

## Fault Tolerance

### Chapter 7

---

---

---

---

---

---

---

---

## Basic Concepts

Dependability Includes

- Availability
- Reliability
- Safety
- Maintainability

---

---

---

---

---

---

---

---

## Failure Models

Type of failure	Description
Crash failure	A server halts, but is working correctly until it halts
Omission failure <i>Receive omission</i> <i>Send omission</i>	A server fails to respond to incoming requests A server fails to receive incoming messages A server fails to send messages
Timing failure	A server's response lies outside the specified time interval
Response failure <i>Value failure</i> <i>State transition failure</i>	The server's response is incorrect The value of the response is wrong The server deviates from the correct flow of control
Arbitrary failure	A server may produce arbitrary responses at arbitrary times

Different types of failures.

---

---

---

---

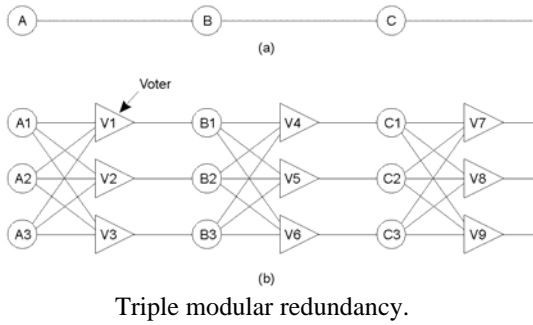
---

---

---

---

## Failure Masking by Redundancy




---

---

---

---

---

---

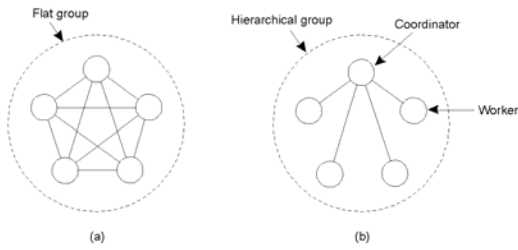
---

---

---

---

## Flat Groups versus Hierarchical Groups



- a) Communication in a flat group.
- b) Communication in a simple hierarchical group

---

---

---

---

---

---

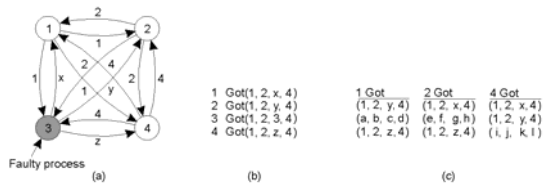
---

---

---

---

## Agreement in Faulty Systems (1)



- The Byzantine generals problem for 3 loyal generals and 1 traitor.
- a) The generals announce their troop strengths (in units of 1 kilosoldiers).
  - b) The vectors that each general assembles based on (a)
  - c) The vectors that each general receives in step 3.

---

---

---

---

---

---

---

---

---

---

## Agreement in Faulty Systems (2)



The same as in previous slide, except now with 2 loyal generals and one traitor.

---

---

---

---

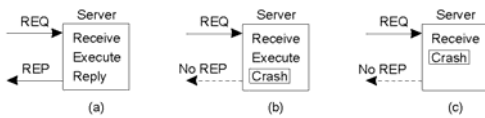
---

---

---

---

## Lost Request Messages Server Crashes (1)



A server in client-server communication

- a) Normal case
- b) Crash after execution
- c) Crash before execution

---

---

---

---

---

---

---

---

## Server Crashes (2)

Client	Server					
	Strategy M → P			Strategy P → M		
Reissue strategy	MPC	MC(P)	C(MP)	PMC	PC(M)	C(PM)
Always	DUP	OK	OK	DUP	DUP	OK
Never	OK	ZERO	ZERO	OK	OK	ZERO
Only when ACKed	DUP	OK	ZERO	DUP	OK	ZERO
Only when not ACKed	OK	ZERO	OK	OK	DUP	OK

Different combinations of client and server strategies in the presence of server crashes.

---

---

---

---

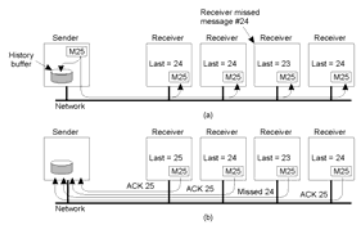
---

---

---

---

## Basic Reliable-Multicasting Schemes



A simple solution to reliable multicasting when all receivers are known and are assumed not to fail

- a) Message transmission
- b) Reporting feedback

---

---

---

---

---

---

---

---

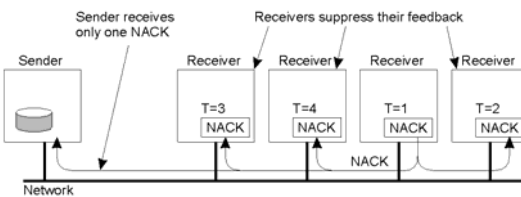
---

---

---

---

## Nonhierarchical Feedback Control



Several receivers have scheduled a request for retransmission, but the first retransmission request leads to the suppression of others.

---

---

---

---

---

---

---

---

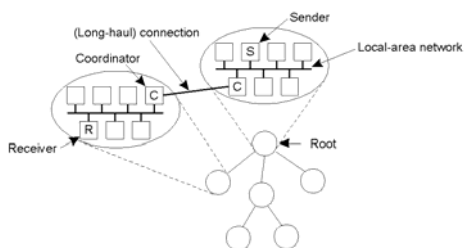
---

---

---

---

## Hierarchical Feedback Control



The essence of hierarchical reliable multicasting.

- a) Each local coordinator forwards the message to its children.
- b) A local coordinator handles retransmission requests.

---

---

---

---

---

---

---

---

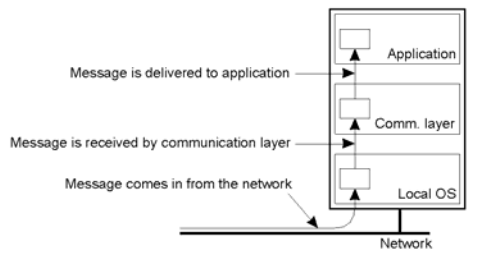
---

---

---

---

## Virtual Synchrony (1)



The logical organization of a distributed system to distinguish between message receipt and message delivery

---

---

---

---

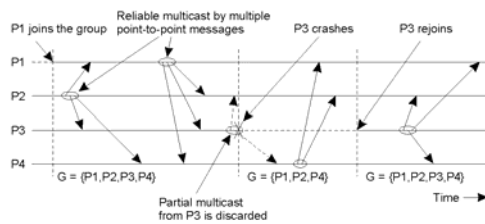
---

---

---

---

## Virtual Synchrony (2)



The principle of virtual synchronous multicast.

---

---

---

---

---

---

---

---

## Message Ordering (1)

Process P1	Process P2	Process P3
sends m1	receives m1	receives m2
sends m2	receives m2	receives m1

Three communicating processes in the same group.  
The ordering of events per process is shown along the vertical axis.

---

---

---

---

---

---

---

---

## Message Ordering (2)

Process P1	Process P2	Process P3	Process P4
sends m1	receives m1	receives m3	sends m3
sends m2	receives m3	receives m1	sends m4
	receives m2	receives m2	
	receives m4	receives m4	

Four processes in the same group with two different senders, and a possible delivery order of messages under FIFO-ordered multicasting

---

---

---

---

---

---

---

---

## Implementing Virtual Synchrony (1)

Multicast	Basic Message Ordering	Total-ordered Delivery?
Reliable multicast	None	No
FIFO multicast	FIFO-ordered delivery	No
Causal multicast	Causal-ordered delivery	No
Atomic multicast	None	Yes
FIFO atomic multicast	FIFO-ordered delivery	Yes
Causal atomic multicast	Causal-ordered delivery	Yes

Six different versions of virtually synchronous reliable multicasting.

---

---

---

---

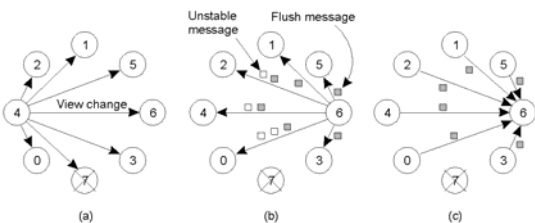
---

---

---

---

## Implementing Virtual Synchrony (2)



- a) Process 4 notices that process 7 has crashed, sends a view change
- b) Process 6 sends out all its unstable messages, followed by a flush message
- c) Process 6 installs the new view when it has received a flush message from everyone else

---

---

---

---

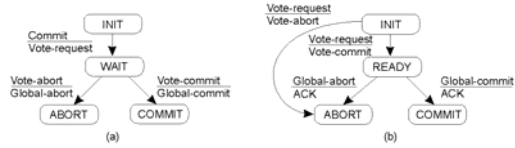
---

---

---

---

## Two-Phase Commit (1)



- a) The finite state machine for the coordinator in 2PC.
- b) The finite state machine for a participant.

---

---

---

---

---

---

---

---

## Two-Phase Commit (2)

State of Q	Action by P
COMMIT	Make transition to COMMIT
ABORT	Make transition to ABORT
INIT	Make transition to ABORT
READY	Contact another participant

Actions taken by a participant *P* when residing in state *READY* and having contacted another participant *Q*.

---

---

---

---

---

---

---

---

## Two-Phase Commit (3)

```

actions by coordinator:
while START_2PC to local log;
multicast VOTE_REQUEST to all participants;
while not all votes have been collected {
  wait for any incoming vote;
  if timeout {
    while GLOBAL_ABORT to local log;
    multicast GLOBAL_ABORT to all participants;
    exit;
  }
  record vote;
}
if all participants sent VOTE_COMMIT and coordinator votes COMMIT{
  write GLOBAL_COMMIT to local log;
  multicast GLOBAL_COMMIT to all participants;
} else {
  write GLOBAL_ABORT to local log;
  multicast GLOBAL_ABORT to all participants;
}
  
```

Outline of the steps taken by the coordinator in a two phase commit protocol

---

---

---

---

---

---

---

---

## Two-Phase Commit (4)

Steps taken by participant process in 2PC.

```

actions by participant:
write INIT to local log;
wait for VOTE_REQUEST from coordinator;
if timeout {
  write VOTE_ABORT to local log;
  exit;
}
if participant votes COMMIT {
  write VOTE_COMMIT to local log;
  send VOTE_COMMIT to coordinator;
  wait for DECISION from coordinator;
  if timeout {
    multicast DECISION_REQUEST to other participants;
    wait until DECISION is received; /* remain blocked */
    write DECISION to local log;
  }
  if DECISION == GLOBAL_COMMIT
    write GLOBAL_COMMIT to local log;
  else if DECISION == GLOBAL_ABORT
    write GLOBAL_ABORT to local log;
} else {
  write VOTE_ABORT to local log;
  send VOTE_ABORT to coordinator;
}
  
```

---

---

---

---

---

---

---

---

---

---

---

---

## Two-Phase Commit (5)

```

actions for handling decision requests: /* executed by separate thread */
while true {
  wait until any incoming DECISION_REQUEST is received; /* remain blocked */
  read most recently recorded STATE from the local log;
  if STATE == GLOBAL_COMMIT
    send GLOBAL_COMMIT to requesting participant;
  else if STATE == INIT or STATE == GLOBAL_ABORT
    send GLOBAL_ABORT to requesting participant;
  else
    skip; /* participant remains blocked */
}
  
```

Steps taken for handling incoming decision requests.

---

---

---

---

---

---

---

---

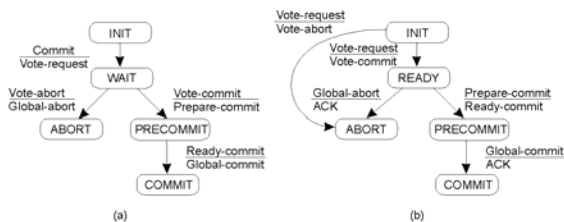
---

---

---

---

## Three-Phase Commit



- a) Finite state machine for the coordinator in 3PC
- b) Finite state machine for a participant

---

---

---

---

---

---

---

---

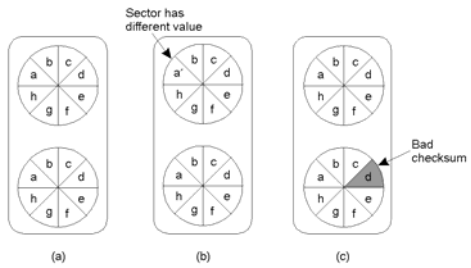
---

---

---

---

## Recovery Stable Storage



- a) Stable Storage
- b) Crash after drive 1 is updated
- c) Bad spot

---

---

---

---

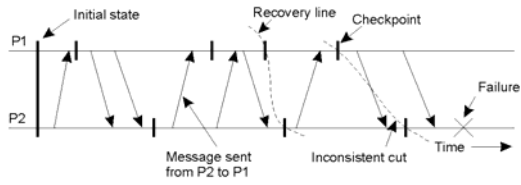
---

---

---

---

## Checkpointing



A recovery line.

---

---

---

---

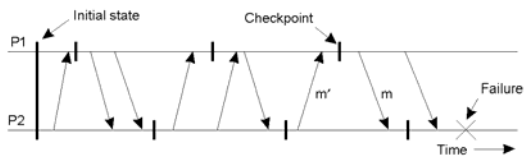
---

---

---

---

## Independent Checkpointing



The domino effect.

---

---

---

---

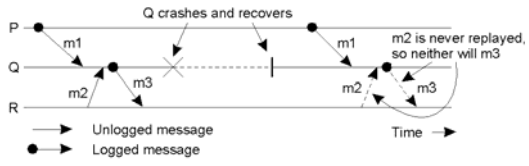
---

---

---

---

## Message Logging



Incorrect replay of messages after recovery,  
leading to an orphan process.

---

---

---

---

---

---

---

---