

Classes of process failures

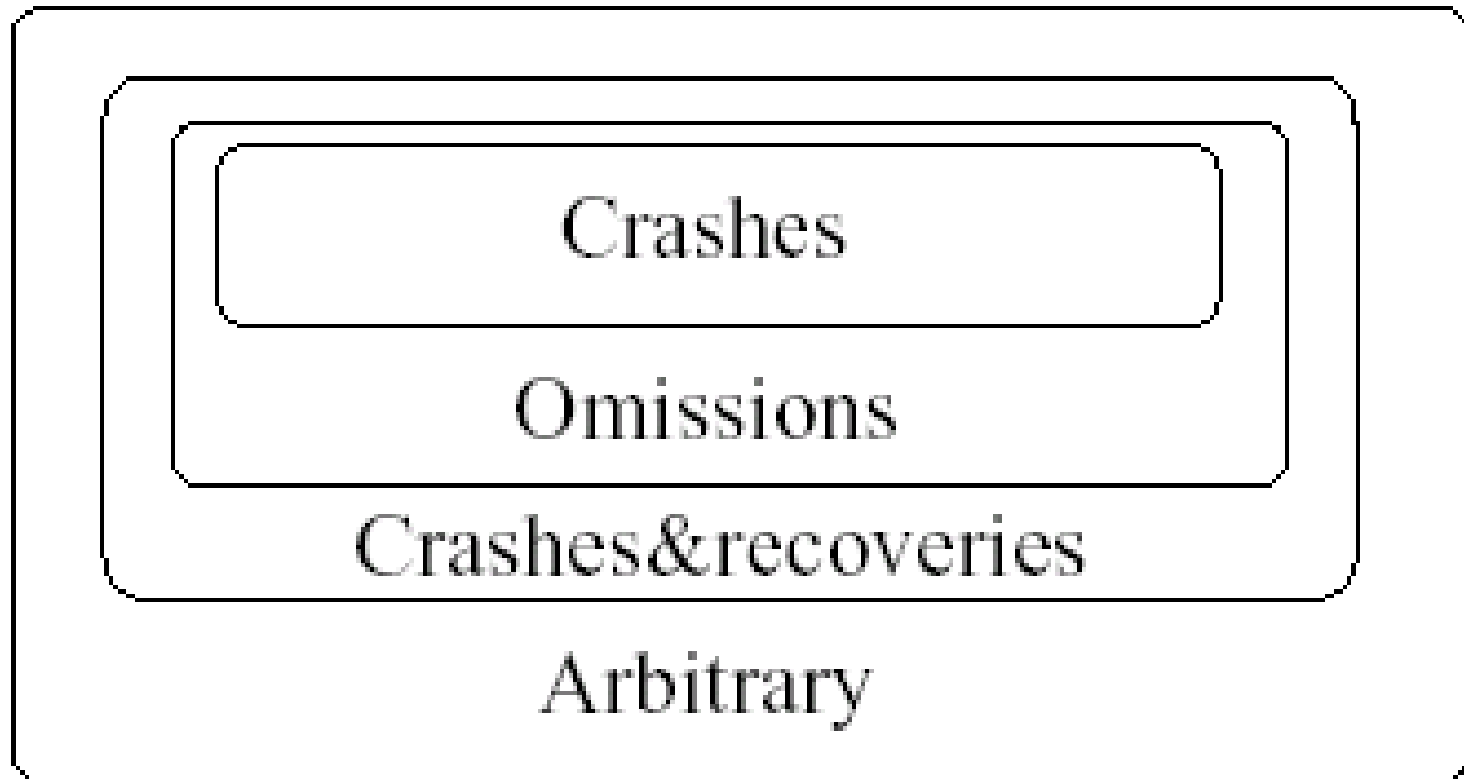


Figure 2.2. Failure modes of a process

Module:

Name: FairLossPointToPointLinks (flp2p).

Events:

Request: $\langle \text{flp2pSend}, \text{dest}, m \rangle$: Used to request the transmission of message m to process dest .

Indication: $\langle \text{flp2pDeliver}, \text{src}, m \rangle$: Used to deliver message m sent by process src .

Properties:

FLL1: *Fair loss:* If a message m is sent infinitely often by process p_i to process p_j , and neither p_i nor p_j crash, then m is delivered an infinite number of times by p_j .

FLL2: *Finite duplication:* If a message m is sent a finite number of times by process p_i to process p_j , then m cannot be delivered an infinite number of times by p_j .

FLL3: *No creation:* If a message m is delivered by some process p_j , then m has been previously sent to p_j by some process p_i .

Module 2.1 Interface and properties of fair-lossy point-to-point links.

Module:

Name: StubbornPointToPointLink (sp2p).

Events:

Request: $\langle \text{sp2pSend}, \text{dest}, m \rangle$: Used to request the transmission of message m to process dest .

Indication: $\langle \text{sp2pDeliver}, \text{src}, m \rangle$: Used to deliver message m sent by process src .

Properties:

SL1: *Stubborn delivery:* Let p_i be any process that sends a message m to a correct process p_j . If p_i does not crash, then p_j delivers m an infinite number of times.

SL2: *No creation:* If a message m is delivered by some process p_j , then m was previously sent to p_j by some process p_i .

Module 2.2 Interface and properties of stubborn point-to-point links.

Algorithm 2.1 Stubborn links using fair-loss links.

Implements:

StubbornPointToPointLink (sp2p).

Uses:

FairLossPointToPointLinks (flp2p).

upon event $\langle \textit{sp2pSend}, \textit{dest}, m \rangle$ **do**
 while (true) **do**
 trigger $\langle \textit{flp2pSend}, \textit{dest}, m \rangle$;

upon event $\langle \textit{flp2pDeliver}, \textit{src}, m \rangle$ **do**
 trigger $\langle \textit{sp2pDeliver}, \textit{src}, m \rangle$;

Module:

Name: PerfectPointToPointLink (pp2p).

Events:

Request: $\langle pp2pSend, dest, m \rangle$: Used to request the transmission of message m to process $dest$.

Indication: $\langle pp2pDeliver, src, m \rangle$: Used to deliver message m sent by process src .

Properties:

PL1: *Reliable delivery*: Let p_i be any process that sends a message m to a process p_j . If neither p_i nor p_j crashes, then p_j eventually delivers m .

PL2: *No duplication*: No message is delivered by a process more than once.

PL3: *No creation*: If a message m is delivered by some process p_j , then m was previously sent to p_j by some process p_i .

Module 2.3 Interface and properties of perfect point-to-point links.

Algorithm 2.2 Perfect links using stubborn links.

Implements:

PerfectPointToPointLinks (pp2p).

Uses:

StubbornPointToPointLinks (sp2p).

upon event $\langle \textit{Init} \rangle$ **do**

$\textit{delivered} := \emptyset$;

upon event $\langle \textit{pp2pSend}, \textit{dest}, m \rangle$ **do**

trigger $\langle \textit{sp2pSend}, \textit{dest}, m \rangle$;

upon event $\langle \textit{sp2pDeliver}, \textit{src}, m \rangle$ **do**

if $m \notin \textit{delivered}$ **then**

$\textit{delivered} := \textit{delivered} \cup \{m\}$;

trigger $\langle \textit{pp2pDeliver}, \textit{src}, m \rangle$;

Partially synchronous systems

- System that eventually is synchronous (without stating when exactly).
- System that might not always be synchronous and there is no bound on the period during which it is asynchronous.
- System characterized by periods during which it is synchronous, and some of these periods are long enough for an algorithm to terminate its execution.

Module:

Name: PerfectFailureDetector (\mathcal{P}).

Events:

Indication: $\langle \text{crash}, p_i \rangle$: Used to notify that process p_i has crashed.

Properties:

PFD1: *Strong completeness:* Eventually every process that crashes is permanently detected by every correct process.

PFD2: *Strong accuracy:* No process is detected by any process before it crashes.

Module 2.4 Interface and properties of the perfect failure detector.

Algorithm 2.3 Perfect failure detector with perfect links and timeouts.

Implements:

PerfectFailureDetector (\mathcal{P}).

Uses:

PerfectPointToPointLinks (pp2p).

upon event $\langle \textit{Init} \rangle$ **do**

$\textit{alive} := \Pi$;

$\textit{suspected} := \emptyset$;

upon event $\langle \textit{TimeDelay} \rangle$ **do**

$\forall p_i \in \Pi :$

if $p_i \notin \textit{alive}$ **and** $p_i \notin \textit{suspected}$ **then**

$\textit{suspected} := \textit{suspected} \cup \{p_i\}$;

trigger $\langle \textit{crash}, p_i \rangle$;

$\textit{alive} := \emptyset$;

$\forall p_i \in \Pi : \text{trigger} \langle \textit{pp2pSend}, p_i, [\textit{DATA}, \textit{heartbeat}] \rangle$;

upon event $\langle \textit{pp2pDeliver}, \textit{src}, [\textit{DATA}, \textit{heartbeat}] \rangle$ **do**

$\textit{alive} := \textit{alive} \cup \{\textit{src}\}$;

Module:

Name: EventuallyPerfectFailureDetector ($\diamond P$).

Events:

Indication: $\langle suspect, p_i \rangle$: Used to notify that process p_i is suspected to have crashed.

Indication: $\langle restore, p_i \rangle$: Used to notify that process p_i is not suspected anymore.

Properties:

EPFD1: *Eventual strong completeness:* Eventually, every process that crashes is permanently suspected by every correct process.

EPFD2: *Eventual strong accuracy:* Eventually, no correct process is suspected by any correct process.

Module 2.5 Interface and properties of the eventually perfect failure detector.

Algorithm 2.4 Eventually perfect failure detector with perfect links and timeouts.

Implements:

EventuallyPerfectFailureDetector ($\diamond\mathcal{P}$).

Uses:

PerfectPointToPointLinks (pp2p).

upon event $\langle \textit{Init} \rangle$ **do**

$\textit{alive} := \Pi$;

$\textit{suspected} := \emptyset$;

upon event $\langle \textit{TimeDelay} \rangle$ **do**

$\forall p_i \in \Pi$:

if $p_i \notin \textit{alive}$ **then**

$\textit{suspected} := \textit{suspected} \cup \{p_i\}$;

trigger $\langle \textit{crash}, p_i \rangle$;

else

if $p_i \in \textit{suspected}$ **then**

$\textit{suspected} := \textit{suspected} \setminus \{p_i\}$;

$\textit{TimeDelay} := \textit{TimeDelay} + \Delta$;

trigger $\langle \textit{restore}, p_i \rangle$;

$\textit{alive} := \emptyset$;

$\forall p_i \in \Pi$: **trigger** $\langle \textit{pp2pSend}, p_i, [\textit{DATA}, \textit{heartbeat}] \rangle$;

upon event $\langle \textit{pp2pDeliver}, \textit{src}, [\textit{DATA}, \textit{heartbeat}] \rangle$ **do**

$\textit{alive} := \textit{alive} \cup \{\textit{src}\}$;

Module:

Name: EventualLeaderDetector (Ω).

Events:

Indication: $\langle trust, p_i \rangle$: Used to notify that process p_i is trusted to be leader.

Properties:

CD1: *Eventual accuracy*: There is a time after which every correct process trusts some correct process.

CD2: *Eventual agreement*: There is a time after which no two correct processes trust different processes.

Module 2.6 Interface and properties of the eventual leader detector.

Algorithm 2.5 Eventually leader election with crash-recovery processes, stubborn links and timeouts .

Implements:

EventualLeaderDetector (Ω).

Uses:

StubbornPointToPointLinks (flp2p).

upon event $\langle \text{Init} \rangle$ **do**

 leader := p_1 ;
 possible := Π ;
 epoch := 0;

upon event $\langle \text{Recovery} \rangle$ **do**

 retrieve(epoch);
 epoch := epoch + 1;
 store(epoch);

upon event $\langle \text{TimeDelay} \rangle$ **do**

if leader \neq select(possible) **then**
 TimeDelay := TimeDelay + Δ ;
 leader := select(possible);
 if {leader} $\neq \emptyset$ **then**
 trigger $\langle \text{trust}, \text{leader} \rangle$;
 possible := \emptyset ;
 $\forall p_i \in \Pi$: **trigger** $\langle \text{sp2pSend}, p_i, [\text{DATA}, \text{heartbeat}, \text{epoch}] \rangle$;

upon event $\langle \text{flp2pDeliver}, \text{src}, [\text{DATA}, \text{heartbeat}, \text{epc}] \rangle$ **do**

 possible := possible $\cup \{(\text{src}, \text{epc})\}$;

Specific system models

- **Fail-stop model.** Processes execute deterministic algorithms, unless they possibly crash and stop executing any computation, links are perfect, and there available the perfect failure detector.
- **Fail-silent model.** Processes execute deterministic algorithms, unless they possibly crash and stop executing any computation, links are perfect.

Specific system models

- **Fail-noisy model.** Processes execute deterministic algorithms, unless they possibly crash and stop executing any computation, links are perfect, and moreover there available the eventually perfect failure detector.
- **Randomized model.** Processes may use a random oracle to choose among several steps to execute.

Module:

Name: BestEffortBroadcast (beb).

Events:

Request: $\langle \text{bebBroadcast}, m \rangle$: Used to broadcast message m to all processes.

Indication: $\langle \text{bebDeliver}, \text{src}, m \rangle$: Used to deliver message m broadcast by process src .

Properties:

BEB1: *Best-effort validity:* For any two processes p_i and p_j . If p_i and p_j are correct, then every message broadcast by p_i is eventually delivered by p_j .

BEB2: *No duplication:* No message is delivered more than once.

BEB3: *No creation:* If a message m is delivered by some process p_j , then m was previously broadcast by some process p_i .

Module 3.1 Interface and properties of best-effort broadcast.

Algorithm 3.1 Basic Broadcast.

Implements:

BestEffortBroadcast (beb).

Uses:

PerfectPointToPointLinks (pp2p).

upon event $\langle \text{bebBroadcast}, m \rangle$ **do**

$\forall_{p_i \in \Pi} :$

trigger $\langle \text{pp2pSend}, p_i, m \rangle;$

upon event $\langle \text{pp2pDeliver}, p_i, m \rangle$ **do**

trigger $\langle \text{bebDeliver}, p_i, m \rangle;$

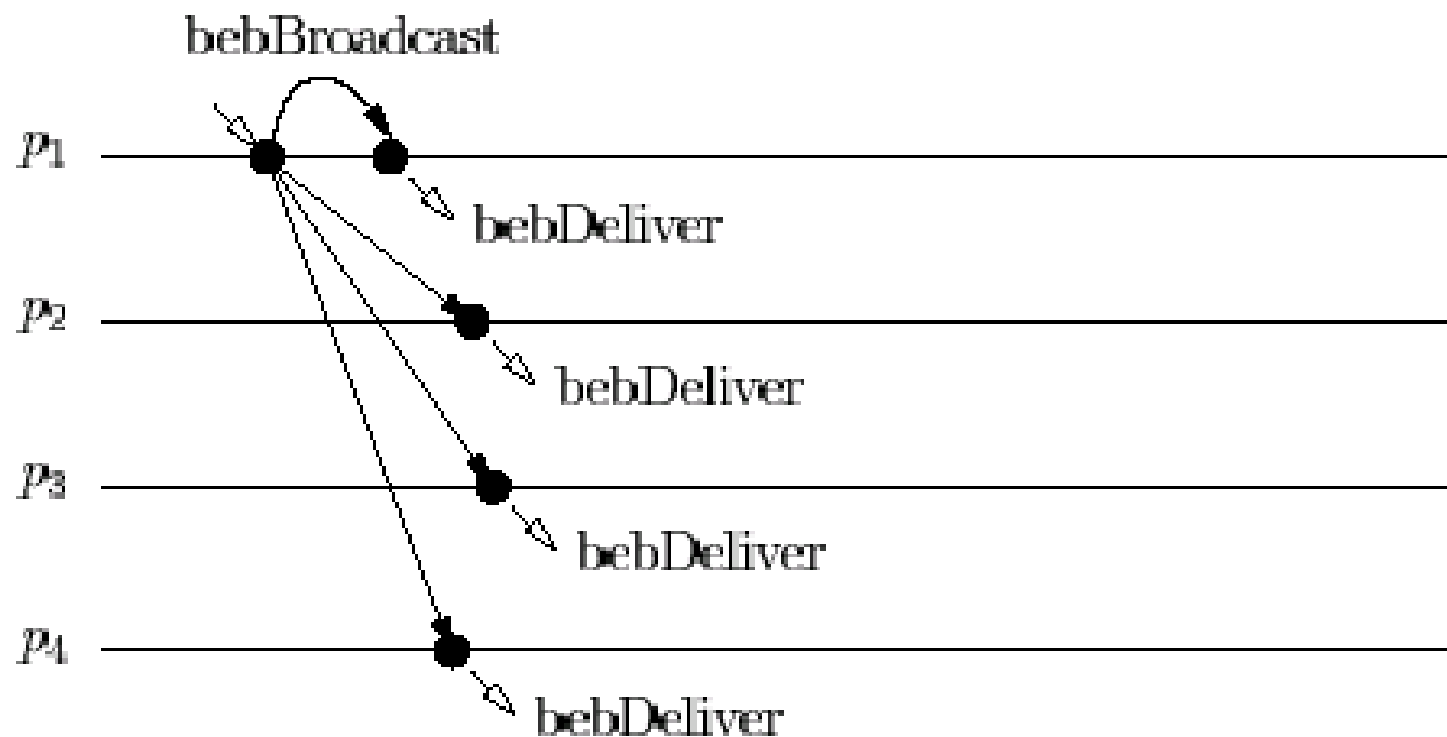


Figure 3.1. Sample execution of Basic Broadcast algorithm.

Module:

Name: (regular)ReliableBroadcast (rb).

Events:

Request: $\langle rbBroadcast, m \rangle$: Used to broadcast message m .

Indication: $\langle rbDeliver, src, m \rangle$: Used to deliver message m broadcast by process src .

Properties:

RB1: *Validity*: If a correct process p_i broadcasts a message m , then p_i eventually delivers m .

RB2: *No duplication*: No message is delivered more than once.

RB3: *No creation*: If a message m is delivered by some process p_j , then m was previously broadcast by some process p_i .

RB4: *Agreement*: If a message m is delivered by some correct process p_i , then m is eventually delivered by every correct process p_j .

Module 3.2 Interface and properties of reliable broadcast.

Algorithm 3.2 Lazy reliable broadcast.

Implements:

ReliableBroadcast (rb).

Uses:

BestEffortBroadcast (beb).

PerfectFailureDetector (\mathcal{P}).**upon event** $\langle \text{Init} \rangle$ **do**delivered $:= \emptyset$;correct $:= \Pi$; $\forall p_i \in \Pi : \text{from}[p_i] := \emptyset$;**upon event** $\langle \text{rbBroadcast}, m \rangle$ **do**trigger $\langle \text{bebBroadcast}, [\text{DATA}, \text{self}, m] \rangle$;**upon event** $\langle \text{bebDeliver}, p_i, [\text{DATA}, s_m, m] \rangle$ **do**if $m \notin \text{delivered}$ thendelivered $:= \text{delivered} \cup \{m\}$ trigger $\langle \text{rbDeliver}, s_m, m \rangle$; $\text{from}[p_i] := \text{from}[p_i] \cup \{[s_m, m]\}$ if $p_i \notin \text{correct}$ thentrigger $\langle \text{bebBroadcast}, [\text{DATA}, s_m, m] \rangle$;**upon event** $\langle \text{crash}, p_i \rangle$ **do**correct $:= \text{correct} \setminus \{p_i\}$ **forall** $[s_m, m] \in \text{from}[p_i]$: **do**trigger $\langle \text{bebBroadcast}, [\text{DATA}, s_m, m] \rangle$;

Algorithm 3.3 Eager reliable broadcast.

Implements:

ReliableBroadcast (rb).

Uses:

BestEffortBroadcast (beb).

upon event $\langle \textit{Init} \rangle$ **do**delivered $:= \emptyset$;**upon event** $\langle \textit{rbBroadcast}, m \rangle$ **do**delivered $:= \text{delivered} \cup \{m\}$ **trigger** $\langle \textit{rbDeliver}, \textit{self}, m \rangle$;**trigger** $\langle \textit{bebBroadcast}, [\textit{DATA}, \textit{self}, m] \rangle$;**upon event** $\langle \textit{bebDeliver}, p_i, [\textit{DATA}, s_m, m] \rangle$ **do****if** $m \notin \text{delivered}$ **do**delivered $:= \text{delivered} \cup \{m\}$ **trigger** $\langle \textit{rbDeliver}, s_m, m \rangle$;**trigger** $\langle \textit{bebBroadcast}, [\textit{DATA}, s_m, m] \rangle$;

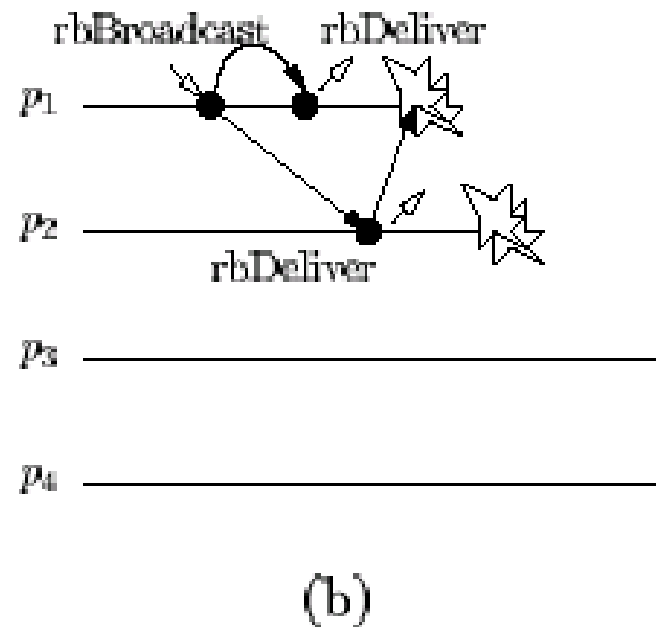
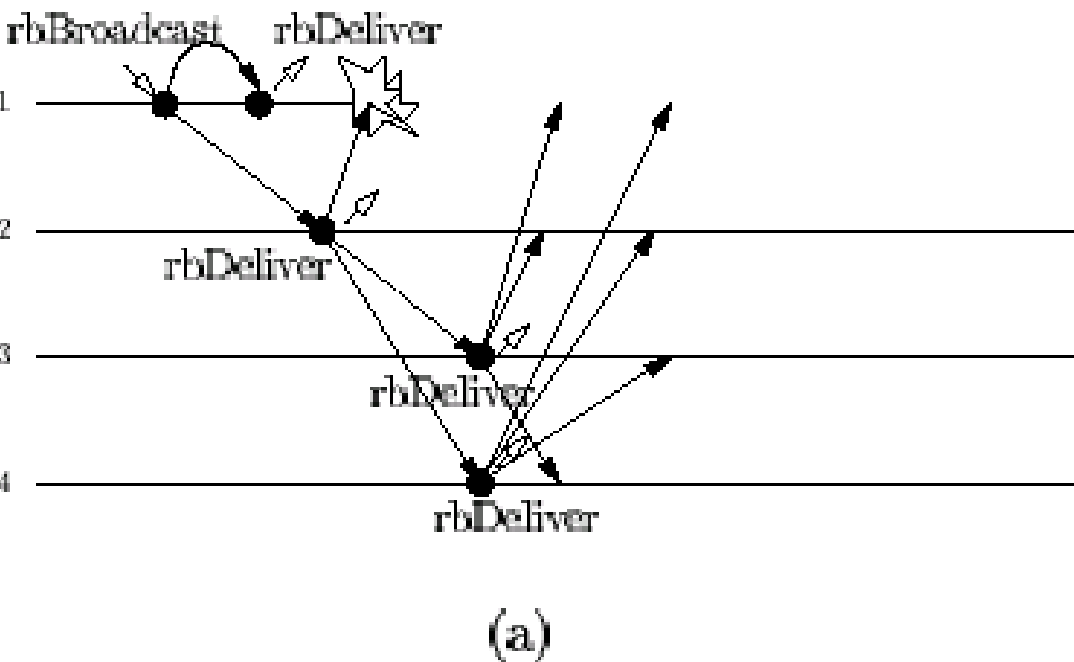


Figure 3.2. Sample executions of eager reliable broadcast.

Module:

Name: UniformReliableBroadcast (urb).

Events:

$\langle \text{urbBroadcast}, m \rangle$, $\langle \text{urbDeliver}, \text{src}, m \rangle$, with the same meaning and interface as in regular reliable broadcast.

Properties:

RB1-RB3: Same as in regular reliable broadcast.

URB4: *Uniform Agreement:* If a message m is delivered by some process p_i (whether correct or faulty), then m is also eventually delivered by every other correct process p_j .

Module 3.3 Interface and properties of uniform reliable broadcast.

Algorithm 3.4 All-ack uniform reliable broadcast.

Implements:

UniformReliableBroadcast (urb).

Uses:

BestEffortBroadcast (beb).

PerfectFailureDetector (\mathcal{P}).

function canDeliver(m) **returns** boolean **is**
 return ($\text{correct} \subset \text{ack}_m$) \wedge ($m \notin \text{delivered}$);

upon event $\langle \text{Init} \rangle$ **do**
 $\text{delivered} := \text{forward} := \emptyset$;
 $\text{correct} := \Pi$;
 $\text{ack}_m := \emptyset, \forall_m$;

upon event $\langle \text{urbBroadcast}, m \rangle$ **do**
 $\text{forward} := \text{forward} \cup \{m\}$
 trigger $\langle \text{bebBroadcast}, [\text{DATA}, \text{self}, m] \rangle$;

upon event $\langle \text{bebDeliver}, p_i, [\text{DATA}, s_m, m] \rangle$ **do**
 $\text{ack}_m := \text{ack}_m \cup \{p_i\}$
 if $m \notin \text{forward}$ **do**
 $\text{forward} := \text{forward} \cup \{m\}$;
 trigger $\langle \text{bebBroadcast}, [\text{DATA}, s_m, m] \rangle$;

upon event $\langle \text{crash}, p_i \rangle$ **do**
 $\text{correct} := \text{correct} \setminus \{p_i\}$;

upon (canDeliver(m)) **do**
 $\text{delivered} := \text{delivered} \cup \{m\}$;
 trigger $\langle \text{urbDeliver}, s_m, m \rangle$;

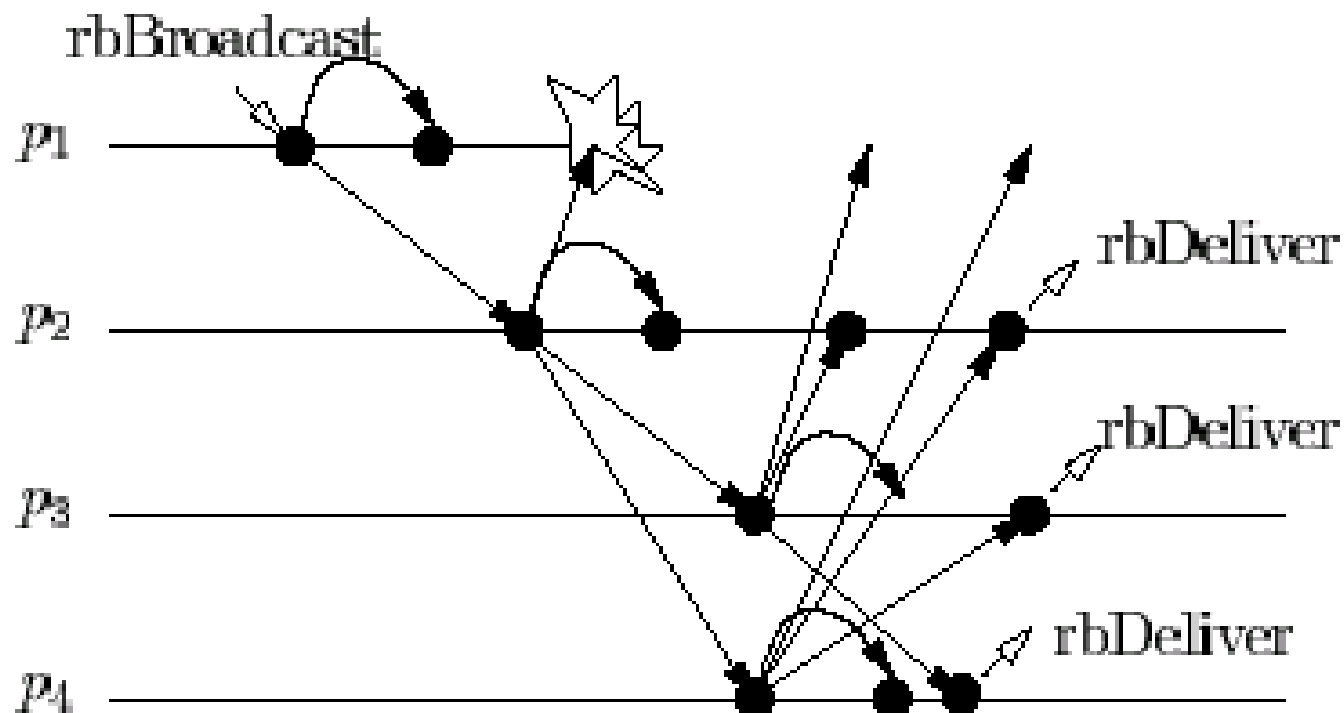


Figure 3.3. Sample execution of uniform reliable broadcast.

Algorithm 3.5 Majority-ack uniform reliable broadcast.

implements:

UniformReliableBroadcast (urb).

uses:

BestEffortBroadcast (beb).

function canDeliver(m) **returns** boolean **is**

return ($|\text{ack}_m| > N/2$) \wedge ($m \notin \text{delivered}$);

Except for the function above, and the non-use of the // perfect failure detector, same as Algorithm 3.

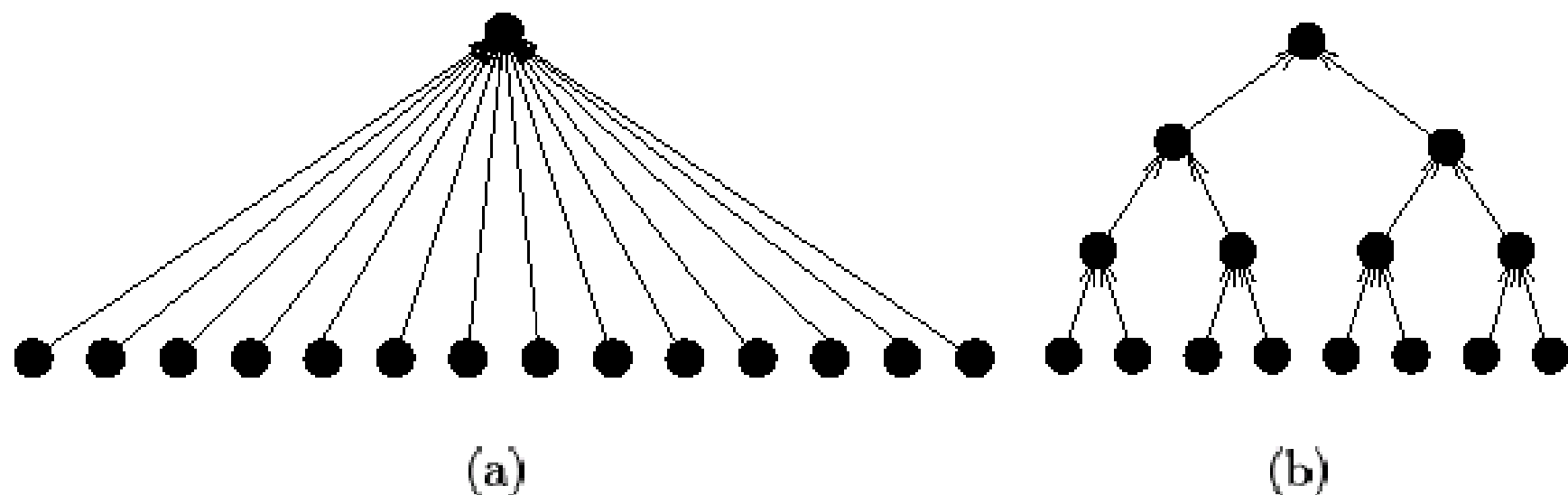


Figure 3.4. Ack implosion and ack tree.

Module:

Name: Probabilistic Broadcast (pb).

Events:

Request: $\langle pbBroadcast, m \rangle$: Used to broadcast message m to all processes.

Indication: $\langle pbDeliver, src, m \rangle$: Used to deliver message m broadcast by process src .

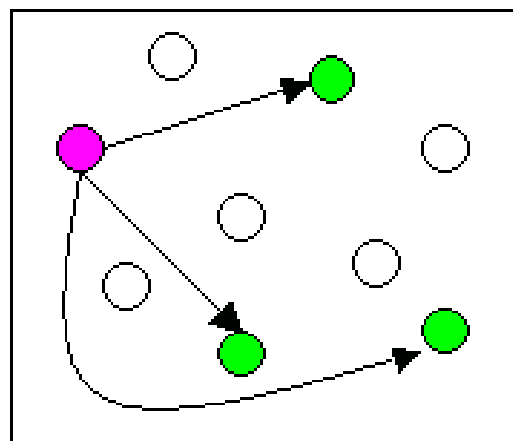
Properties:

PB1: *Probabilistic validity:* There is a given probability such that for any p_i and p_j that are correct, every message broadcast by p_i is eventually delivered by p_j with this probability.

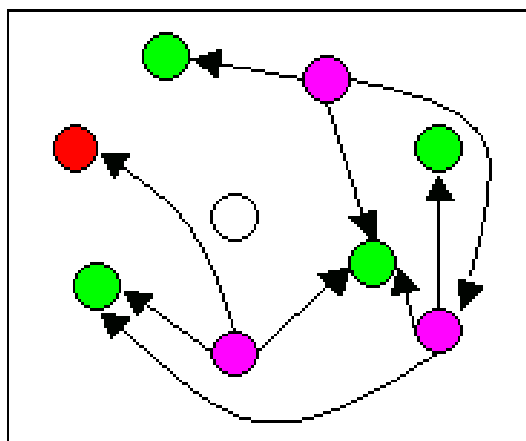
PB2: *No duplication:* No message is delivered more than once.

PB3: *No creation:* If a message m is delivered by some process p_j , then m was previously broadcast by some process p_i .

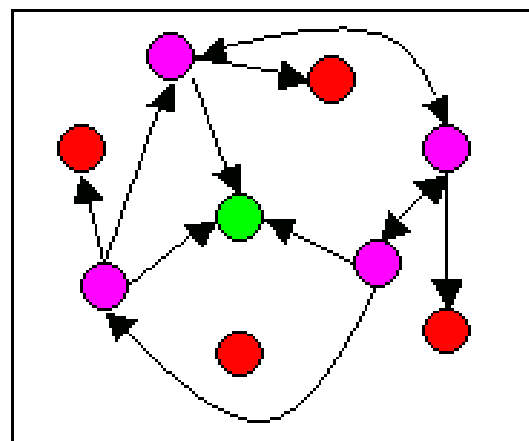
Module 3.7 Interface and properties of probabilistic broadcast.



(a)



(b)



(c)

Figure 3.5. Gossip Dissemination.

Algorithm 3.9 Eager Probabilistic Broadcast.

Implements:

ProbabilisticBroadcast (pb).

Uses:

unreliablePointToPointLinks (up2p).

upon event $\langle \textit{Init} \rangle$ **do**delivered $:= \emptyset$;**function** pick-targets (fanout) **returns** set of processes **do**targets $:= \emptyset$;**while** $|\text{targets}| < \text{fanout}$ **do**candidate $:= \text{random}(II)$;**if** candidate $\notin \text{targets} \wedge \text{candidate} \neq \text{self}$ **then**targets $:= \text{targets} \cup \{\text{candidate}\}$;**return** targets;**procedure** gossip (msg) **do****forall** $t \in \text{pick-targets}(\text{fanout})$ **do**trigger $\langle \textit{up2pSend}, t, \text{msg} \rangle$;**upon event** $\langle \textit{pbBroadcast}, m \rangle$ **do**gossip ($[\text{GOSSIP}, s_m, m, \text{maxrounds}-1]$);**upon event** $\langle \textit{up2pDeliver}, p_i, [\text{GOSSIP}, s_m, m, r] \rangle$ **do****if** $m \notin \text{delivered}$ **then**delivered $:= \text{delivered} \cup \{m\}$ trigger $\langle \textit{pbDeliver}, s_m, m \rangle$;**if** $r > 0$ **then** gossip ($[\text{GOSSIP}, s_m, m, r-1]$);

Implements:

ProbabilisticBroadcast (pb).

Uses:

UnreliablePointToPointLinks (up2p), UnreliableBroadcast (unb).

upon event $\langle \text{Init} \rangle$ **do**

$\forall p_i \in \Pi$ delivered $[p_i] := 0$; lsn := 0; pending := \emptyset ; stored := \emptyset ;

procedure deliver-pending (s) **do**

$\forall_x : [\text{DATA}, s, x, sn_x]$ in pending such that $sn_x = \text{delivered}[s] + 1$ **do**
 delivered $[s] := \text{delivered}[s] + 1$; pending := pending $\setminus \{ [\text{DATA}, s, x, sn_x] \}$;
 trigger $\langle pb\text{Deliver}, s, x \rangle$;

procedure gossip (msg) **do**

forall $t \in \text{pick-targets}(\text{fanout})$ **do**
 trigger $\langle up2p\text{Send}, t, \text{msg} \rangle$;

upon event $\langle pb\text{Broadcast}, m \rangle$ **do**

lsn := lsn + 1;
 trigger $\langle un\text{Broadcast}, [\text{DATA}, \text{self}, m, \text{lsn}] \rangle$;

upon event $\langle un\text{Deliver}, p_i, [\text{DATA}, s_m, m, sn_m] \rangle$ **do**

if random > store-threshold **then** stored := stored $\cup \{ [\text{DATA}, s_m, m, sn_m] \}$;

if $sn_m = \text{delivered}[s_m] + 1$ **then**
 delivered $[s_m] := \text{delivered}[s_m] + 1$;
 trigger $\langle pb\text{Deliver}, s_m, m \rangle$;

else

pending := pending $\cup \{ [\text{DATA}, s_m, m, sn_m] \}$;
 $\forall \text{seqnb} \in [s_m - 1, \text{delivered}[s_m] + 1]$ **do**
 gossip $\langle [\text{REQUEST}, \text{self}, s_m, \text{seqnb}, \text{maxrounds} - 1] \rangle$;

upon event $\langle up2p\text{Deliver}, p_j, [\text{REQUEST}, p_i, s_m, sn_m, r] \rangle$ **do**

if $[\text{DATA}, s_m, m, sn_m] \in \text{stored}$ **then**
 trigger $\langle up2p\text{Send}, p_i, [\text{DATA}, s_m, m, sn_m] \rangle$;
else if $r > 0$ **then** gossip $\langle [\text{REQUEST}, p_i, s_m, sn_m, r - 1] \rangle$;

upon event $\langle up2p\text{Deliver}, p_j, [\text{DATA}, s_m, m, sn_m] \rangle$ **do**

if $sn_m = \text{delivered}[s_m] + 1$ **then**
 delivered $[s_m] := \text{delivered}[s_m] + 1$;
 trigger $\langle pb\text{Deliver}, s_m, m \rangle$;
 deliver-pending (s_m);

else

pending := pending $\cup \{ [\text{DATA}, s_m, m, sn_m] \}$;

Implements:

Probabilistic Partial Membership (ppm).

Uses:

unreliablePointToPointLinks (up2p).

upon event $\langle \text{Init} \rangle$ **do**

view := set of known group members;

subs := \emptyset ; unsubs := \emptyset ;

every T units of time **do**

for 1 to fanout **do**

 target := random (view);

trigger $\langle \text{upp2pSend}, \text{target}, [\text{GOSSIP}, \text{subs}, \text{unsubs}] \rangle$;

upon $\langle \text{ppmJoin} \rangle$ **do**

 subs := subs \cup { self };

upon $\langle \text{ppmLeave} \rangle$ **do**

 unsubs := unsubs \cup { self };

upon event $\langle \text{up2pDeliver}, p_i, [\text{GOSSIP}, s, u] \rangle$ **do**

 view := view \setminus u;

 view := view \cup s \setminus { self };

 unsubs := unsubs \cup u;

 subs := subs \cup s \setminus { self };

 //trim variables

while | view | > viewsz **do**

 target := random (view);

 view := view \setminus { target };

 subs := subs \cup { target };

while | unsubs | > unsubsz **do** unsubs := unsubs \setminus { random(unsubs) };

while | subs | > subssz **do** subs := subs \setminus { random(subs) };

Module:

Name: (regular) Consensus (c).

Events:

Request: $\langle cPropose, v \rangle$: Used to propose a value for consensus.

Indication: $\langle cDecide, v \rangle$: Used to indicate the decided value for consensus.

Properties:

C1: *Termination*: Every correct process eventually decides some value.

C2: *Validity*: If a process decides v , then v was proposed by some process.

C3: *Integrity*: No process decides twice.

C4: *Agreement*: No two correct processes decide differently.

Module 5.1 Interface and properties of consensus.

Implements:

Consensus (c);

Uses:

BestEffortBroadcast (beb);

PerfectFailureDetector (\mathcal{P});**upon event** $\langle \text{Init} \rangle$ **do**correct := correct-last-round := \perp ;proposal-set := correct-this-round := \emptyset ;decided := \perp ;

round := 1;

upon event $\langle \text{crash}, p_i \rangle$ **do**correct := correct $\setminus \{p_i\}$;**upon event** $\langle cPropose, v \rangle$ **do**proposal-set := $\{v\}$;**trigger** $\langle \text{bebBroadcast}, [\text{MYSET}, \text{round}, \text{proposal-set}] \rangle$;**upon event** $\langle \text{bebDeliver}, p_i, [\text{MYSET}, \text{round}, \text{set}] \rangle$ **do**correct-this-round := correct-this-round $\cup \{p_i\}$;proposal-set := proposal-set $\cup \text{set}$;**upon** correct \subset correct-this-round **do**

round := round + 1;

if (correct-this-round = correct-last-round) \wedge (decided = \perp) **then** **trigger** $\langle cDecide, \min(\text{proposal-set}) \rangle$; **trigger** $\langle \text{bebBroadcast}, [\text{DECIDED}, \text{round}, \min(\text{proposal-set})] \rangle$;**else**

correct-last-round := correct-this-round;

 correct-this-round := \emptyset ; **trigger** $\langle \text{bebBroadcast}, [\text{MYSET}, \text{round}, \text{proposal-set}] \rangle$;**upon event** $\langle \text{bebDeliver}, p_i, [\text{DECIDED}, \text{round}, v] \rangle \wedge$ (decided = \perp) **do**

decided := v;

trigger $\langle cDecide, v \rangle$; **trigger** $\langle \text{bebBroadcast}, [\text{DECIDED}, \text{round} + 1, \min(\text{proposal-set})] \rangle$;

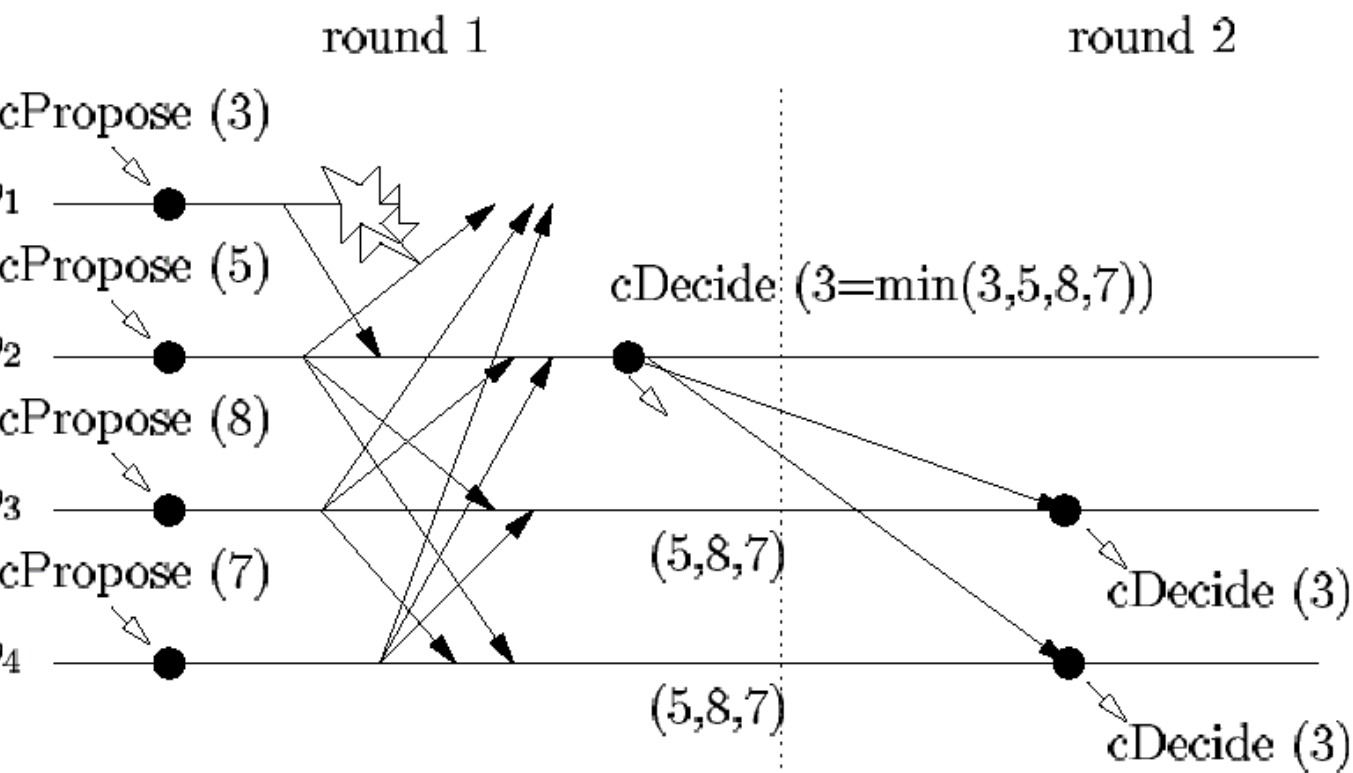


figure 5.1. Sample execution of the flooding consensus algorithm

Implements:Consensus (c);**Uses:**BestEffortBroadcast (beb);PerfectFailureDetector (\mathcal{P});**upon event** $\langle \text{Init} \rangle$ **do**suspected $:= \emptyset$;round $:= 1$;proposal $:= nil$;**for** $i = 1$ **to** N **do** pset[i] $:= p_i$;**for** $i = 1$ **to** N **do** delivered[round] $:= false$;**for** $i = 1$ **to** N **do** broadcast[round] $:= false$;**upon event** $\langle \text{crash}, p_i \rangle$ **do**suspected $:= suspected \cup \{p_i\}$;**upon event** $\langle cPropose, v \rangle$ **do**proposal $:= v$;**upon** (pset[round] = self) \wedge (proposal $\neq nil$) \wedge (broadcast[round] = false) **do****trigger** $\langle cDecide, proposal \rangle$;**trigger** broadcast[round] $:= true$;**trigger** $\langle bebBroadcast, proposal \rangle$;**upon** (pset[round] \in suspected) \vee (delivered[round] = true) **do**round $:= round + 1$;**upon event** $\langle bebDeliver, pset[round], value \rangle$ **do****if** self.id $>$ round **then**proposal $:= value$;delivered[round] $:= true$;

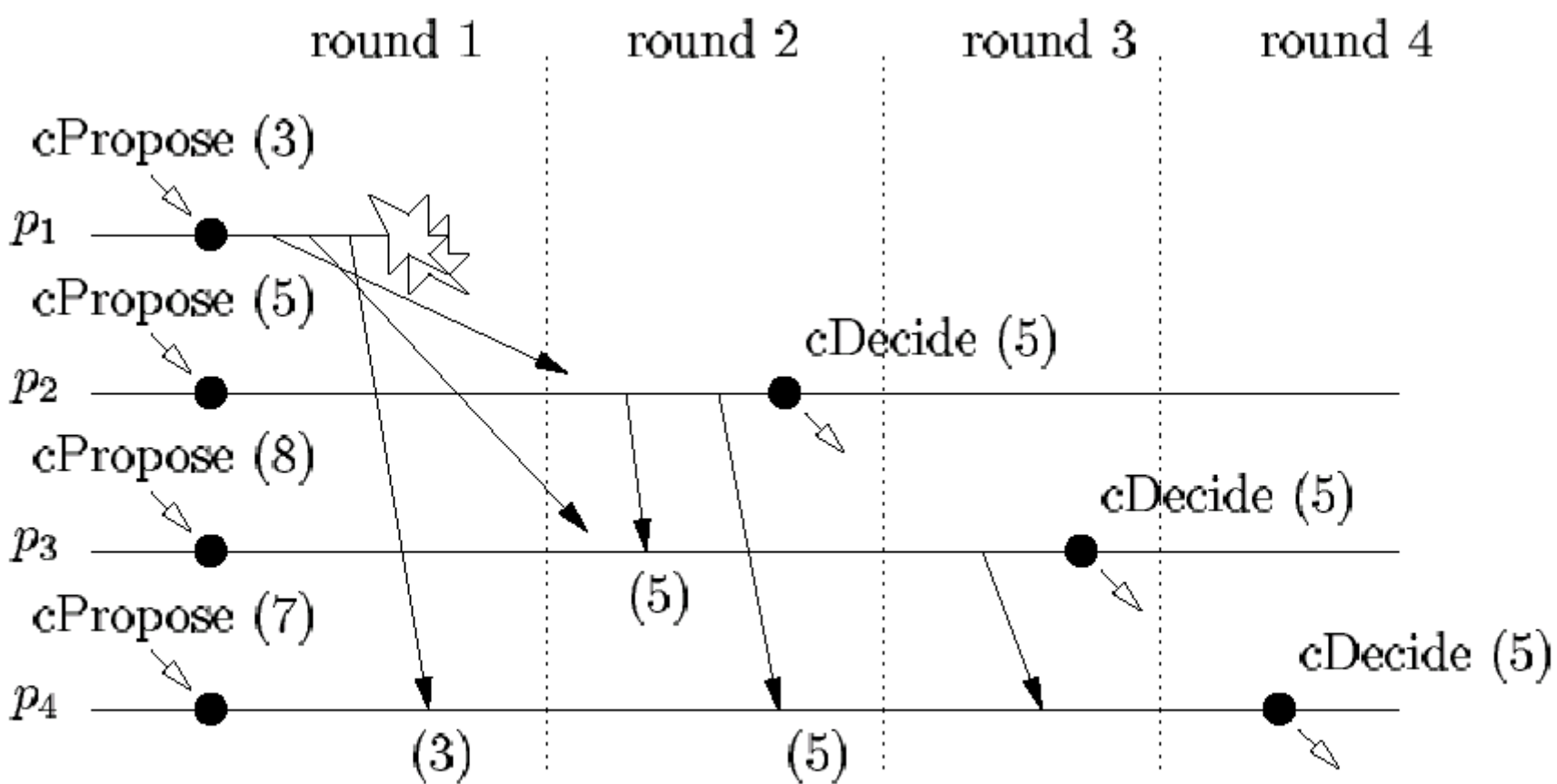


Figure 5.2. Sample execution of hierarchical consensus.

Module:

Name: UniformConsensus (uc).

Events:

$\langle \text{ucPropose}, v \rangle, \langle \text{ucDecide}, v \rangle$: with the same meaning and interface of the consensus interface.

Properties:

C1-C3: from consensus.

C4': *Uniform Agreement*: no two processes decide differently..

Module 5.2 Interface and properties of uniform consensus.

Implements:

UniformConsensus (c);

Uses:

BestEffortBroadcast (beb).

PerfectFailureDetector (\mathcal{P});**upon event** $\langle \text{Init} \rangle$ **do**correct := Π ;

round := 1;

for $i = 1$ **to** N **do** set[i] := delivered[i] := \emptyset ;proposal-set := \emptyset ;

decided := false;

upon event $\langle \text{crash}, p_i \rangle$ **do**correct := correct $\setminus \{p_i\}$;**upon event** $\langle \text{ucPropose}, v \rangle$ **do**proposal-set := $\{v\}$;**trigger** $\langle \text{bebBroadcast}, [\text{MYSET}, \text{round}, \text{proposal-set}] \rangle$;**upon event** $\langle \text{bebDeliver}, p_i, [\text{MYSET}, \text{round}, \text{newSet}] \rangle \wedge (p_i \in \text{correct})$ **do**set[round] := set[round] \cup newSet;delivered[round] := delivered[round] $\cup \{p_i\}$;**upon** (correct \subseteq delivered[round]) \wedge (decided = false) **do****if** round = N **then**

decided := true;

trigger $\langle \text{ucDecide}, \min(\text{proposal-set} \cup \text{set}[\text{round}]) \rangle$;**else**proposal-set := proposal-set \cup set[round];

round := round + 1;

trigger $\langle \text{bebBroadcast}, [\text{MYSET}, \text{round}, \text{proposal-set}] \rangle$;

Implements:

UniformConsensus (uc);

Uses:

PerfectPointToPointLinks (pp2p);

ReliableBroadcast (rb).

BestEffortBroadcast (beb).

PerfectFailureDetector (\mathcal{P});

upon event $\langle \text{Init} \rangle$ do

proposal := decided := \perp ;

round := 1;

suspected := ack-set := \emptyset ;

for $i = 1$ to N do pset[i] := p_i ;

upon event $\langle \text{crash}, p_i \rangle$ do

suspected := suspected $\cup \{p_i\}$;

upon event $\langle \text{ucPropose}, v \rangle$ do

proposal := v ;

upon $(\text{pset}[\text{round}] = \text{self}) \wedge (\text{proposal} \neq \perp) \wedge (\text{decided} = \perp)$ do

trigger $\langle \text{bebBroadcast}, [\text{PROPOSE}, \text{round}, \text{proposal}] \rangle$;

upon event $\langle \text{bebDeliver}, p_i, [\text{PROPOSE}, \text{round}, v] \rangle$ do

proposal := v ;

trigger $\langle \text{pp2pSend}, p_i, [\text{ACK}, \text{round}] \rangle$;

round := round + 1;

upon event $(\text{pset}[\text{round}] \in \text{suspected})$ do

round := round + 1;

upon event $\langle \text{pp2pDeliver}, p_i, [\text{ACK}, \text{round}] \rangle$ do

ack-set := ack-set $\cup \{p_i\}$;

upon event $(\text{ack-set} \cup \text{suspected} = \Pi)$ do

trigger $\langle \text{rbBroadcast}, [\text{DECIDED}, \text{proposal}] \rangle$;

upon event $\langle \text{rbDeliver}, p_i, [\text{DECIDED}, v] \rangle \wedge (\text{decided} = \perp)$ do

decided := v ;

trigger $\langle \text{ucDecide}, v \rangle$;

Module:

Name: Randomized Consensus (rc).

Events:

Request: $\langle rcPropose, v \rangle$: Used to propose a value for consensus.

Indication: $\langle rcDecide, v \rangle$: Used to indicate the decided value for consensus.

Properties:

RC1: *Termination:* With probability 1, every correct process decides some value.

RC2: *Validity:* If a process decides v , then v was proposed by some process.

RC3: *Integrity:* No process decides twice.

RC4: *Agreement:* No two correct processes decide differently.

Module 5.6 Interface and properties of randomized consensus.

Implements:

Randomized Consensus (rc);

Uses:

ReliableBroadcast (rb).

BestEffortBroadcast (beb).

upon event $\langle \text{Init} \rangle$ do

decided := \perp ; estimate := \perp ; round := 0;

for i = 1 to N **do** val[i] := \perp ;

val := \emptyset ;

upon event $\langle \text{rcPropose}, v \rangle$ do

trigger $\langle \text{bebBroadcast}, [\text{INVALUE}, v] \rangle$;

estimate := v; round := round + 1;

val := val $\cup \{v\}$;

trigger $\langle \text{bebBroadcast}, [\text{PHASE1}, \text{round}, v] \rangle$;

upon event $\langle \text{bebDeliver}, p_i, [\text{INIVAL}, v] \rangle$ do

val := val $\cup \{v\}$;

upon event $\langle \text{bebDeliver}, p_i, [\text{PHASE1}, r, v] \rangle$ do

phase1[r] := phase1[r] $\oplus v$;

upon (decided = $\perp \wedge |\text{phase1}[\text{round}]| > N/2$) do

if $\exists v : \forall x \in \text{phase1}[\text{round}] : x = v$ **then** estimate := v;

else estimate := \perp ;

trigger $\langle \text{bebBroadcast}, [\text{PHASE2}, \text{round}, \text{estimate}] \rangle$;

upon event $\langle \text{bebDeliver}, p_i, [\text{PHASE2}, r, v] \rangle$ do

phase2[r] := phase2[r] $\oplus v$;

upon (decided = $\perp \wedge |\text{phase2}[\text{round}]| > N/2$) do

if $\exists v \neq \perp : \forall x \in \text{phase1}[\text{round}] : x = v$ **then**

decided := v;

trigger $\langle \text{rbBroadcast}, [\text{DECIDED}, \text{round}, \text{decided}] \rangle$;

else

if $\exists v \in \text{phase2}[\text{round}] : v \neq \perp$ **then** estimate := v;

else estimate := random(val);

round := round + 1; // start one more round

trigger $\langle \text{rbBroadcast}, [\text{PHASE1}, \text{round}, \text{estimate}] \rangle$;

upon event $\langle \text{rbDeliver}, p_i, [\text{PHASE2}, r, v] \rangle$ do

decided := v;

trigger $\langle \text{rcDecided}, \text{decided} \rangle$;

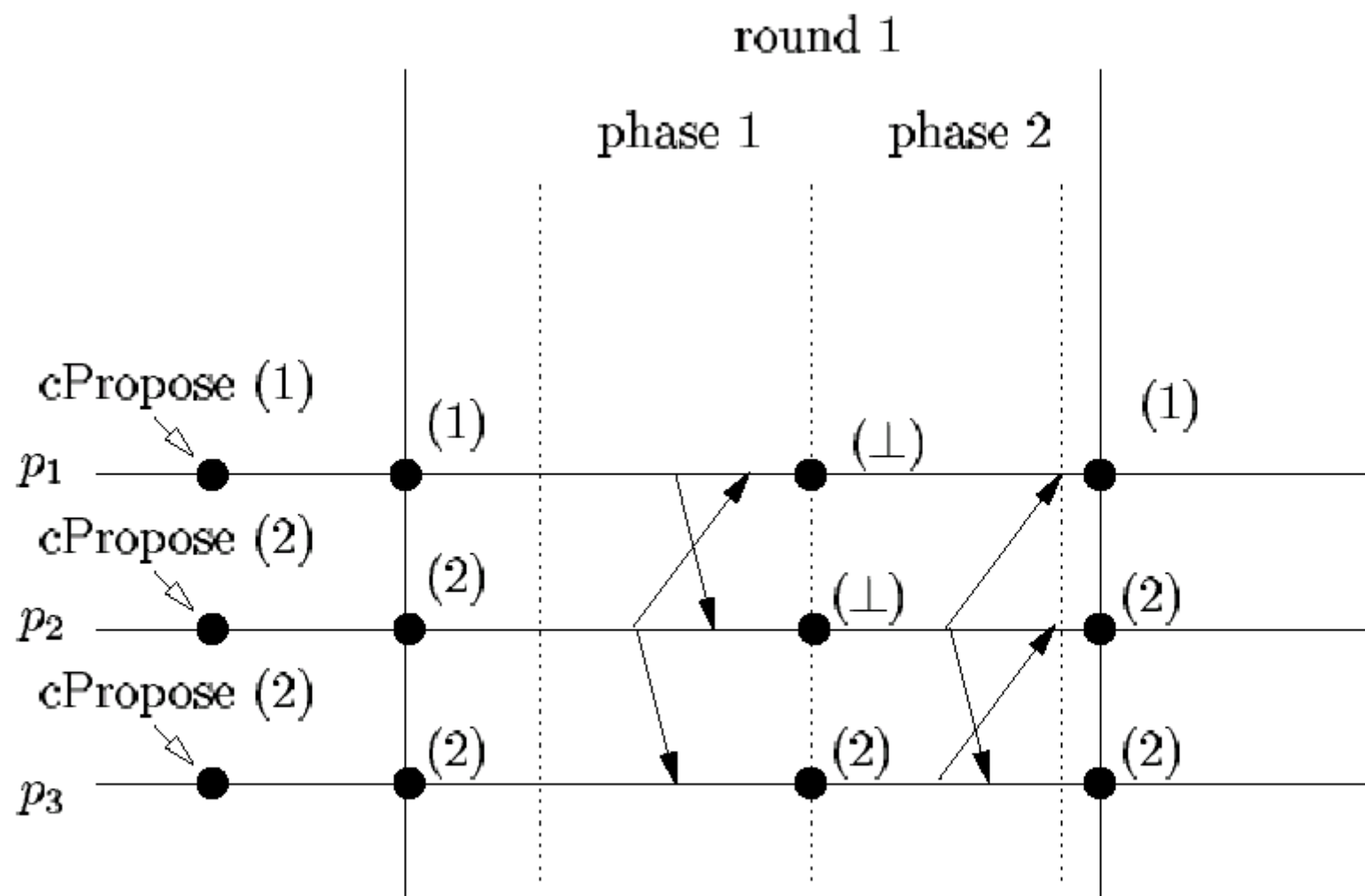


Figure 5.3. Role of randomization.