH stores this checkpoint as its local checkpoint and resets all except the ith component of dependency vector $R_i$ before processing the messages.
- Process all the received buffered messages.
- The dependency vector is updated.

Now this reconnect algorithm ends and MH makes normal communication.MSS1 removes the disconnect_checkpointi at the end of disconnect interval. In such a way mobility and disconnection of MH get managed.

### V. THE CHECKPOINTING ALGORITHM

In this section, we present a mobility based efficient checkpoint algorithm for mobile distributed systems. The algorithm forces a minimum set of nodes to take local checkpoints. Thus overhead of checkpoint collection get minimized. After the coordinated snapshot collection terminates, the nodes that did not participate in snapshot collection can take their local checkpoints in lazy phase approach. When a node initiates a request for snapshot collection to another node then that node takes its local snapshot and propagating the request to neighbouring nodes. A global snapshot is collection of all the local nodes which participates for snapshot initiation. The snapshot thus generated is latest than each of the snapshot thus collected independently. Thus amount of lost work during rollback, after the node failure is minimized. The underlying computation need not have to be suspended during snapshot collection.

### A. *Maintain Minimal Dependency Information:*

For maintaining dependency information we here use the same approach as in [2]. The dependency is created by means of messages between nodes. Node Pi maintains a Boolean vector $D_i$ of n components. At $P_i$, the vector initialized as follows:



(a)                                                            (b)

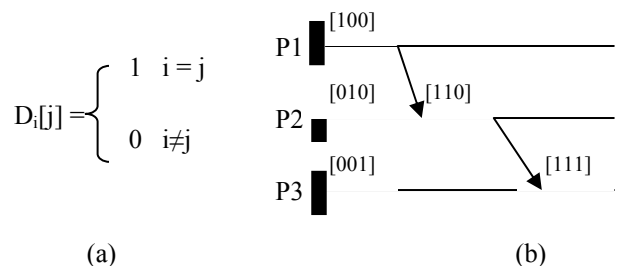*Figure 2(a, b) Dependency vector forming equation 2(b) Propagation of dependency information via D[i] vector.*

When a node Pi sends a message to Pj it then modifies vector Di. This informs Pj about the nodes that have affected Pi. While processing a message M Pj extracts Boolean vector M.D from the message and uses it to update Dj as follows: $Dj[k] \leftarrow Dj[k]$ *OR* $M.D[k]$, where $1 \leq k \leq n$.
Fig. 2 (b) shows the dependency information through messages: Since P2 was dependent on P1 before sending M2 to P3; P3 becomes transitively dependent on P1 on receiving M2.
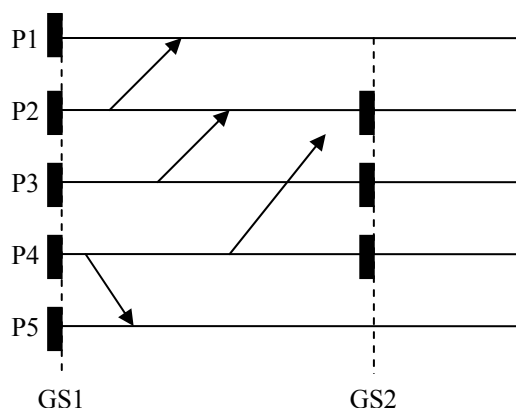


GS1                                    GS2

*Figure 3.  Minimum no. of checkpoints and Consistent Global State (CGS) or consistent cut.*

Fig.3 shows the vertical line $G_1$ the global checkpoint at the beginning of the computation. Let Process $P_2$ initiates a new snapshot collection.  only $P_3$ and $P_4$ need to take their local snapshot because they depends upon node $P_2$ .But the nodes $P_1$ and $P_5$ need not take their snapshot because they do not have dependencies on to process P2.Dotted line $G_2$ shows the global recovery line or current global checkpoint.

There must be a great need to avoid the orphan message to maintain the global consistent state checkpoint in global checkpoints collection.  There are three processes P, Q and R. Let Q initiates checkpoint request to processes P and R. Let P and R take their local checkpoints. If R sends message M to P before receiving checkpoint request then message M will become an orphan message which creates a problem during snapshot collection. To avoid such problem

a concept of checkpoint sequence number get arise. We call this ckpt_num in our algorithm.

### B. Data Structure

*Di* : a Boolean vector Di of n components. At Pi, the vector initialized as follows as in [2]:
Di[i] = 1; Di[j] = 0 if i ≠ j; When a node Pi sends a message to Pj it then changes vector Di . This tells Pj about the nodes that are dependent on Pi. While processing a message M Pj extracts Boolean vector M.R from the message and uses it to update Rj as follows: Dj[k] ← Dj[k] OR M.D[k], where 1 ≤ k ≤ n.

*ckpt_num*: when the node takes its local checkpoint then this integer number is increased.

*Weight*: A nonnegative real variable with maximum value 1 used to detect the termination of snapshot collection or checkpointing algorithm.

*Transmit:* a Boolean array of size n maintained by each node in its stable storage. This array is initialized to all zeros. It is used to keep the trail of those nodes to which checkpoint requests were sent by node. If in this array each element has all 0s then response message is sent to the snapshot initiator with a weight equal to weight received in the request. If in this array some elements are put to 1 then for all i such that transmit[i] = 1, a request is sent to Pi with a non zero segment of weight received in request message and rest part of weight is sent to initiator with a response message.

*Trigger:* A set of 2-tuples (init_id, init_ckptnum) maintained by each node, where init_id indicates the identifier of checkpointing initiator. init_ckptnum shows the checkpoint number of the initiator node when it took its own local snapshot on initiating the snapshot collection. trigger is changed for all system messages and the first computation message that a node sends to every other node after taking a local snapshot.

*ckpt_array:* This is an array of n integer maintained at each node, ckpt_array[i] indicates the ckpt_num of the next message expected from node $P_i$.

*self_trigger:* The trigger tuple of a node receiving computation message

*msg_trigger:* Trigger tuple of computation message

*get_weight:* The weight received by dependent nodes

*forward_weight:* The weight sends by the node which further spread checkpoint request

### C. The Algoritm

To minimize the number of MHs taking checkpoints, the checkpoint initiator has to identigy which MHs must take a checkpoint. The detailed checkpointing algorithm is described as follows:

*(a)Checkpoint initiation process:* If the checkpoint initiator process, say $P_{in}$ runs on an MH, it first takes soft checkpoint and sends the checkpoint initiation request to its local MSS, say $MSS_{ini}$ and then $MSS_{ini}$ become the checkpoint proxy. This proxy MSS handle all the computation and location information regarding the MH and act as follows:
(1) Increments its ckpt_num (2) initialized weight to 1 (3) It sets init_id and init_ckptnum in its trigger tuple (4) It sends checkpoint request message to all its dependent nodes. The request message now includes: weight, initiator's trigger and dependency vector $D_i$.

*(b)Response of a node receiving of checkpoint request:*
Let node $P_i$ receives checkpoint request

if (reqst_msg.trigger ≠ $P_i$.trigger) then
  {
    $P_i$ takes soft checkpoint;
    $P_i$ propagates reqst_msg to all the dependent nodes but not $M.D_i$; /* Explained Checkpoint module */
    Send portion of the received weight with its reqst_msg;
    Update initiator trigger tuple;
    Send response_msg to the initiator;
  }
else if ( transmit[i] = 0 ) then
  {
    $P_i$ send response_msg with weight received with reqst_msg to initiator;
  }
 else
  {
    $P_i$ send reqst_msg to nodes for which transmit[j] =1 with portion of weight;
    $P_i$ sends response_msg with remaining weight to initiator;
  }
}

*(c)Response of a node receiving of computational message:* Let a node Pj receives a computation message M from other node $P_i$ then following action occurs:
If (ckpt_num$_i$ ≤ ckpt_array[i]) then
  {
    Pj will not take any checkpoint;
    Restart the computation by processing message M;
  }
else
  {
  Set   ckpt_array[i] = ckpt_num$_i$ ;
/* $P_i$ has already taken a checkpoint before sending M and this is the first computation message sent from $P_i$ to Pj.
  So M carries a trigger (init_id, init_ckptnum)*/
  }

 if (msg_trigger = self_trigger) then
   Update dependency vector Di[j];
   // Pi and Pj has taken checkpoints w.r.t same initiator
else (msg_trigger.pid ≠ self_trigger.pid)
  {
   If (Pj had processed a message from node $P_k$) then
     Return;
     Pj takes tentative checkpoint;
     Set msg_trigger=self_trigger;
     Propagate snapshot request to dependent processes;
  }
 }

*(d)Further Checkpoint request propagation:* Let node $P_i$ take checkpoint request. It propagates checkpoint request to its dependent processes as follows:
Take soft checkpoint;
Update ckpt_num$_i$ and transmit[i] ;
 Self_trigger=msg_trigger;
 transmit[i] = R$_i$ − M.R;
 for all k dependent nodes set transmit[k]=1
 {
 get_weight = get_weight/2
 forward_weight  = get_weight;
// Send following module to dependent nodes //

send(P$_i$ ,request_msg,chkpt_num,self_trigger,forward_weight);
}

(e) ***Termination of Checkpointing:*** After receiving the checkpoint request message, each MH sends a reply to the MSSini along with weight. When the initiator receives reply from all the relevant MHs means the sum of received weight becomes equal to 1. it make a decision whether it commit or abort to all the relevant processes. This is same as any other two –phase checkpointing strategies. On receiving COMMIT or ABORT message from the initiator MSS the process P$_i$ act as follows:

*On receiving COMMIT()*
If(Soft$_i$)
{
Discard old permanent checkpoint, if any;
Convert the tentative checkpoint into permanent one;
}
On receiving ABORT()
{
if (soft$_i$ )
Discard the soft checkpoint;
}
If a process receives commit message, if it taken soft checkpoint already   The previous permanent local checkpoints at these nodes are discarded. Now if further recovery is required the nodes will rollback to current checkpoint.

### D.  Working Example

Following example (Fig. 4) clarifies the concepts used in our algorithm with the help of Fig node P$_3$ initiates snapshot collection by taking its local checkpoint. The node P$_2$ and P$_4$ shows dependencies to P$_3$. The broken arrows shows request messages sent to P$_2$ and P$_4$ to take their snapshots on their timeline. P$_4$ takes first snapshot and then sends a message M3 to P$_2$. When M3 reaches P$_2$ it is the first message reached at P$_2$ such that msg_trigger.pid $\neq$ self_trigger.pid. Hence P$_2$ takes its snapshot before processing M3. Node P$_1$ takes its local independent snapshot before sending a message M4 to P$_2$. The interval number of M4 is greater than the value expected by P$_2$ from P$_1$.
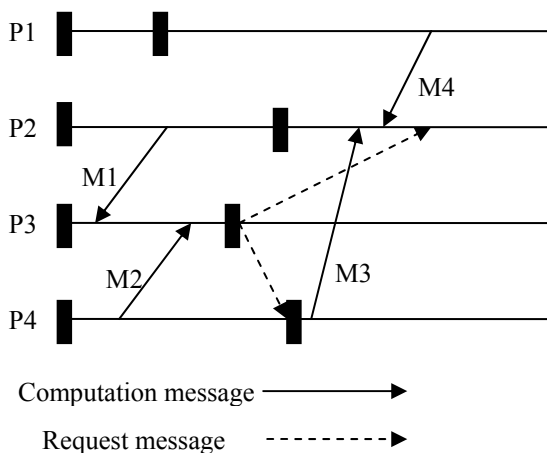


*Figure 4. Global checkpoint collection*

But when M4 reaches P$_2$ it is not the first computation message received by P$_2$ with a higher interval number than expected whose msg_trigger.pid is different from P$_2$'s self_trigger.pid. So a snapshot is not taken because it will create inconsistency: The reception of M3 will be recorded if P$_2$ takes a snapshot just before it processes M4, but the transmission of M3 will not have been recorded by P$_4$ and now M3 becomes orphan. Also here msg_trigger of M3 = self_trigger of request message to P2 . Hence no need to take further checkpoint.

## VI.   RESULT AND DISCUSSION

In this algorithm we have described a efficient coordinated checkpointing algorithm to survive a failure in mobile distributed system which is able to save consistent global states with small overheads by reducing the searching cost as well as power and bandwidth requirements. The main feature of our algorithm are: (1) it is non-bocking ;(2) it is orphan free as it takes checkpointing decision on the basis of checkpoint sequence number; (3) it is less power consuming as it is forces to take minimum processes which are directly and transitively dependent on the initiator process and all the information related to the MHs are handled by there corresponding MSS; (4) it utilize the wireless bandwidth efficiently as the only minimum required information are transfer through between MHs and MSSs. (5) it minimize the coordinated message as MSS maintain a minimal dependency information related on its own level.

## VII.   REFERENCES

[1]  K.M. Chandy and L.Lamport. "Distributed Snapshots: Determining Global States of Distributed Systems" ACM Transactions Computer systems vol. 3, no.1.pp.63-75, Feb.1985.

[2]  Prakash R. and Singhal M., "Low-Cost Checkpointing and Failure Recovery in Mobile Computing Systems" ,IEEE Transaction On Parallel and Distributed Systems, vol. 7, no. 10, pp. 1035-1048, October1996

[3]  Guohong Cao and Mukesh Singhal, "On Coordinated Checkpointing in Distributed Systems" IEEE Transaction On Parallel and Distributed Systems, vol. 9, no. 12, pp. 1213-1224, December 1998.

[4]  Guohong Cao and Mukesh Singhal, "Mutable Checkpoints: A New Checkpointing Approach for Mobile Computing Systems", IEEE Transaction On Parallel and Distributed Systems, vol. 12, no. 2, pp. 157-171, February 2001.

[5]  Weigang Ni, Susan V. Vrbsky and Sibabrata Ray "Pitfalls in Distributed Non blocking Checkpointing", University of Alabama

[6]  Prakash R. and Singhal M. "Maximal Global Snapshot with concurrent initiators," Proc. Sixth IEEE Symp. Parallel and Distributed Processing, pp.344-351, Oct.1994.

[7]  Koo. R. and S.Toueg. "Checkpointing and Rollback-Recovery for Distributed Systems" .IEEE Transactions on Software Engineering, SE-13(1):23-31, January 1987.

[8]  Bidyut Gupta, S.Rahimi and Z.Lui. "A New High Performance Checkpointing Approach for Mobile Computing Systems". IJCSNS International Journal of Computer Science and Network Security, Vol.6 No.5B, May 2006.

[9] Acharya A. and Badrinath B. R., "Checkpointing Distributed Applications on Mobile Computers," Proceedings of the 3rd International Conference on Parallel and Distributed Information Systems, pp. 73-80, September 1994.

[10] Ch.D.V. Subba Rao and M.M.Naidu. "A New, Efficient Coordinated Checkpointing Algorithm Combined with Selective Sender-Based Message Logging"

[11] Nuno Neves and W. Kent Fuchs. "Adaptive Recovery for Mobile Environments",in Proc.IEEE High-Assurance Systems Engineering Workshop,October 21-22,1996,pp.134-141.

[12] Y.Manable. "A Distributed Consistent Global Checkpoint Algorithm With minimum number of Checkpoints". Technical Report of IEICE, COMP97-6(April1997)

[13] J.L.Kim and T.Park. "An efficient algorithm for checkpointing recovery in Distributed Systems" IEEE Transaction On Parallel and Distributed Systems,4(8):pp.955-960, Aug 1993.

[14] Elnozahy E.N., Alvisi L., Wang Y.M. and Johnson D.B., "A Survey of Rollback-Recovery Algorithms in Message-Passing Systems," ACM Computing Surveys, vol. 34, no. 3, pp. 375-408, 2002.

[15] S.Venkatesan and T.T.-Y.Juang , " Low Overhead Optimistic Crash Recovery:", Preliminary version appears in Proc. 11th Int'l Conf. Distributed Computing Systems as "Crash Recovery with Little Overhead,"pp.454-461,1991.