

## Chapter 07

(version 26th October 2001)

Maarten van Steen

Vrije Universiteit Amsterdam, Faculty of Science

Dept. Mathematics and Computer Science

Room R4.20. Tel: (020) 444 7784

E-mail: [steen@cs.vu.nl](mailto:steen@cs.vu.nl), URL: [www.cs.vu.nl/~steen](http://www.cs.vu.nl/~steen)

01	Introduction
02	Communication
03	Processes
04	Naming
05	Synchronization
06	Consistency and Replication
07	Fault Tolerance
08	Security
09	Distributed Object-Based Systems
10	Distributed File Systems
11	Distributed Document-Based Systems
12	Distributed Coordination-Based Systems

00 – 1

/

## Introduction

- Basic concepts
- Process resilience
- Reliable client-server communication
- Reliable group communication
- Distributed commit
- Recovery

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible]

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

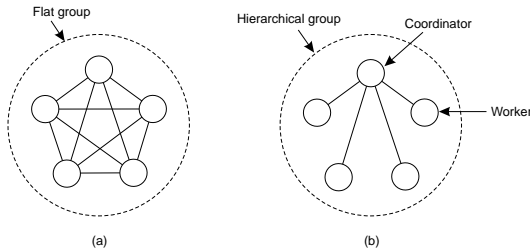


# Process Resilience

**Basic issue:** Protect yourself against faulty processes by replicating and distributing computations in a group.

**Flat groups:** Good for fault tolerance as information exchange immediately occurs with all group members; however, may impose more overhead as control is completely distributed (hard to implement).

**Hierarchical groups:** All communication through a single coordinator  $\Rightarrow$  not really fault tolerant and scalable, but relatively easy to implement.



07 – 6

Fault Tolerance/7.2 Process Resilience

## Groups and Failure Masking (1/3)

**Terminology:** when a group can mask any  $k$  concurrent member failures, it is said to be  **$k$ -fault tolerant** ( $k$  is called degree of fault tolerance).

**Problem:** how large does a  $k$ -fault tolerant group need to be?

- Assume crash/performance failure semantics  $\Rightarrow$  a total of  $k + 1$  members are needed to survive  $k$  member failures.
- Assume arbitrary failure semantics, and group output defined by voting  $\Rightarrow$  a total of  $2k + 1$  members are needed to survive  $k$  member failures.

**Assumption:** all members are identical, and process all input in the same order  $\Rightarrow$  only then are we sure that they do exactly the same thing.

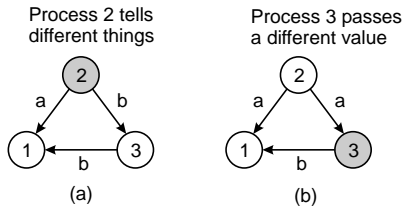
07 – 7

Fault Tolerance/7.2 Process Resilience

## Groups and Failure Masking (2/3)

**Assumption:** Group members are not identical, i.e., we have a distributed computation

**Problem:** Nonfaulty group members should reach agreement on the same value



**Observation:** Assuming arbitrary failure semantics, we need  $3k + 1$  group members to survive the attacks of  $k$  faulty members

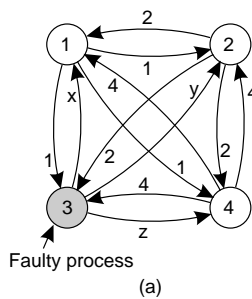
**Note:** This is also known as **Byzantine failures**.

**Essence:** We are trying to reach a majority vote among the group of loyalists, in the presence of  $k$  traitors  $\Rightarrow$  need  $2k + 1$  loyalists.

07 – 8

Fault Tolerance/7.2 Process Resilience

## Groups and Failure Masking (3/3)



1 Got(1, 2, x, 4)	1 Got	2 Got	4 Got
2 Got(1, 2, y, 4)	(1, 2, y, 4)	(1, 2, x, 4)	(1, 2, x, 4)
3 Got(1, 2, 3, 4)	(a, b, c, d)	(e, f, g, h)	(1, 2, y, 4)
4 Got(1, 2, z, 4)	(1, 2, z, 4)	(1, 2, z, 4)	(i, j, k, l)

(b)

(c)

- (a) what they send to each other
- (b) what each one got from the other
- (c) what each one got in second step

07 – 9

Fault Tolerance/7.2 Process Resilience

# Reliable Communication

**So far:** Concentrated on **process resilience** (by means of process groups). What about reliable communication channels?

## Error detection:

- Framing of packets to allow for bit error detection
- Use of frame numbering to detect packet loss

## Error correction:

- Add so much redundancy that corrupted packets can be automatically *corrected*
- Request retransmission of lost, or last  $N$  packets

**Observation:** Most of this work assumes point-to-point communication

## Reliable RPC (1/3)

### What can go wrong?:

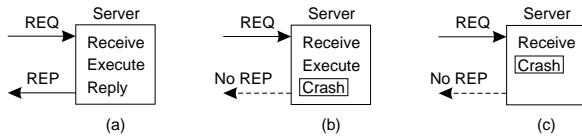
- 1: Client cannot locate server
- 2: Client request is lost
- 3: Server crashes
- 4: Server response is lost
- 5: Client crashes

[1:] Relatively simple – just report back to client

[2:] Just resend message

## Reliable RPC (2/3)

[3:] Server crashes are harder as you don't what it had already done:



**Problem:** We need to decide on what we expect from the server

- **At-least-once-semantics:** The server guarantees it will carry out an operation at least once, no matter what
- **At-most-once-semantics:** The server guarantees it will carry out an operation at most once.

## Reliable RPC (3/3)

[4:] Detecting lost replies can be hard, because it can also be that the server had crashed. You don't know whether the server has carried out the operation

**Solution:** None, except that you can try to make your operations **idempotent**: repeatable without any harm done if it happened to be carried out before.

[5:] **Problem:** The server is doing work and holding resources for nothing (called doing an **orphan** computation).

- Orphan is killed (or rolled back) by client when it reboots
- Broadcast new epoch number when recovering  
⇒ servers kill orphans
- Require computations to complete in a  $T$  time units. Old ones are simply removed.

**Question:** What's the rolling back for?

## Reliable Multicasting (1/2)

**Basic model:** We have a **multicast channel**  $c$  with two (possibly overlapping) groups:

- **The sender group**  $SND(c)$  of processes that *submit* messages to channel  $c$
- **The receiver group**  $RCV(c)$  of processes that can receive messages from channel  $c$

**Simple reliability:** If process  $P \in RCV(c)$  at the time message  $m$  was submitted to  $c$ , and  $P$  does not leave  $RCV(c)$ ,  $m$  should be delivered to  $P$

**Atomic multicast:** How can we ensure that a message  $m$  submitted to channel  $c$  is delivered to process  $P \in RCV(c)$  only if  $m$  is delivered to *all* members of  $RCV(c)$

## Reliable Multicasting (2/2)

**Observation:** If we can stick to a local-area network, reliable multicasting is “easy”

**Principle:** Let the sender log messages submitted to channel  $c$ :

- If  $P$  sends message  $m$ ,  $m$  is stored in a **history buffer**
- Each receiver acknowledges the receipt of  $m$ , or requests retransmission at  $P$  when noticing message lost
- Sender  $P$  removes  $m$  from history buffer when everyone has acknowledged receipt

**Question:** Why doesn't this scale?



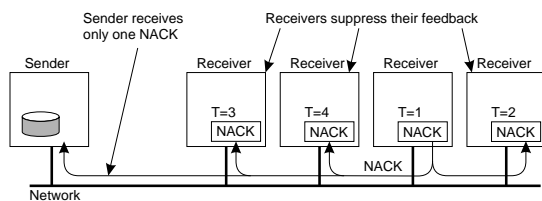
# Scalable Reliable Multicasting: Feedback Suppression

**Basic idea:** Let a process  $P$  suppress its own feedback when it notices another process  $Q$  is already asking for a retransmission

## Assumptions:

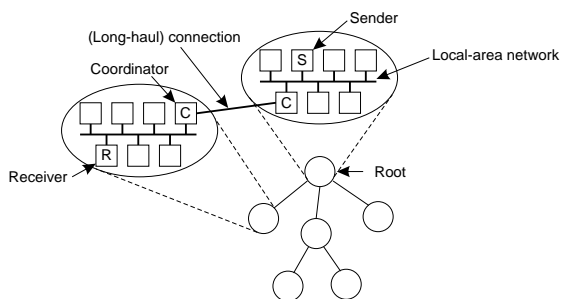
- All receivers listen to a common **feedback channel** to which feedback messages are submitted
- Process  $P$  schedules its own feedback message *randomly*, and suppresses it when observing another feedback message

**Question:** Why is the random schedule so important?



# Scalable Reliable Multicasting: Hierarchical Solutions

**Basic solution:** Construct an hierarchical feedback channel in which all submitted messages are sent only to the root. Intermediate nodes aggregate feedback messages before passing them on.

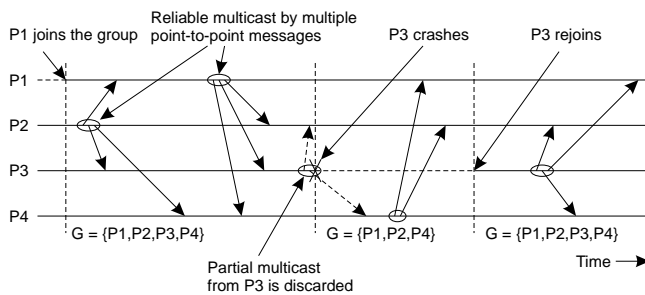


**Question:** What's the main problem with this solution?

**Observation:** Intermediate nodes can easily be used for retransmission purposes

# Atomic Multicast

**Idea:** Formulate reliable multicasting in the presence of process failures in terms of process groups and changes to group membership:



**Guarantee:** A message is delivered only to the non-faulty members of the current group. All members should agree on the current group membership.

**Keyword:** Virtually synchronous multicast

## Virtual Synchrony (1/2)

**Essence:** We consider **views**  $V \subseteq RCV(c) \cup SND(c)$

Processes are added or deleted from a view  $V$  through **view changes** to  $V^*$ ; a view change is to be executed *locally* by each  $P \in V \cap V^*$

- (1) For each consistent state, there is a **unique view** on which all its members agree. **Note:** implies that all nonfaulty processes see all view changes in the same order
- (2) If message  $m$  is sent to  $V$  before a view change  $vc$  to  $V^*$ , then either all  $P \in V$  that execute  $vc$  receive  $m$ , or no processes  $P \in V$  that execute  $vc$  receive  $m$ . **Note:** all nonfaulty members in the same view get to see the same set of multicast messages.
- (3) A message sent to view  $V$  can be delivered only to processes in  $V$ , and is discarded by successive views

A reliable multicast algorithm satisfying (1)–(3) is **virtually synchronous**

## Virtual Synchrony (2/2)

- A sender to a view  $V$  need not be member of  $V$
- If a sender  $S \in V$  crashes, its multicast message  $m$  is *flushed* before  $S$  is removed from  $V$ :  $m$  will never be delivered after the point that  $S \notin V$

**Note:** Messages from  $S$  may still be delivered to all, or none (nonfaulty) processes in  $V$  before they all agree on a new view to which  $S$  does not belong

- If a receiver  $P$  fails, a message  $m$  may be lost but can be recovered as we know exactly what has been received in  $V$ . Alternatively, we may decide to deliver  $m$  to members in  $V - \{P\}$

**Observation:** Virtually synchronous behavior can be seen independent from the ordering of message delivery. The only issue is that messages are delivered to an *agreed upon* group of receivers.

07 – 20

Fault Tolerance/Reliable Communication

## Virtual Synchrony Implementation (1/3)

- The current view is known at each  $P$  by means of a delivery list  $\text{dest}[P]$
- If  $P \in \text{dest}[Q]$  then  $Q \in \text{dest}[P]$
- Messages received by  $P$  are queued in  $\text{queue}[P]$
- If  $P$  fails, the group view must change, but not before all messages from  $P$  have been flushed
- Each  $P$  attaches a (stepwise increasing) **time-stamp** with each message it sends
- Assume FIFO-ordered delivery; the highest numbered message from  $Q$  that has been received by  $P$  is recorded in  $\text{rcvd}[P][Q]$
- The vector  $\text{rcvd}[P][\ ]$  is sent (as a control message) to all members in  $\text{dest}[P]$
- Each  $P$  records  $\text{rcvd}[Q][\ ]$  in  $\text{remote}[P][Q]$

07 – 21

Fault Tolerance/Reliable Communication

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible]This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible]

- 
- This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

07-24

## Fault Tolerance/7.5 Distributed Commit

**Model:** The client who initiated the computation acts as coordinator; processes required to commit are the participants

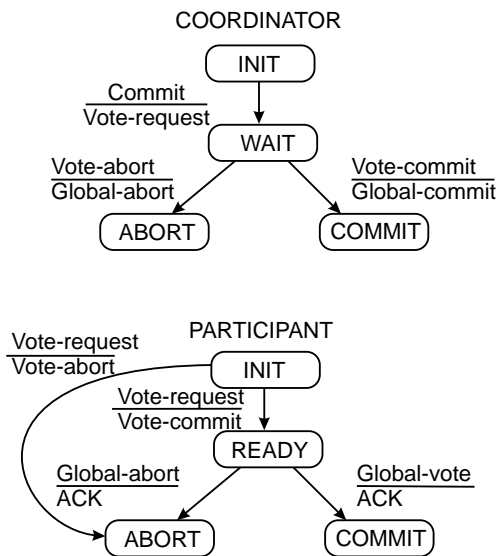
**Phase 1b:** When participant receives VOTE\_REQUEST it returns either YES or NO to coordinator. If it sends NO, it aborts its local computation

**Phase 2a:** Coordinator collects all votes; if all are YES, it sends COMMIT to all participants, otherwise it sends ABORT

07-25

## Fault Tolerance/7.5 Distributed Commit

## Two-Phase Commit (2/2)



07 – 26

Fault Tolerance/7.5 Distributed Commit

## 2PC – Failing Participant

**Observation:** Consider participant crash in one of its states, and the subsequent recovery to that state:

**Initial state:** No problem, as participant was unaware of the protocol

**Ready state:** Participant is waiting to either commit or abort. After recovery, participant needs to know which state transition it should make  $\Rightarrow$  log the coordinator's decision

**Abort state:** Merely make entry into abort state *idempotent*, e.g., removing the workspace of results

**Commit state:** Also make entry into commit state *idempotent*, e.g., copying workspace to storage.

**Observation:** When distributed commit is required, having participants use temporary workspaces to keep their results allows for simple recovery in the presence of failures.

07 – 27

Fault Tolerance/7.5 Distributed Commit

## 2PC – Failing Coordinator

**Observation:** The real problem lies in the fact that the coordinator's final decision may not be available for some time (or actually lost)

**Alternative:** Let a participant  $P$  in the ready state timeout when it hasn't received the coordinator's decision;  $P$  tries to find out what other participants know.

**Question:** Can  $P$  *not* succeed in getting the required information?

**Observation:** Essence of the problem is that a recovering participant cannot make a **local** decision: it is dependent on other (possibly failed) processes

## Three-Phase Commit (1/2)

**Phase 1a:** Coordinator sends VOTE\_REQUEST to participants

**Phase 1b:** When participant receives VOTE\_REQUEST it returns either YES or NO to coordinator. If it sends NO, it aborts its local computation

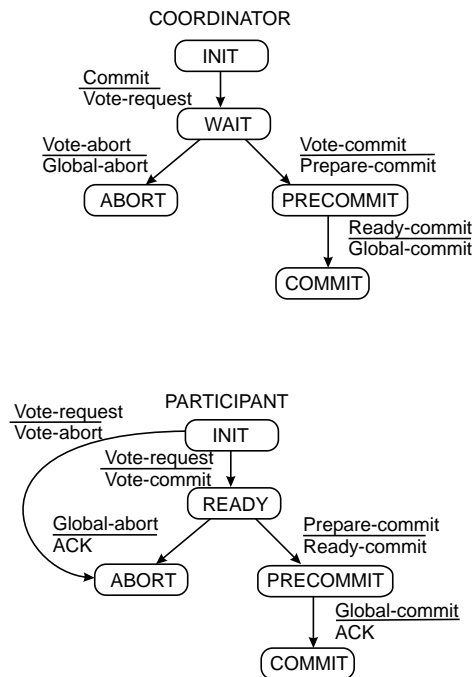
**Phase 2a:** Coordinator collects all votes; if all are YES, it sends PREPARE to all participants, otherwise it sends ABORT, and halts

**Phase 2b:** Each participant waits for PREPARE, or waits for ABORT after which it halts

**Phase 3a:** (Prepare to commit) Coordinator waits until all participants have ACKed receipt of PREPARE message, and then sends COMMIT to all

**Phase 3b:** (Prepare to commit) Participant waits for COMMIT

## Three-Phase Commit (2/2)



07 – 30

Fault Tolerance/7.5 Distributed Commit

## 3PC – Failing Participant

**Basic issue:** Can  $P$  find out what it should do after crashing in the ready or pre-commit state, even if other participants or the coordinator failed?

**Essence:** Coordinator and participants on their way to commit, never differ by more than one state transition

**Consequence:** If a participant timeouts in ready state, it can find out at the coordinator or other participants whether it should abort, or enter pre-commit state

**Observation:** If a participant already made it to the pre-commit state, it can always safely commit (but is not allowed to do so for the sake of failing other processes)

**Observation:** We may need to elect another coordinator to send off the final COMMIT

07 – 31

Fault Tolerance/7.5 Distributed Commit



## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

- 
- This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

- 
- This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

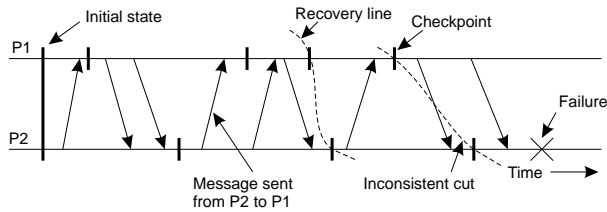
This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

# Consistent Recovery State

**Requirement:** Every message that has been received is also shown to have been sent in the state of the sender

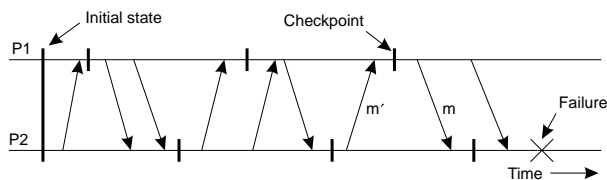
**Recovery line:** Assuming processes regularly **checkpoint** their state, the most recent **consistent global checkpoint**.



**Observation:** If and only if the system provides *reliable* communication, should sent messages also be received in a consistent state

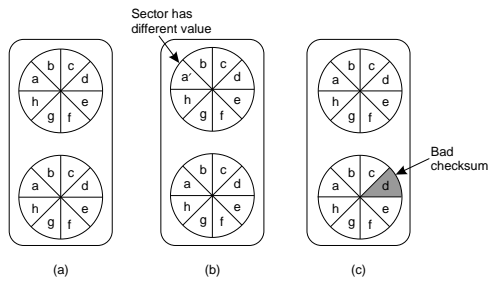
## Cascaded Rollback

**Observation:** If checkpointing is done at the “wrong” instants, the recovery line may lie at system startup time ⇒ **cascaded rollback**



# Checkpointing: Stable Storage

**Principle:** Replicate all data on at least two disks, and keep one copy “correct” at all times.



## After a crash:

- If both disks are identical: you're in good shape.
- If one is bad, but the other is okay (checksums): choose the good one.
- If both seem okay, but are different: choose the main disk.
- If both aren't good: you're **not** in a good shape.

## Independent Checkpointing

**Essence:** Each process independently takes checkpoints, with the risk that a cascaded rollback to system startup.

- Let  $CP[i](m)$  denote  $m^{\text{th}}$  checkpoint of process  $P_i$  and  $INT[i](m)$  the interval between  $CP[i](m-1)$  and  $CP[i](m)$
- When process  $P_i$  sends a message in interval  $INT[i](m)$ , it piggybacks  $(i, m)$
- When process  $P_j$  receives a message in interval  $INT[j](n)$ , it records the dependency  $INT[i](m) \rightarrow INT[j](n)$
- The dependency  $INT[i](m) \rightarrow INT[j](n)$  is saved to stable storage when taking checkpoint  $CP[j](n)$

**Observation:** If process  $P_i$  rolls back to  $CP[i](m-1)$ ,  $P_j$  must roll back to  $CP[j](n-1)$ . **Question:** How can  $P_j$  find out where to roll back to?

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

---

### Fault Tolerance/Recovery

## This image shows a single sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

---

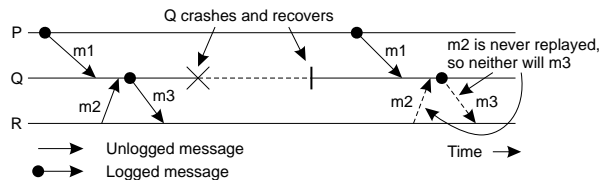
### Fault Tolerance/Recovery

# Message Logging and Consistency

**Problem:** When should we actually log messages?

**Issue:** Avoid **orphans**:

- Process  $Q$  has just received and subsequently delivered messages  $m_1$  and  $m_2$
- Assume that  $m_2$  is never logged.
- After delivering  $m_1$  and  $m_2$ ,  $Q$  sends message  $m_3$  to process  $R$
- Process  $R$  receives and subsequently delivers  $m_3$



**Goal:** Devise message logging schemes in which orphans do not occur

07 – 40

Fault Tolerance/Recovery

## Message-Logging Schemes (1/2)

**HDR[m]:** The header of message  $m$  containing its source, destination, sequence number, and delivery number

The header contains all information for resending a message and delivering it in the correct order (assume data is reproduced by the application)

A message  $m$  is **stable** if HDR[m] cannot be lost (e.g., because it has been written to stable storage)

**DEP[m]:** The set of processes to which message  $m$  has been delivered, as well as any message that causally depends on delivery of  $m$

**COPY[m]:** The set of processes that have a copy of HDR[m] in their volatile memory

If  $C$  is a collection of crashed processes, then  $Q \notin C$  is an orphan if there is a message  $m$  such that  $Q \in \text{DEP}[m]$  and  $\text{COPY}[m] \subseteq C$

07 – 41

Fault Tolerance/Recovery

[illegible][illegible][illegible][illegible][illegible][illegible]