

Unit-4 - check pointing & Rollback

Recovery & Consensus.

In Distributed System check pointing & Rollback recovery are used for fault Tolerance.

→ Fault tolerance increases the performance of a System in terms of Availability & Reliability.

Fault tolerance is an ability or capability of a System to continue its operation without interruption, when one or more components of the System fail.

Rollback Recovery - Restore the System back to a consistent state after a failure. It achieve fault tolerance by periodically saving the state of a process during the failure free execution.

→ It treats distributed System as a collection of processes that communicate through a network.

- up on failure, a failed process restarts from one of its Saved states. which reduces the amount of lost computation.

→ The saved state of a process at an instance of a time is called check point

→ In distributed system, all processes save their local states at certain instances of time. This saved state is called local check point.

→ The contents of a check point depend upon the application content and the check pointing method.

→ Each process periodically saves its state on stable storage.

→ The saved state contains sufficient information to restart the process execution.

→ Global check point is a set of n local check points.

check points & Different types of messages. ①

The Saved states of a process at a particular instance of time is check point.

All processes save their local states at certain instants of time.

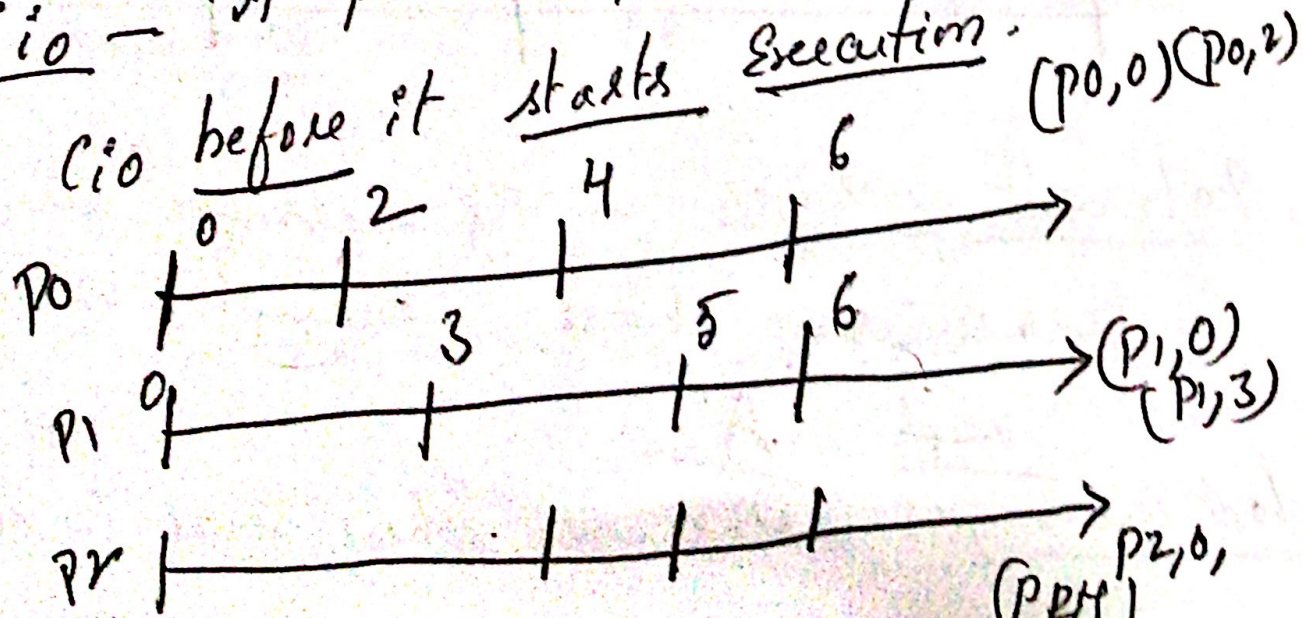
→ A local check point is a Snapshot of the state of the process at a given instance.

→ A process stores all local checkpoints on the stable storage.

→ A process is able to roll back to any of its existing local checkpoints.

C_{ik} - The k^{th} local check point at process P_i .

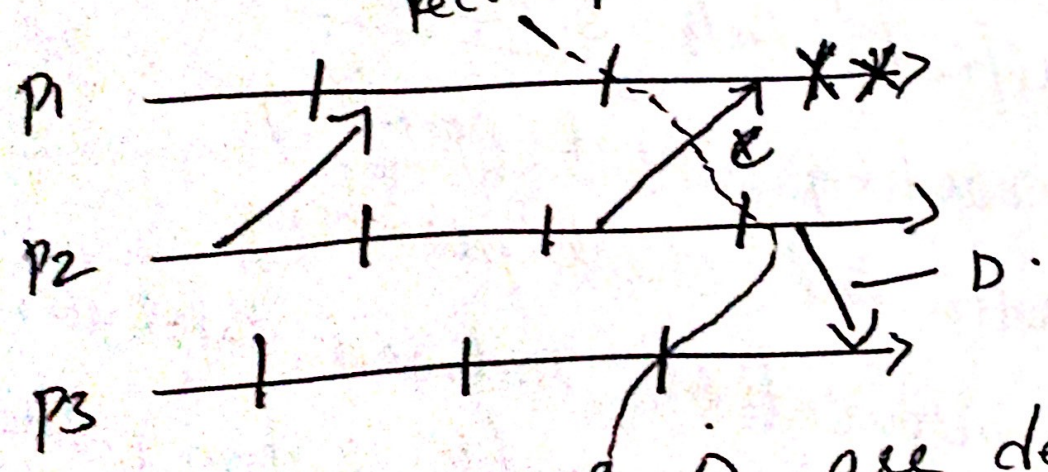
C_{i0} - A process P_i takes a checkpoint before it starts execution.



but receive is undone due to roll back ⁽³⁾ are called lost messages.

This type of messages occur when the process rolls back to a check point prior to the reception of the message. Message 'B' is an example.

Delayed messages. Messages whose receive is not recorded, because the receiving process was either down or messages arrived after the roll back.
recovery line.



messages c & d are delayed

orphan messages - Messages receive is recorded, but sent not recorded are called orphan messages.

Duplicate messages.

Duplicate messages arise due to message logging & replaying during process recovery.

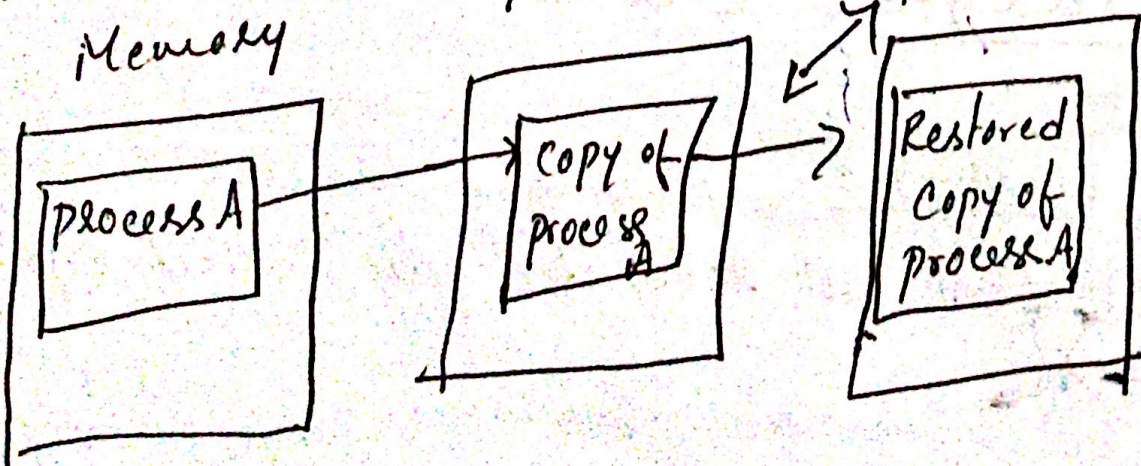
Check-point based Rollback Recovery.

upon a failure, check point based roll back recovery restores the system state to the most recent consistent set of check points.

check point based recovery techniques are classified into 3 categories. They are

- 1) Uncoordinated check pointing.
- 2) Coordinated check pointing.
- 3) Communication - Induced check pointing.

Non-volatile. System failure.



Any consistent global check point ^(H) can be used to re-start the process execution upon a failure.

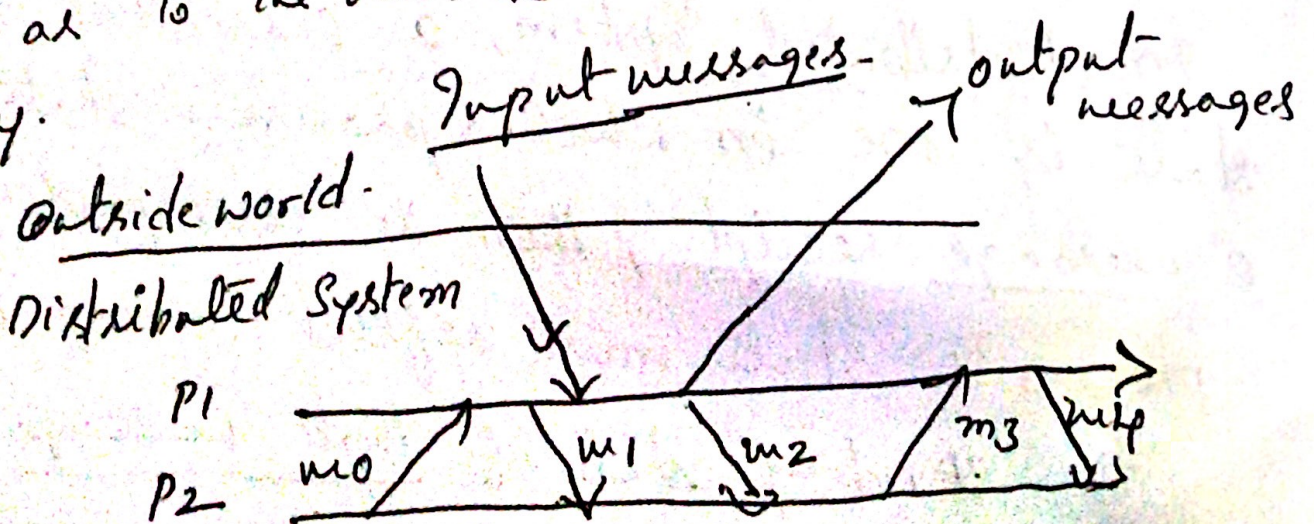
→ The most recent consistent global checkpoint is termed as the recovery line.

→ The fundamental goal of the recovery roll back protocol is to bring the system into a consistent state whenever inconsistencies occur because of a failure.

Background. System model : / issues in failures.

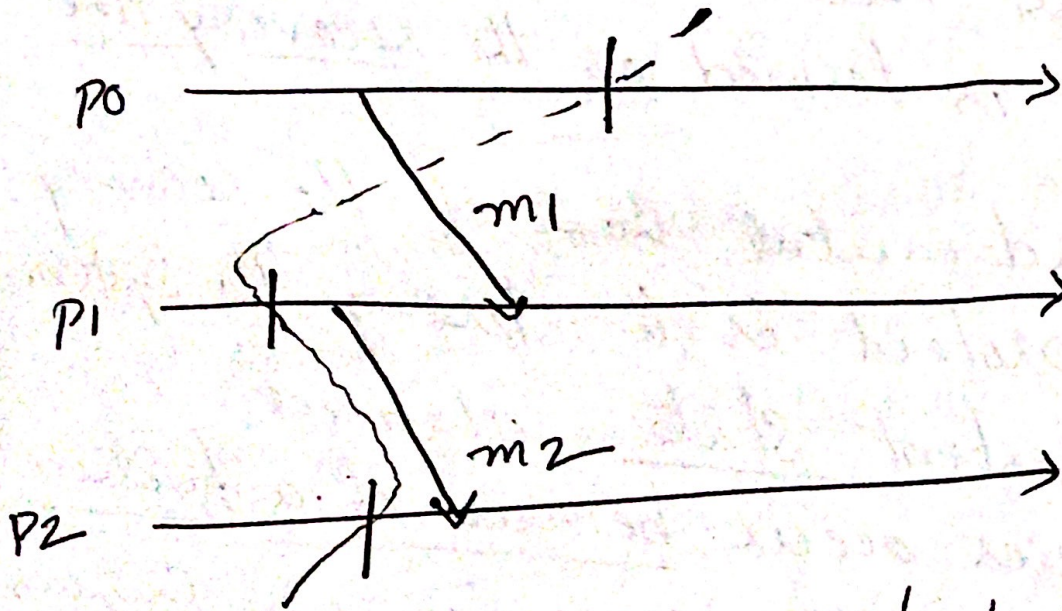
A distributed system consists of a fixed no. of processes P_1, P_2, \dots, P_N , which communicate through messages.

processes cooperate with each other as well as to the outside world through messages only.

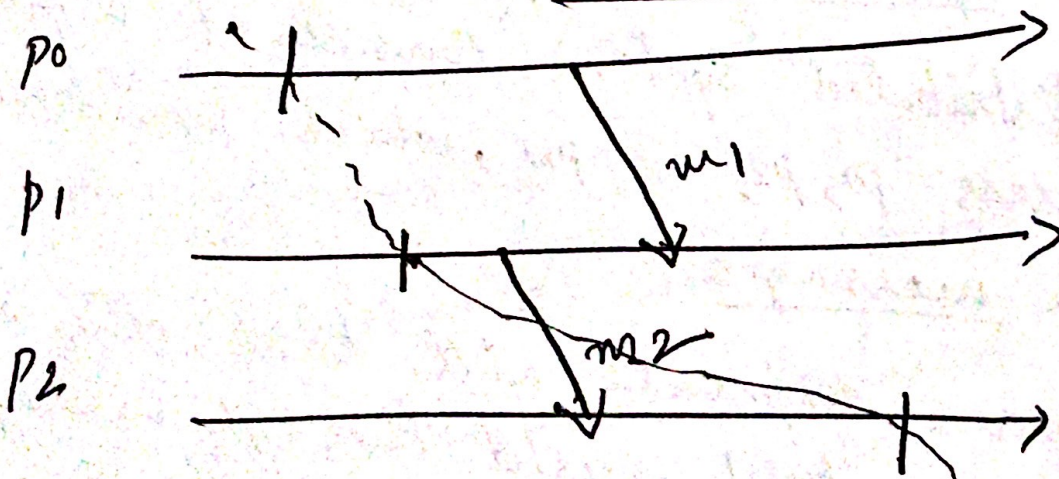


Always consistent state is needed for failure free execution in distributed system.

consistent state



inconsistent state



In distributed system, a consistent system state is one in which a process's state is a message receipt, then the corresponding sender process's state must reflect ~~corresp~~ the sending state of that message.

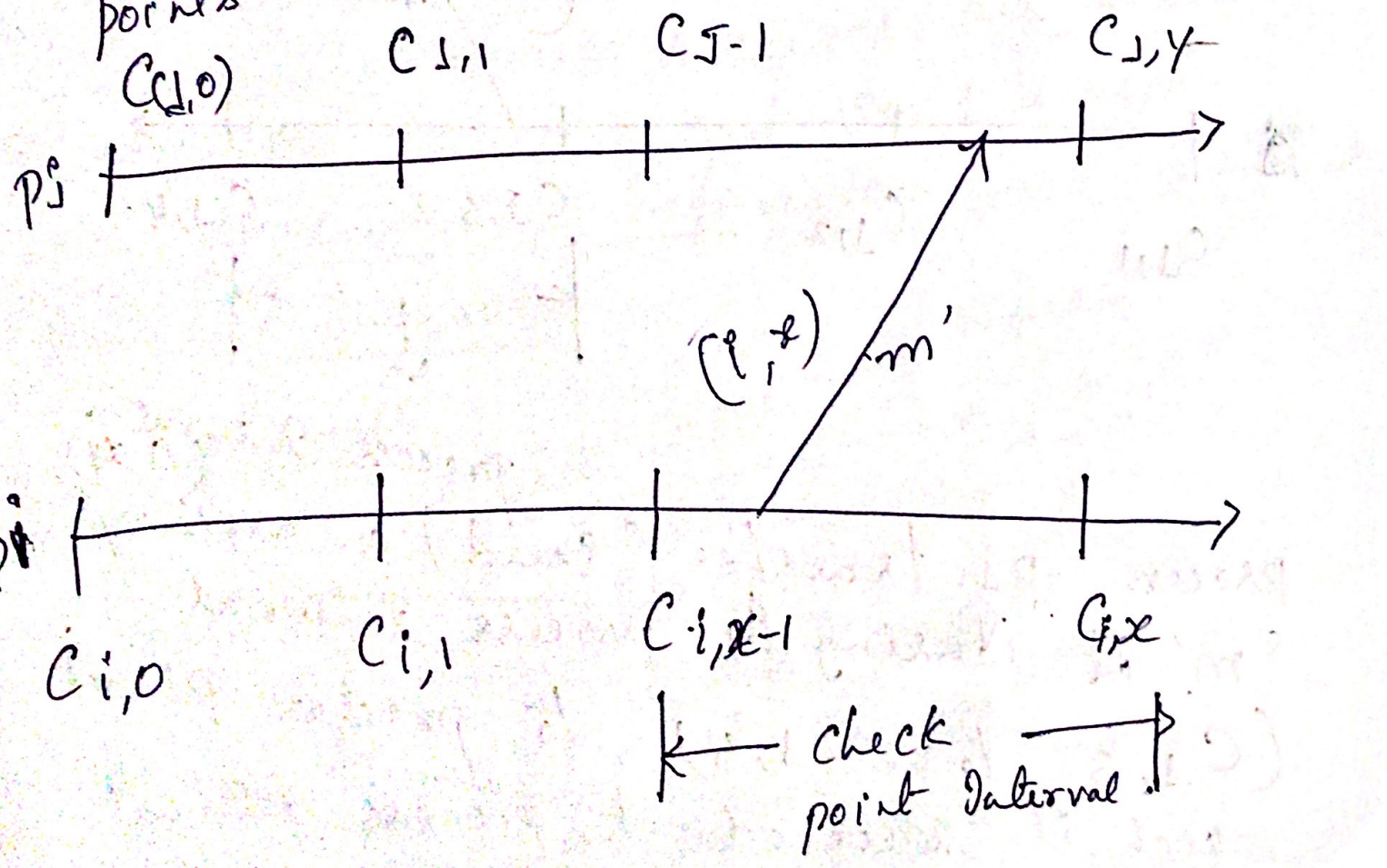
check point based Roll back Recovery (5)

Continuation from previous -

1. Un Coordinated check pointing.

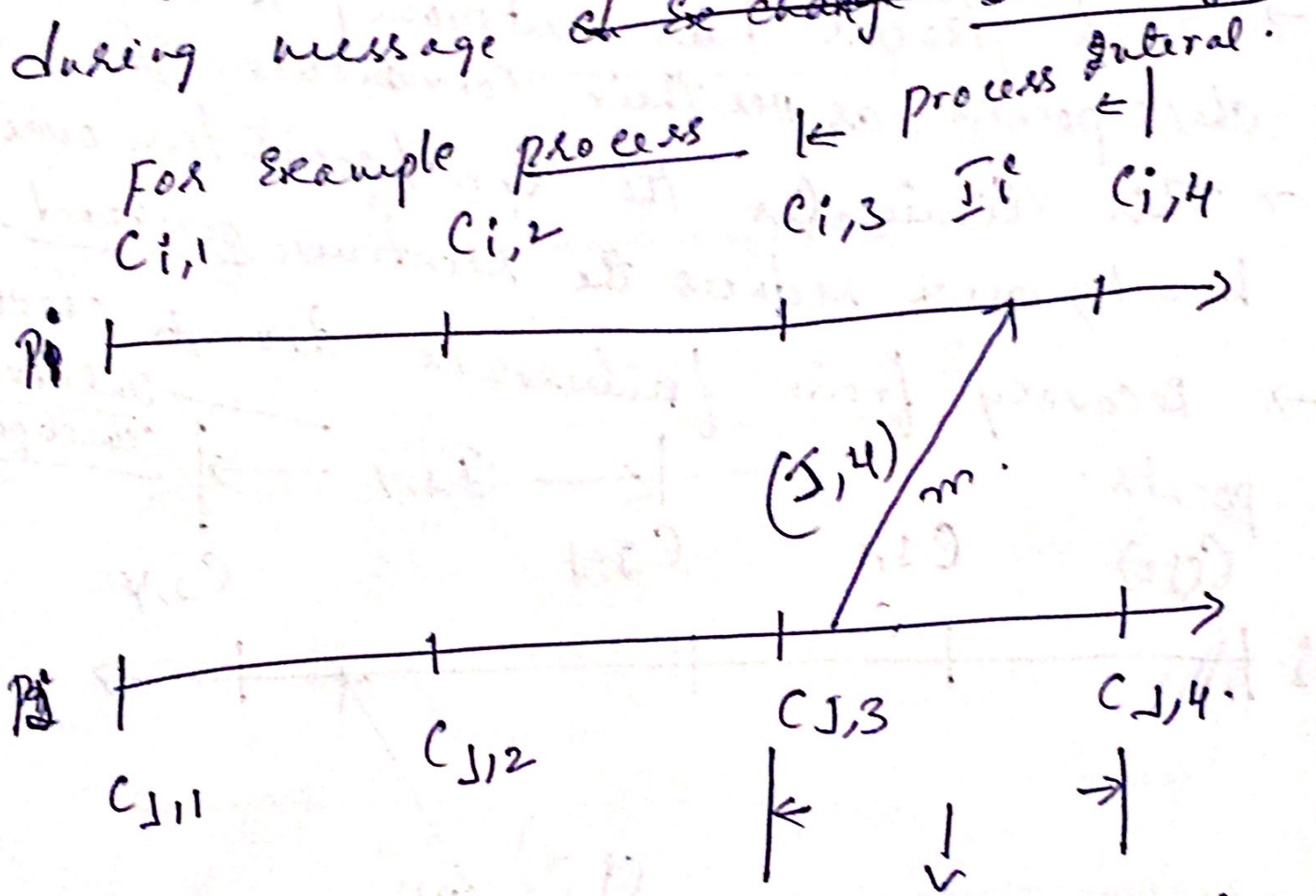
- In uncoordinated check pointing there is no need of coordination b/w processes.
- Each process has autonomy to take the check points as per their convenience.
- This eliminates the synchronization overhead, which reduces the runtime overhead.

→ Recovery from failure is through check points.



→ Recovery through this method is using a consistent Global check point.

→ A consistent global check point can be determined by recording the dependencies among the check points, caused during ~~message exchange~~ Exchange during message ~~exchange~~.



process P_j transfer (send) a message 'm' at during the process interval $(j,3)$ to process P_i at $(C_{j,3}$ to $C_{i,4})$ which is received at process interval $(i,3)$.

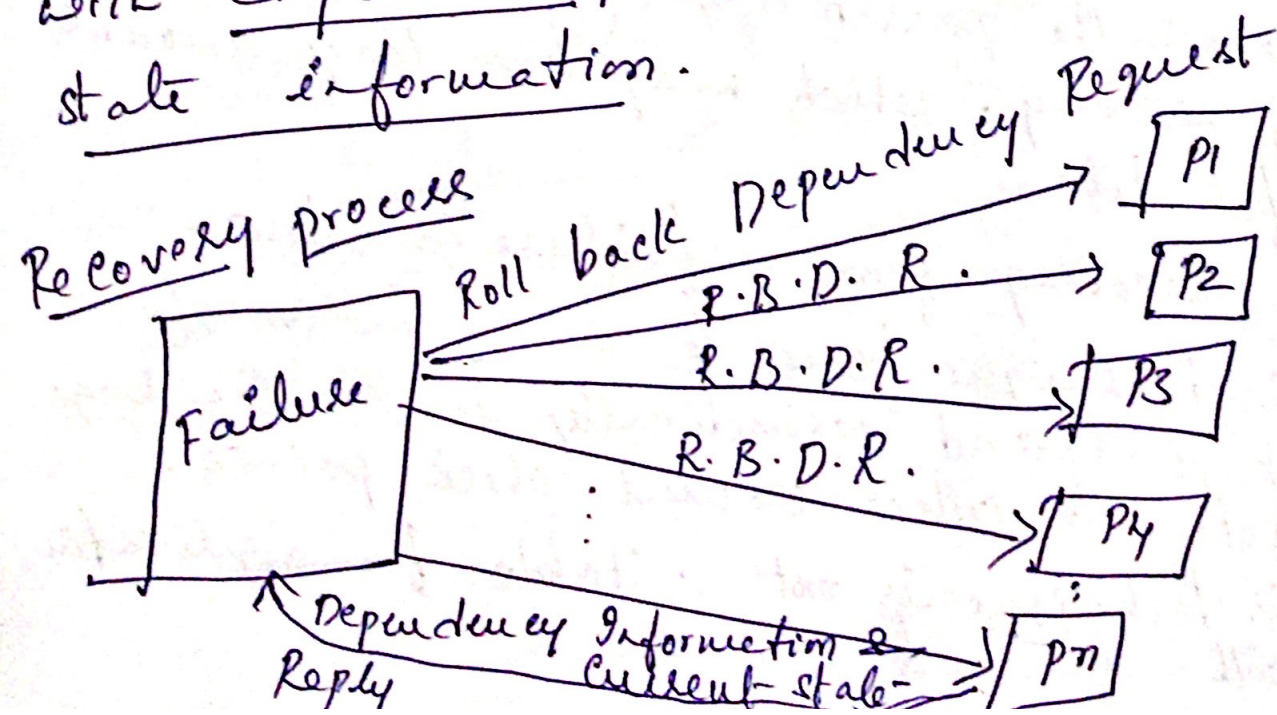
interval $(I, 4)$ i.e. between $C_{i,3}$ to $C_{i,4}$ ⑧

Along with the message 'm', its dependency
~~(I, 4)~~ also 'piggybacks' from processor P_j to P_i .

These dependencies are also saved on
stable storage when process P_i takes check
point $(C_{i,4})$.

→ When ever a failure occurs, the recover-
ing process initiates a rollback by broad-
casting a dependency request message to
all all dependencies information.

When a process receives this message,
it stops its execution and send a reply
with dependency information and its current
state information.



→ The recovery process then calculates the recovery line based up on the dependencies information from all processes (Global Dependency information).

→ After that it broadcasts a Roll back request message with the recovery line information.

→ up on receiving this information message, a process whose current state belongs to the recovery line resumes its execution or roll backs to an earlier checkpoint as indicated by the recovery line.

Disadvantages

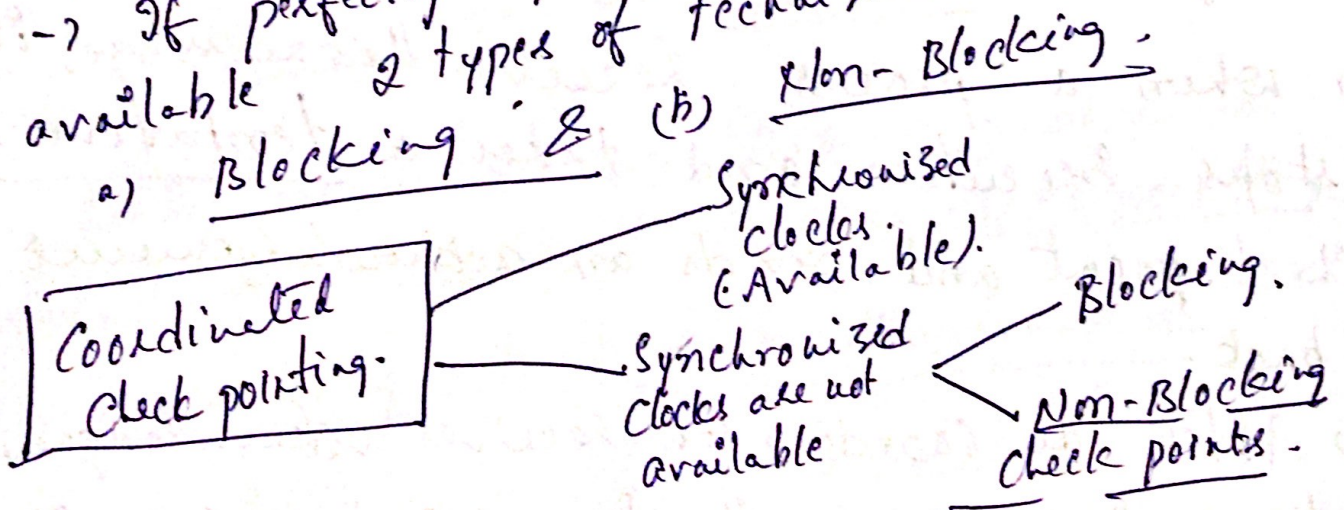
- 1) There is the possibility of Domino effect during recovery, which may lose large amounts of useful data.
- 2) The recovery from a failure is slow.
- 3) This technique forces to use more no. of check points and periodically invokes Garbage collector to collect unused check points.
- 4) This technique is not suitable for applications with frequent commits.

Coordinated check pointing - In Coordinated ⁽⁷⁾ check pointing all processes coordinate their check points information to form a Global Consistent State.

- > Coordinated check pointing simplifies recovery and it is not susceptible to 'Domino effect' maintains
- > Coordinated check pointing technique, only one check point on the stable storage.
- > This technique reduces the storage overhead and ^{there is no} need for Garbage collection.
- > Large latency is required in this technique.

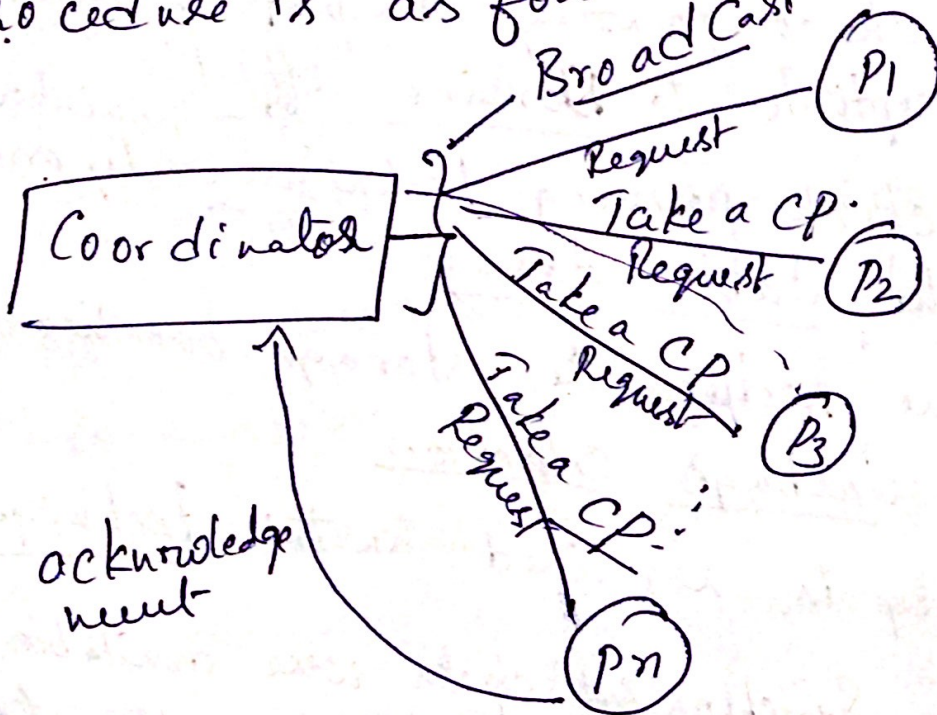
If perfectly Synchronized clocks were available at processes, All processes agree to take check points at particular instance of time.

-> If perfectly Synchronized clocks are not available 2 types of techniques are used. They are



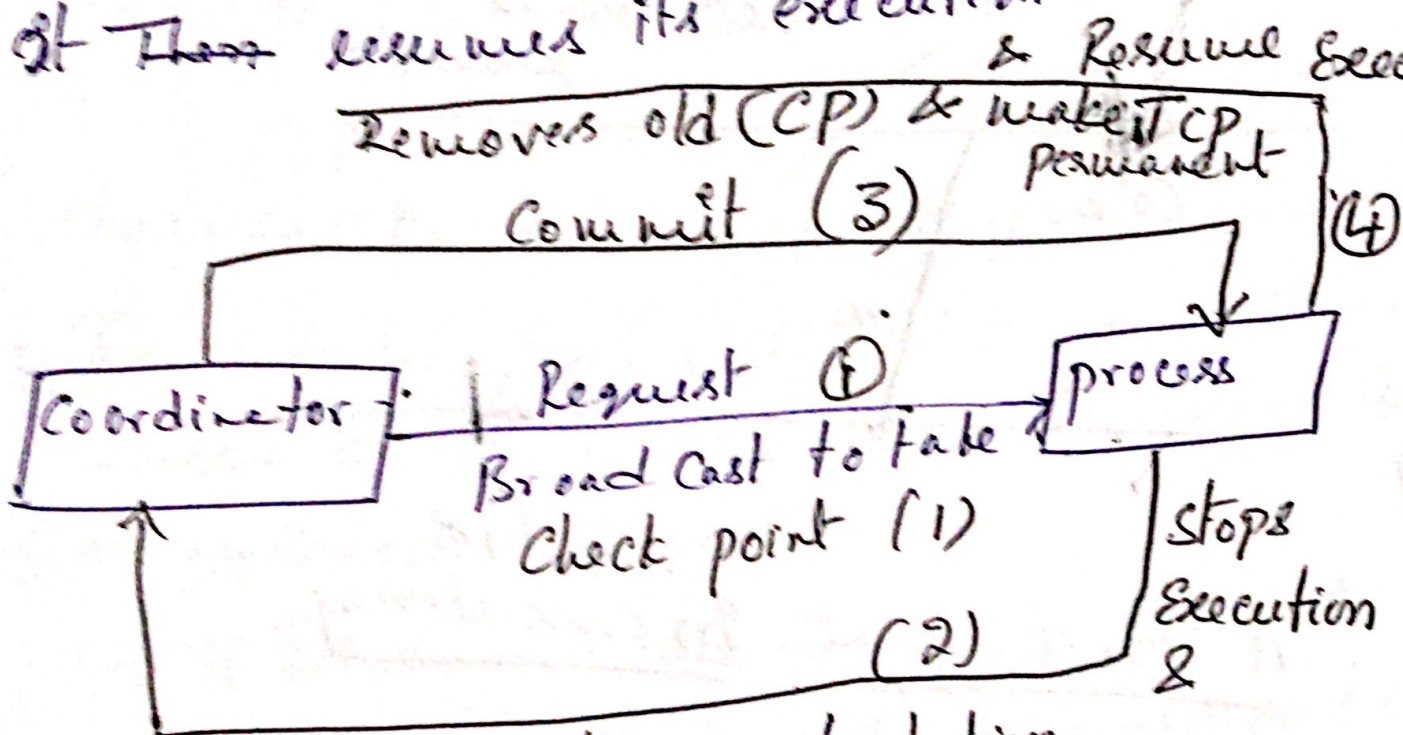
Blocking Coordinated Checkpointing.

In this technique, communications are blocked while checkpointing protocol executes. The procedure is as follows... Broadcast message.



- a) The coordinator takes a check point and broadcasts a request message to all processes asking them to take a check point.
- b) When a process receives this message, it stops execution and takes a tentative check point and sends an acknowledgment back.
- c) After, the coordinator receives acknowledgments from all processes, it broadcasts a commit message.

d) After, receiving the commit message the process removes old check point and makes the tentative check point permanent & then it ~~then~~ resumes its execution.



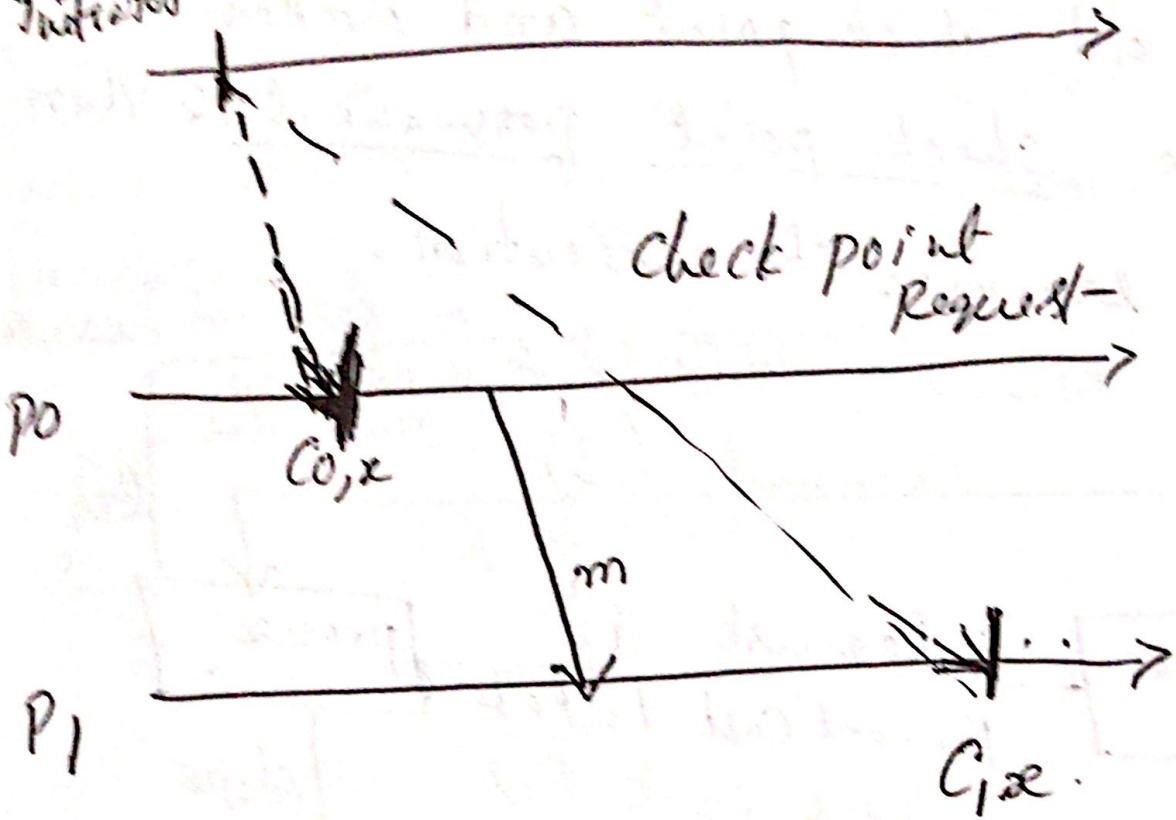
Sends an Acknowledge & takes a tentative check point
 (Same for all processes)

→ Only problem with this approach is Computation blocked during check pointing.

Non-blocking checkpoint coordination -

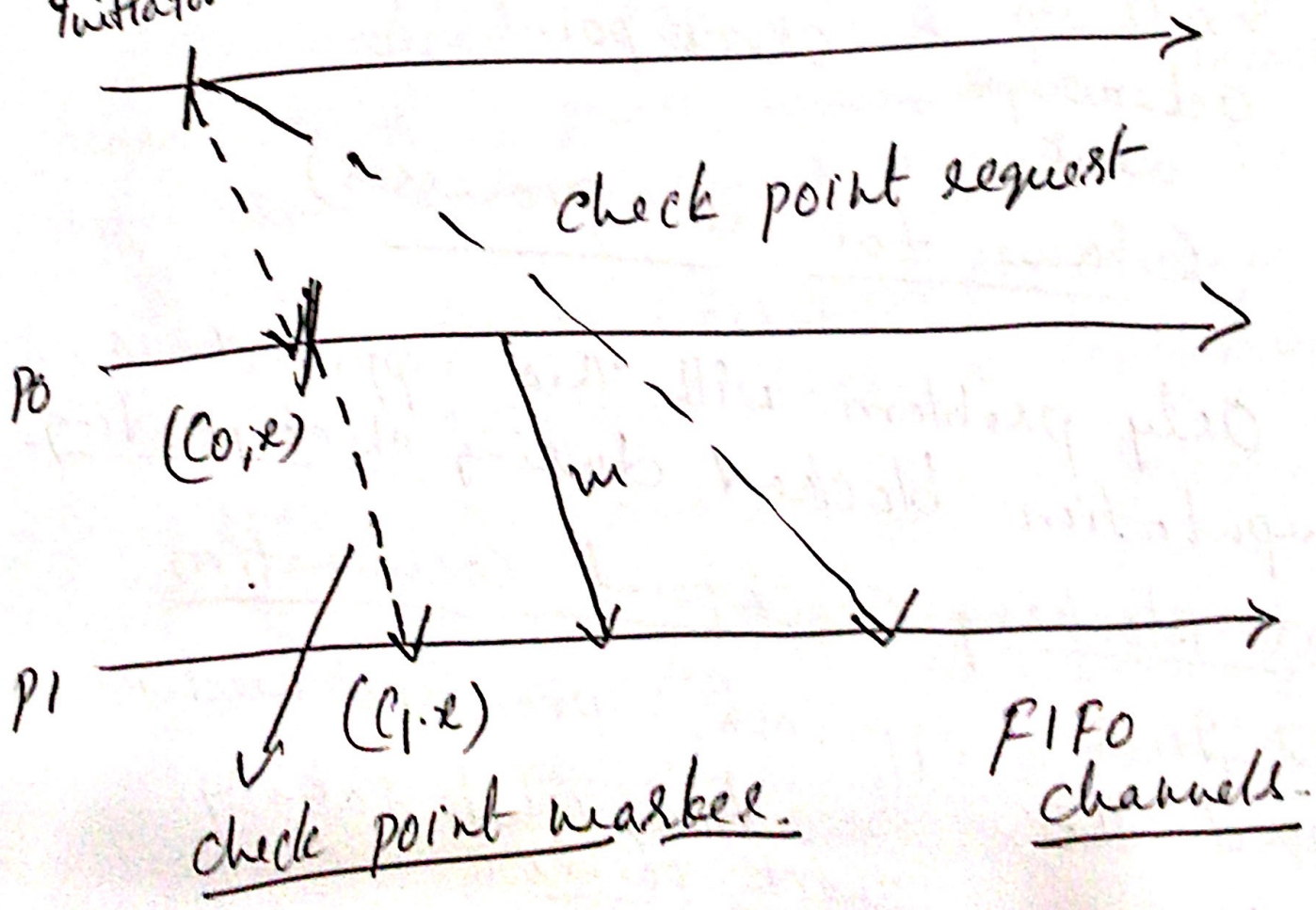
→ In this approach, processes need not stopped their execution, while taking check points.

Quintatos.



A check point in consistency.

Quintatos



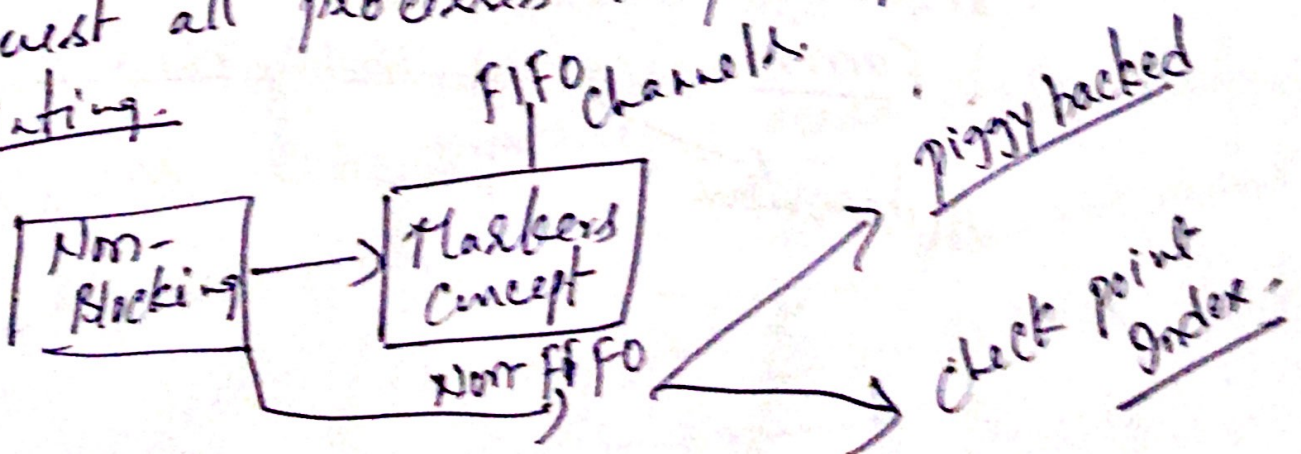
A fundamental problem in Non-blocking ^⑨ coordinated check point is that because of non blocking of messages, an application message can make check point inconsistent.

To avoid this problem, in FIFO channels the concept of "Markers" are used. The initiator takes a check point and sends a marker to all outgoing channels before sending an application message.

→ If the channels are non-FIFO, then the markers are "piggy backed" on every post-check point message.

(OR)

check point indices can serve this when the receiver's local check point index is lower than the piggy backed check point index, then coordinator request all processes to participate in check pointing.



Communication-induced check pointing Communication induced check pointing (CIC) protocols avoid the domino effect without requiring all checkpoints to be coordinated. Communication-induced check pointing reduces or completely eliminates the useless checkpoints.

In communication-induced check pointing, processes take two types of checkpoints, namely Autonomous and Forced checkpoints.

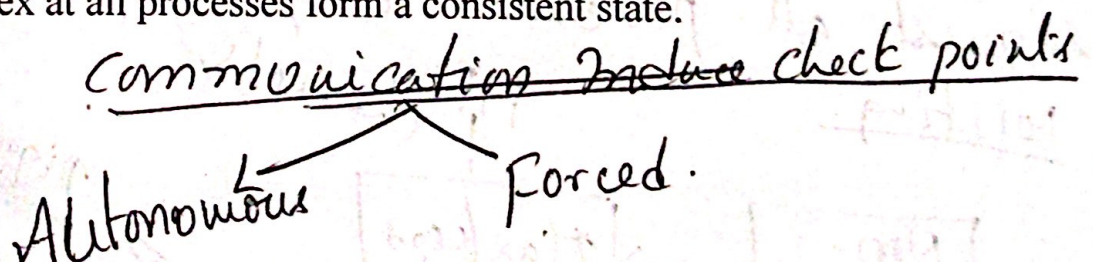
The checkpoints that a process takes independently are called local checkpoints, / Autonomous check Points, while the check points those are taken forcibly are called forced checkpoints.

The forced checkpoint must be taken before the application may process the contents of the message. It is therefore desirable in these systems to minimize the number of forced checkpoints.

Two types of communication-induced check pointing techniques are used. They are Model based check pointing and Index-based check pointing.

In model-based check pointing, the system maintains checkpoints and communication structures to prevent the domino effect.

In Index-based check pointing, the system uses an indexing scheme for the local and forced checkpoint, such that the checkpoints of the same index at all processes form a consistent state.



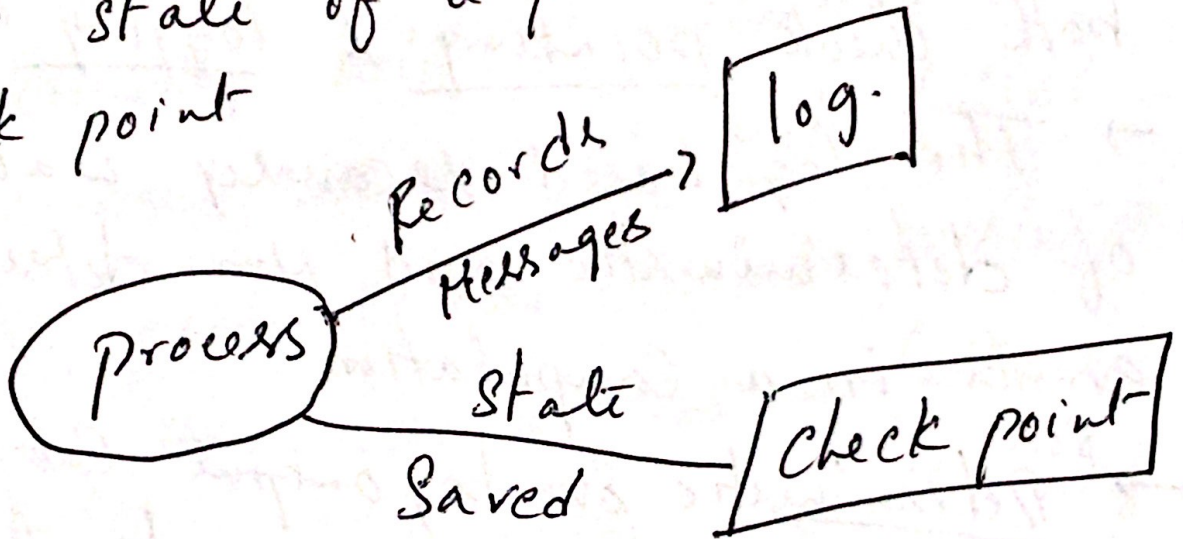
Log based Rollback Recovery.

⑩ ⑩

In general message logging & checkpointing are used for fault tolerance in Distributed Environment.

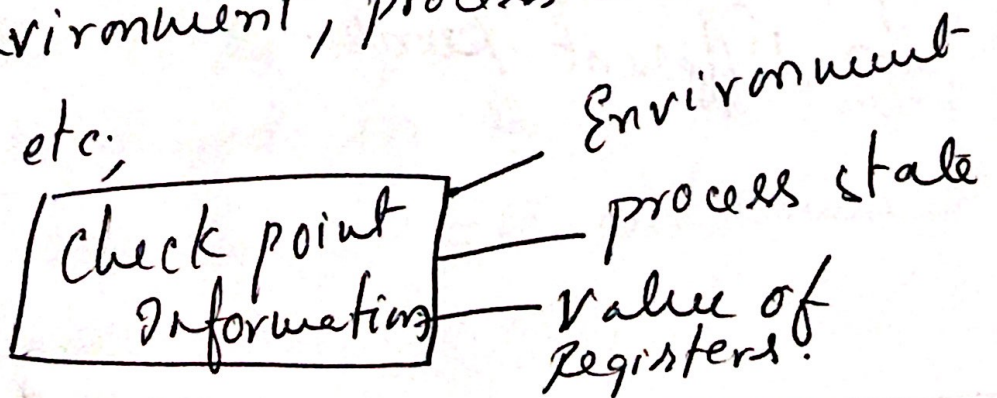
Each message received by a process is recorded in a 'log'.

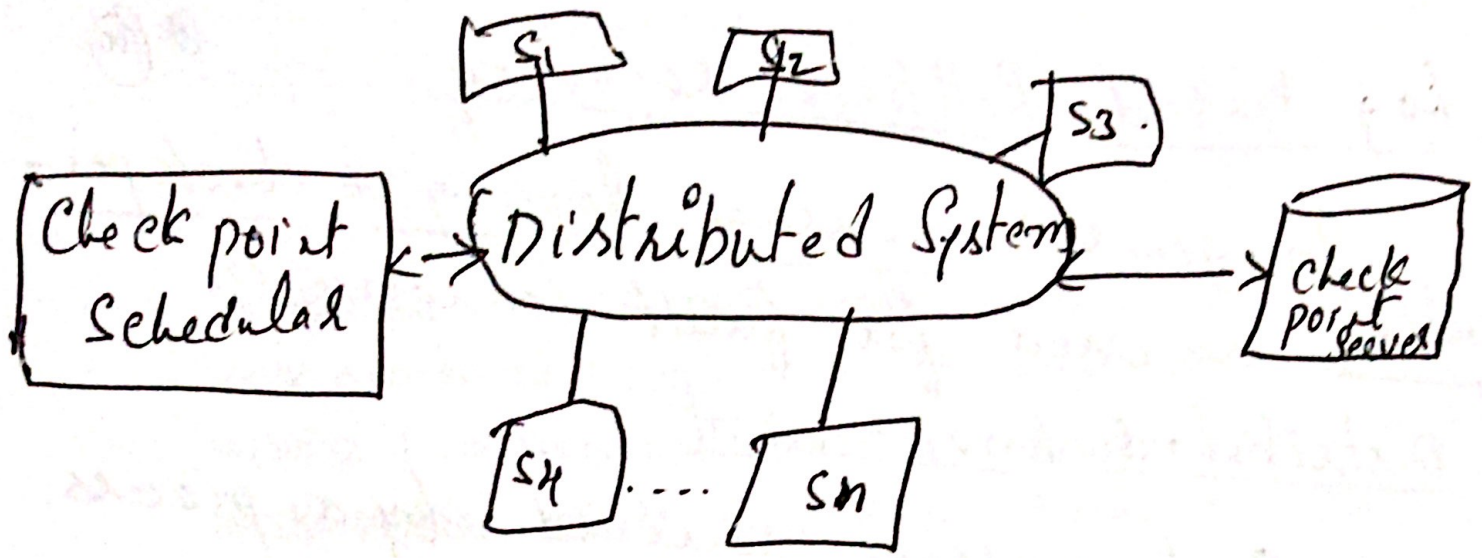
→ The state of a process is saved as check point



logged messages & checkpoints are saved in memory - which is useful during recovery.

The check point information that is stored includes Environment, process state, value of the registers etc;





→ Log Based recovery techniques combines both check pointing & logging of Events.

→ The log based recovery makes use of deterministic and non-deterministic events in a computation.

→ For Deterministic Events, ~~output can be~~
 For a given particular input the computer will always produce the same output.

In non deterministic events - for the same input, different outputs will be generated in different runs.

pessimistic message logging. (Believe that ¹⁷ the worst will happen).

→ pessimistic → the worst will happen.

→ pessimistic message logging protocols
log messages As Synchronously

→ These protocols guarantee that any failed processes can be recovered individually without affecting the states of processes, which did not failed.

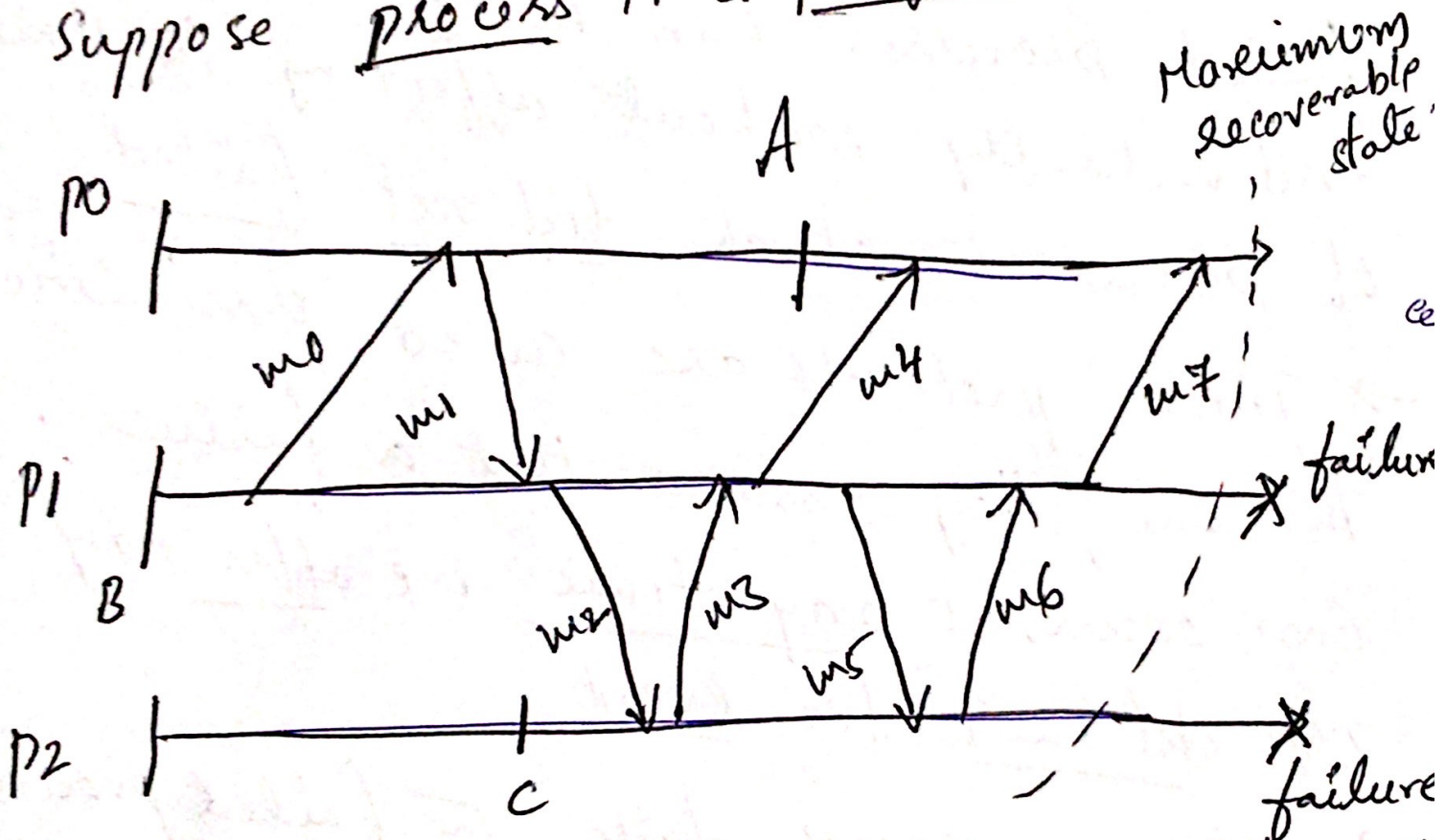
→ These protocols are called pessimistic because, they assume that a failure can occur, at any time i.e. after any non deterministic event.

→ When a process fails, the failed process is always restarted from its most recent checkpoint, and all the messages recorded ~~into~~ in the log after that-checkpoint are replayed in the same order.

→ Based on this the process deterministically reexecuted.

pessimistic message logging protocols
 have the advantage of being able to
 restart the system without affecting
 the states of other processes which
did not fail.

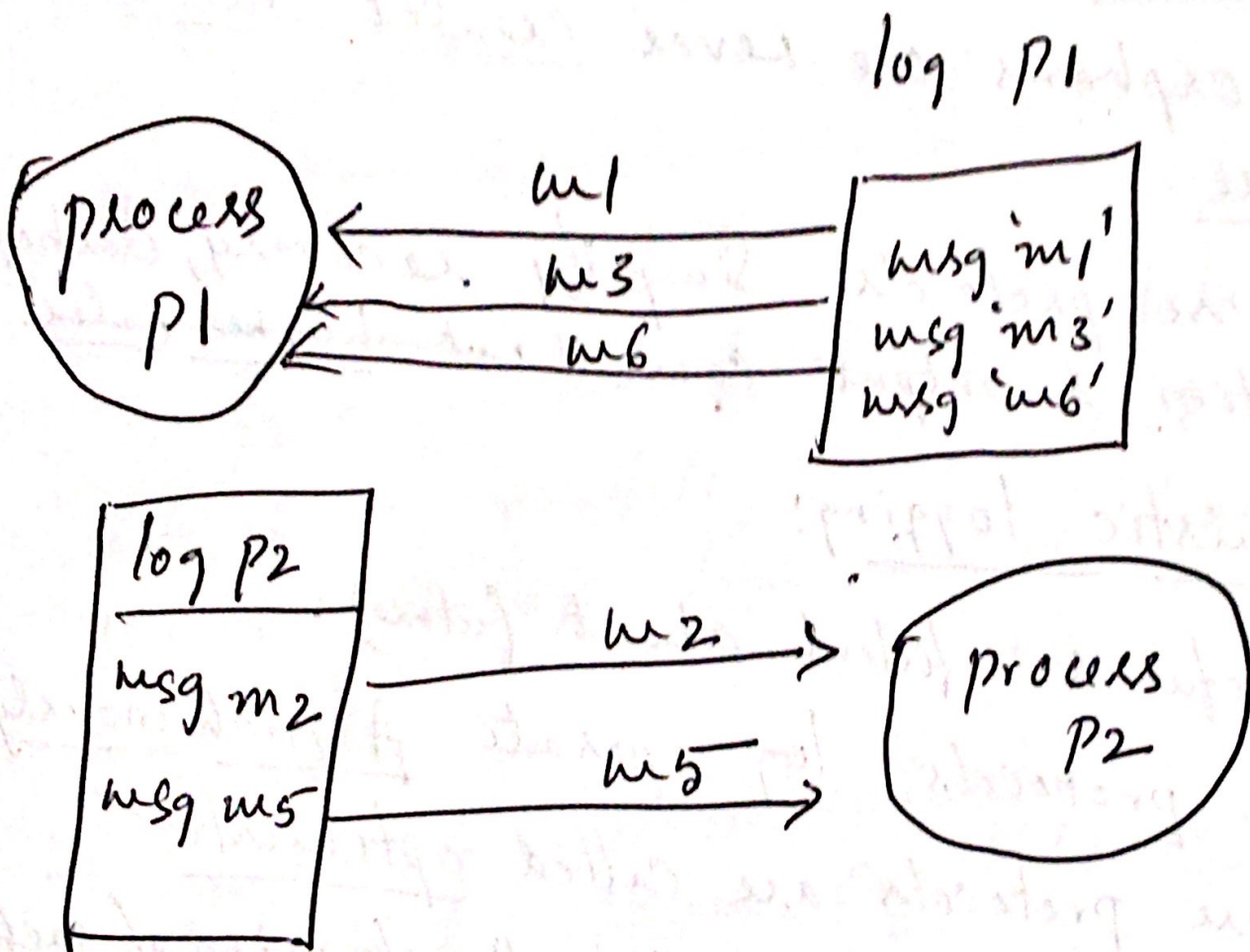
Suppose process P1 & P2 fail as shown



During the failure free operation, the logs of
 processes P0, P1, & P2 contains determini-
 nants needed to restart the system like
 as follows

$P_0 - \log - \{m_0, m_4, m_7\}$
 $P_1 - \{m_1, m_3, m_6\}$
 $P_2 - \{m_2, m_5\}$

If Suppose P1 & P2 fail and restarted from the check points B & C, their concerned determinant logs deliver the same messages in the same sequence as like in pre-failure execution.



This guarantees that P1 & P2 will repeat execution as a part of recovery.

After recovery both the processes will be consistent with the state of PO that include the receipt of message 'm7' from P1.

In pessimistic logging system, the observable state of each process is always recoverable.
→ Synchronous logging of results high performance.

Pessimistic logging - recovery protocols guarantee that orphans are never created due to failure.

→ These protocols simplify recovery, Garbage collection & output commit, at higher rates.

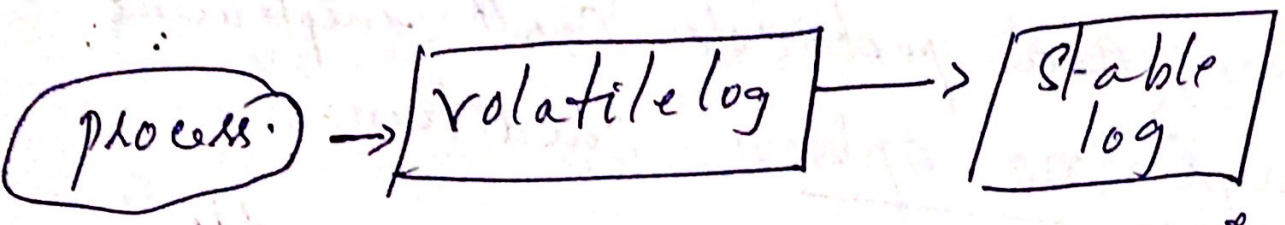
Optimistic logging:-

↓
(Hopeful & confident about future).

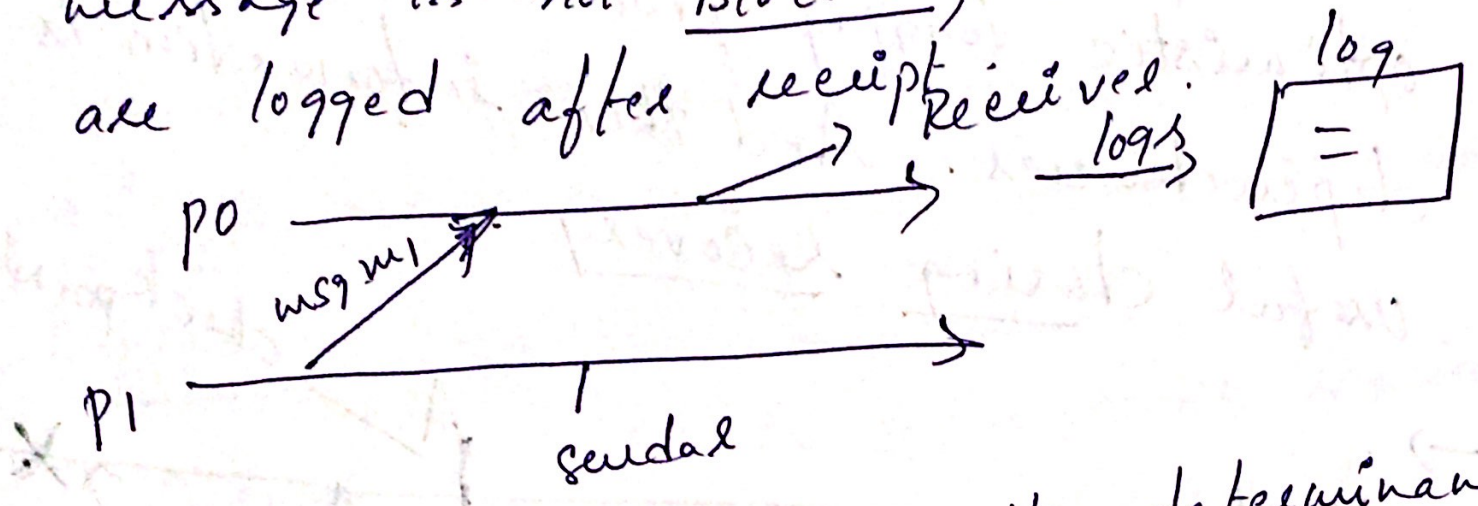
These protocols log asynchronously

→ These protocols are called 'optimistic' because, they assume that the logging of each message received by a process will be completed before the process fails.

→ Determinants are first kept in volatile log and are periodically flushed to the stable storage.



In this technique, the receiver of a message is not Blocked, and messages are logged after receipt received.



→ If a process fails, the determinants in its volatile log are lost. So roll back through logs can be done from stable log.

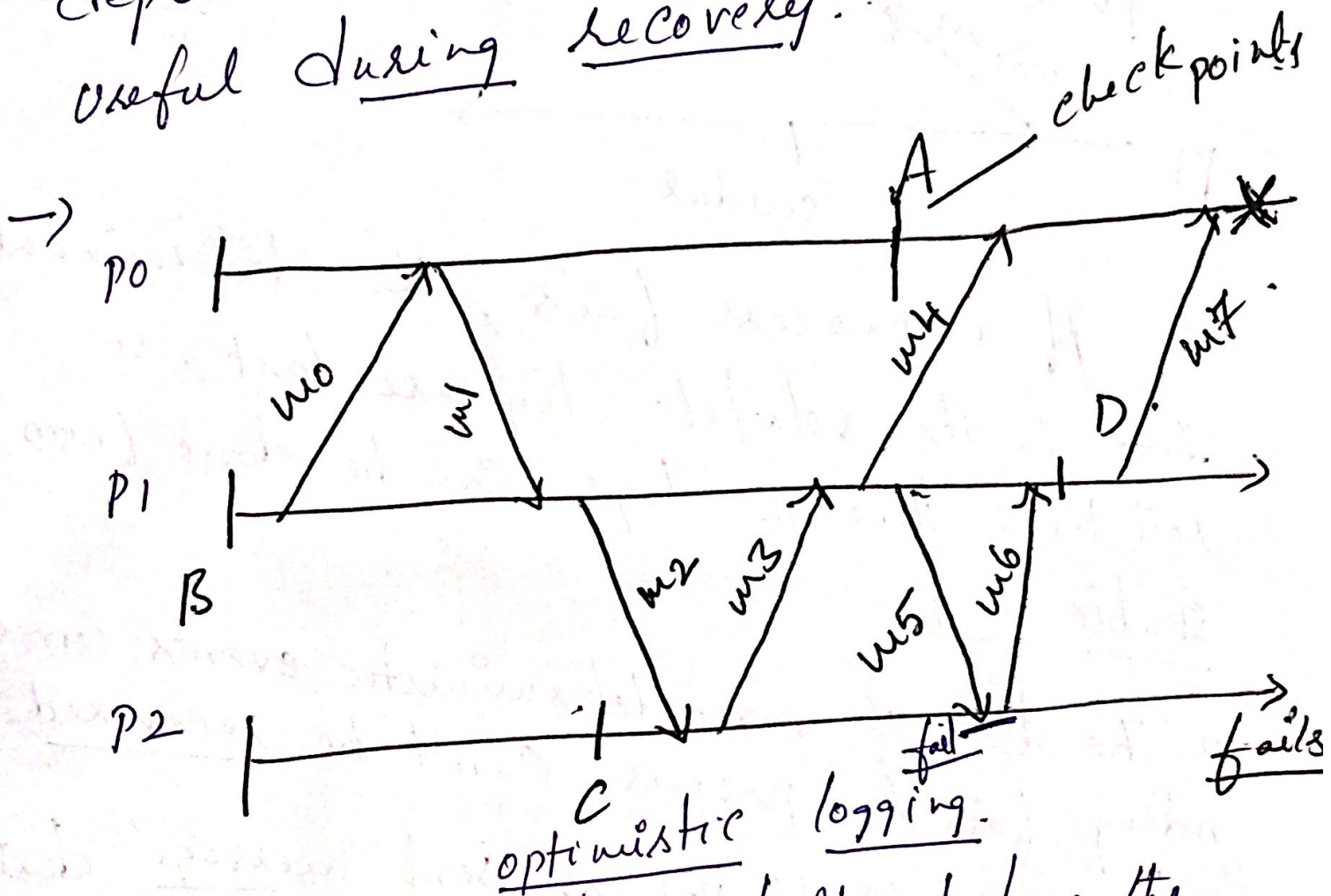
→ The state of non-deterministic events corresponding failed process can't be recovered.

→ If the failed process 'sent message' during state intervals can't be recovered.

→ The receiver of a message becomes an orphan process and it must also roll back to undo the effects of receiving message.

→ Hence, these protocols can't implement.
'Always - no - orphan' condition:

→ To perform roll back correctly, optimistic logging protocols use causal dependencies. And this information is useful during recovery.



Suppose process P2 fails before the determinant for m_5 is logged to the stable storage, process P1 becomes an 'orphan' process & must roll back to undo the effects of message m_6

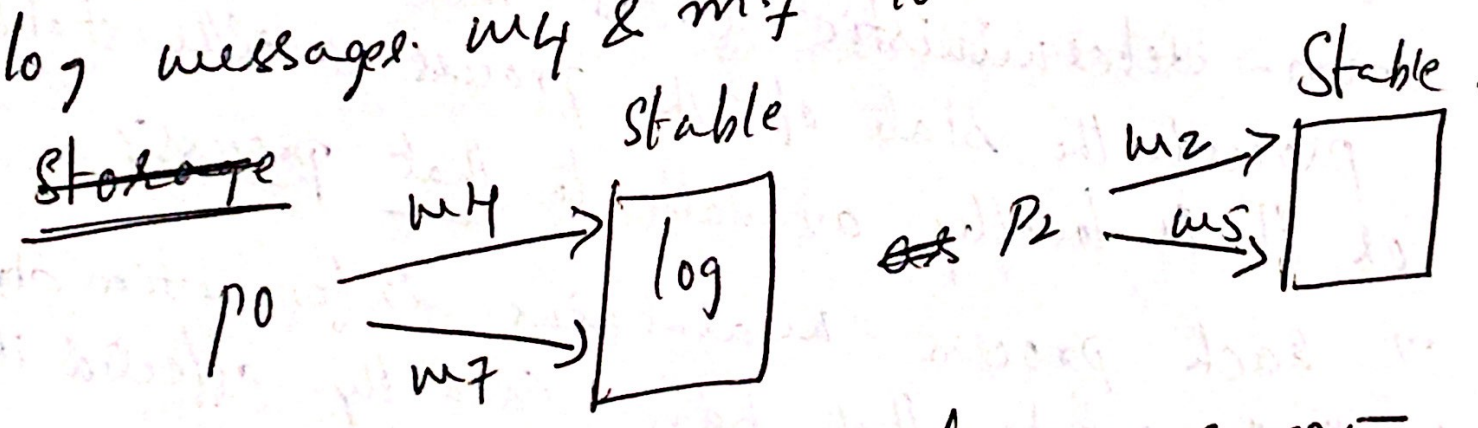
-> The roll back of P_1 further forces P_0 to roll back to undo the effects of receiving message m_7 .

-> Optimistic protocols need to keep multiple checkpoints for each process.

-> Optimistic logging protocols require a "non-trivial garbage collection scheme"

-> Output Commit in OL protocols requires a guarantee that no failure scenario can revoke the output.

For example if process P_0 needs to commit output at state X , it must log message m_4 & m_7 to the stable



Storage & ask P_2 to log m_2 & m_5 .
In this case if process fails, The computation can be reconstructed up to state X'

Casual logging

- Casual logging combines the advantages of both pessimistic & optimistic.
- Like optimistic, it does not require Synchronous Access to the stable storage except during output commit.
- Like pessimistic logging, it allows each process to commit output independently & never creates orphans.
- Rollback is limited to the most recent check point on the stable storage.
- [On this approach, the determinant of each non-deterministic event that casually proceeds the state of the process is either stable or it is locally available to that process]
- Each process maintains information about all the events that have casually affected its state.
- This information is useful during failures

The check pointing algorithm

The checkpoint algorithm works on the following assumptions about the distributed system:

- Processes communicate by exchanging messages through communication channels.
- Communication channels are FIFO.
- Communication failures do not partition the network.
- No process fails during the execution of the algorithm.

The checkpoint algorithm takes two kinds of checkpoints on the stable storage, **They are Permanent and Tentative.**

A permanent checkpoint is a local checkpoint at a process and is a part of a consistent global checkpoint.

A tentative checkpoint is a temporary checkpoint

In case of a failure, processes roll back only to their permanent checkpoints for recovery.

First phase

An initiating process P_i takes a tentative checkpoint and requests all other processes to take tentative checkpoints.

After taking the tentative check points, Each process informs to the initiator process P_i whether it succeeded in taking a tentative checkpoint.

A process says “no” to a request if it fails to take a tentative checkpoint.

If P_i learns that all the processes have successfully taken tentative checkpoints, P_i decides that all tentative checkpoints should be made permanent; otherwise, P_i decides that all the tentative checkpoints should be discarded.

Second phase

In second phase, after all check points the initiator process P_i informs to all that we have reached at the end of the first phase. And do not send any messages until there is an instruction / Decision.

[The algorithm requires that after a process has taken a tentative checkpoint, it cannot send messages related to the underlying computation until it is informed of P_i 's decision.]

Correctness

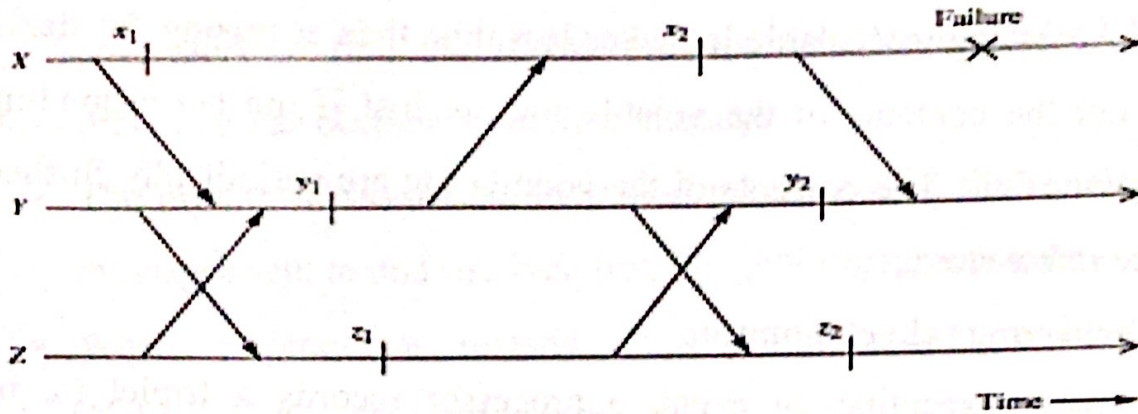
A set of permanent checkpoints taken by this algorithm is consistent because of the following two reasons:

First, Either all or none of the processes take permanent checkpoints, second, no process sends a message after taking a tentative checkpoint until the receipt of the initiating process's decision.

An optimization : Taking unnecessary check points are avoided.

Algorithm for Asynchronous Check pointing and Recovery

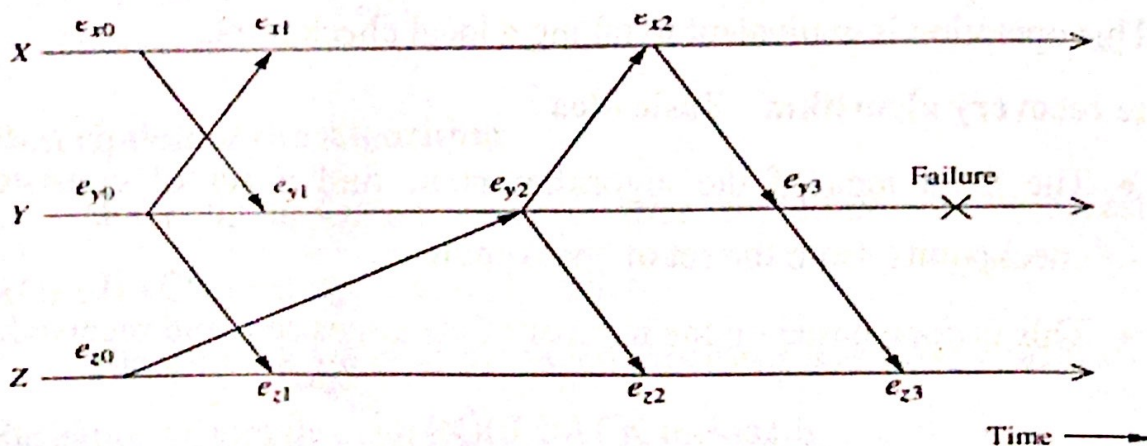
(Juang-Venkatesan algorithm for asynchronous checkpointing and recovery)



System model and assumptions

The algorithm makes the following assumptions about the underlying system:

- The communication channels are reliable, deliver the messages in FIFO order, and have infinite buffers.
- The message transmission delay is arbitrary, but finite.
- The processors directly connected to a processor via communication channels are called its neighbors.
- The underlying computation or application is assumed to be event-driven.



To facilitate recovery after a process failure and restore the system to a consistent state, two types of log storage are maintained, volatile log and stable log.

Accessing the volatile log takes less time than accessing the stable log, but the contents of the volatile log are lost if the corresponding processor fails. The contents of the volatile log are periodically flushed to the stable storage.

Asynchronous check pointing

After executing an event, a processor records a triplet (s, m, msgs_sent) in its volatile storage,

where 's' is the state of the processor before the event,

'm' is the message

'msgs_sent' is the set of messages that were sent by the processor during the event.

- A local checkpoint at a processor consists of the record of an event occurring at the processor and it is taken without any synchronization with other processors.
- Periodically, a processor independently saves the contents of the volatile log in the stable storage and clears the volatile log.
- This operation is equivalent to taking a local checkpoint.

The recovery algorithm - Basic idea

- The main idea of the algorithm is to find a set of consistent checkpoints, from the set of checkpoints.
- This is done based on the number of messages sent and received.

Consensus & Agreement Algorithm

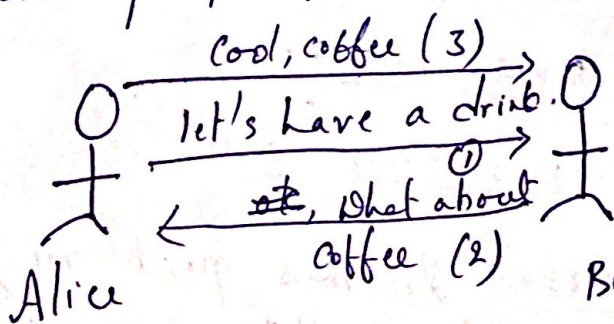
Problem Definition

The Byzantine General problem is an analogy in Computer Science used to describe the challenge of establishing & maintaining security on distributed Network. For this all the nodes (majority of nodes) must establish a set of rules & come to an agreement.

Agreement in distributed system among processes is a fundamental problem requirement for a wide range of applications.

Agreement is a Common Understanding. A classical example is 'Commit' or 'Abort' a transaction must be done with a Common Agreement.

For example, 2 persons wanted to have a drink.



- 1) Alice & Bob wanted to have a drink.
- 2) Bob suggested 'Coffee'
- 3) Alice 'Agreed' to have a drink.

They both agreed on a 'value' that is suggested by one of them. This is Consensus.

In simple words a Consensus is a general agreement on a value.

Hence this is known as Agreement variable or Value. This may be a boolean value or Multivalued and need not be an integer.

The Byzantine & other problems. The Byzantine Agreement problem requires a designated process called "source process" with an initial value to reach an agreement with other processes about its "initial value" with the following conditions.

Agreement: All Non-Faulty processes must agree on the same value.

Validity: Validity says that the value ^{Suggested} ~~proposed~~ by one of the node should be Accepted by all non-faulty nodes.

Termination: Every non-faulty node will eventually decide some value.

Other 2 Important problems are

The consensus problem: In this problem, each process has an initial value. And this value ^{is} ~~is~~ broadcasted to all other processes. This ^{initial} value may be different.

The following are the conditions.

Agreement: All non-faulty processors agree on the same single value.

Validity: If all the non-faulty processors have the same initial value, it must be agreed by all non-faulty processors.

Termination: Each non-faulty processor must eventually decide on a value.

The iterative consistency problem
In this problem, that each processor has an initial value, and all the correct processes must agree upon a set of values, with one value for each process $\{v_1, \dots, v_n\}$.

Agreement: All non-faulty processes must agree on the same array of values $A[v_1, \dots, v_n]$

Validity: If process i is non-faulty and its initial value is v_i , then all the non-faulty processes agree on v_i , as the i th element of the array A . If process j is faulty, then non-faulty processes can agree on any value for $A[j]$

Termination: Each non-faulty process must eventually decide on the array A .

The 3 problems defined are equivalent in the sense. A solution to any of one of them can be used as a solution to other 2 problems.

Problem	Byzantine Agreement	Consensus	Interactive Consistency
Initiation of a value	one processor	All processors	All processors.
Final Agreement	Single value	Single value	A vector of values

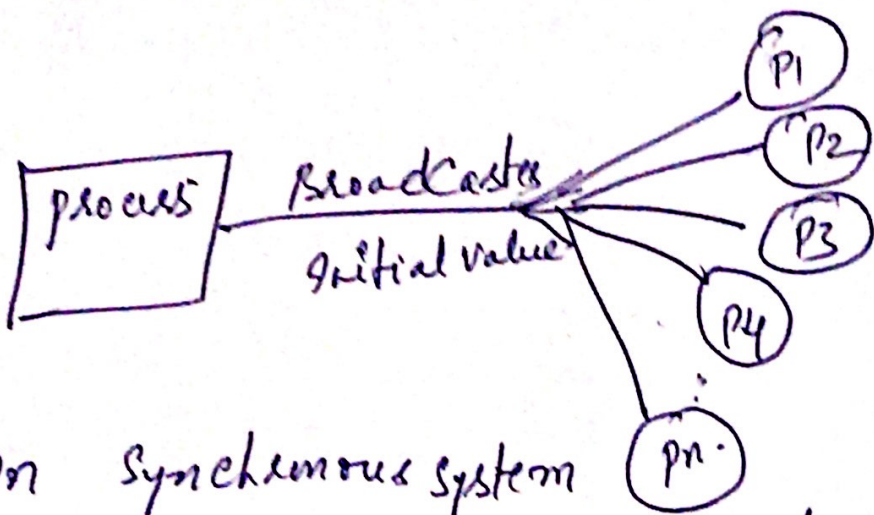
Agreement in a failure free system.

In a failure free system, consensus can be reached by collecting information from the different processes, and performs a decision.

In this each process broadcasts its value to others.

The decision can be reached by using some system functions like 'win', 'max()', 'majority'

To collect the initial values a token is circulated in logical ring (or) through direct communication.



→ In Synchronous system
 this can be done with simply using constant
 no. of rounds (Based up on the topology)

→ In another round decision value can be
 obtained.

→ In an Asynchronous System, Consensus can
 similarly be reached in a constant no. of
message hops.

From these hops Common decision value
 is obtained.