

Fault-tolerance in Message-Passing Distributed Systems

Annette Bieniusa

R
P **TU** Rheinland-Pfälzische
Technische Universität
Kaiserslautern
Landau

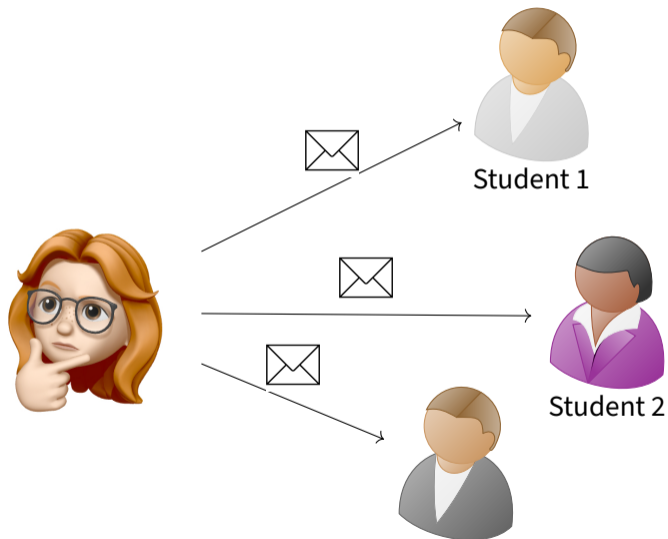


In the morning before I go to university

I communicate

- A. with just myself
- B. with my plants and/or pets
- C. by talking to at least one other human being
- D. by sending messages to some selected other human being
- E. by broadcasting information to “the world” using social media

Broadcasting Information



The Need for Distributed Algorithms

- Distributed algorithms are at the core of any distributed systems
- Implemented as middleware between network and application
- Services beyond network protocols (e.g. TCP, UDP)
 - Group communication
 - Shared memory abstractions
 - State machine replication

Overview

- Formal models for specifying and analyzing distributed algorithms
- Composability of distributed algorithms
- The Broadcast Problem
 - Best-effort broadcast
 - Reliable broadcast
 - FIFO broadcast
 - Causal broadcast
 - Total-Order broadcast

Goals of this Lecture

In this lecture, you will learn

- to formally specify safety and liveness properties of several broadcast problem ([1])
- to define fault-tolerant algorithms for Best-effort, Reliable, FIFO and Causal Broadcast in an asynchronous system with reliable channels
- to prove the correctness of some algorithms
- to use space-time diagrams to visualize executions
- to implement these algorithms in Elixir

Relation to this Summer School

- Practical: Implementing a chat application in Elixir with different broadcast variants
- Testing Distributed System Implementations
- Verification with TLA+ [9]
- Foundations and mind-set for all the other lectures

The Broadcast Problem

Informally: A process needs to transmit a message to other processes.

$\text{broadcast}(m) \approx \mathbf{forall} \ j \in \{1, \dots, n\} : \text{send } m \text{ to } p_j$

System model

- Asynchronous system
- Static set of processes $\Pi = \{p_1, \dots, p_n\}$
 - crash-stop fault model
- Sending and receiving messages through reliable channels (*perfect point-to-point links*)
 - no message loss / creation / modification / duplication
 - bidirectional
 - infinite capacity
- Messages are uniquely identifiable
 - e.g. tagged with `<sender, seq_number>`

Only a subset $\Pi' \subseteq \Pi$ receives messages in arbitrary order at distinct, independent time instants.

What is the simplest solution that you can think of?

What is the simplest solution that you can think of?

Just go ahead and send the message to everyone, one at a time.

Specifying Distributed Algorithms

Deterministic I/O automaton with send/receive operations

- Event-driven programming model
- Events triggered by messages, timers, conditions, ...

```
Upon Event(arg1, arg2, ...) do:  
  // local computation  
  trigger Event(arg1', arg2', ...)
```

- Correctness properties
 - Safety: Nothing bad ever happens
 - Liveness: Something good eventually happens

The Anatomy of a Broadcast Algorithm

For the broadcast algorithms:

Upon `Init` **do:** ...

Upon `Broadcast(m)` **do:** ...

Upon `Receive(p_k , m)` **do:** ...

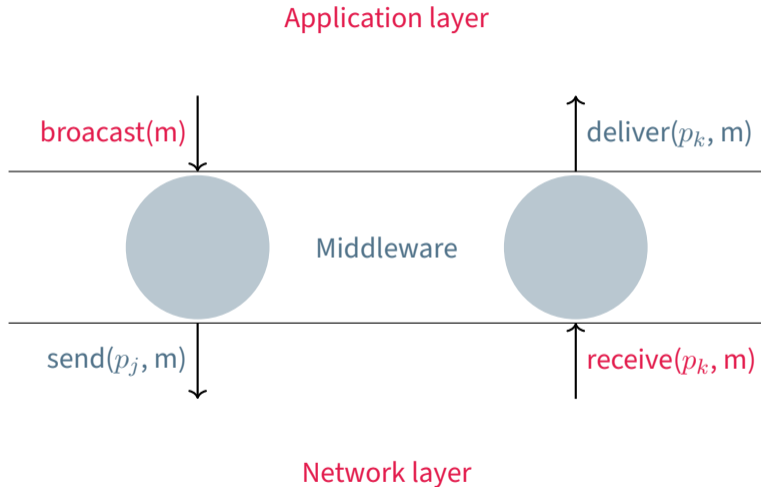
- You can trigger an event on another component, e.g.:

trigger `Send(p_j , m)`

trigger `Deliver(p_k , m)`

- There is a special event called `Init` for initializing the local state.
- p_j denotes the target process when sending a message
- p_k denotes the process where the message originated from

At Process p_i



Best-effort Broadcast (BEB): Specification

- *BEB-Validity*: If a correct process p_j beb-delivers a message m , then m has previously been beb-broadcast to p_j by some process p_i .
 - No creation, no alteration of messages
- *BEB-Integrity*: A process beb-delivers a message m at most once.
 - No duplication of messages
- *BEB-Termination*: For any two **correct** processes p_i and p_j , every message that has been beb-broadcast by p_i is eventually beb-delivered by p_j .

Best-effort Broadcast: Algorithm

Idea:

- Just go ahead and send the message to every other process.
- When you get one of these messages, you deliver it to the upper layer.
- Intuition: No guarantees if sender crashes

State: --

Upon Init do: --

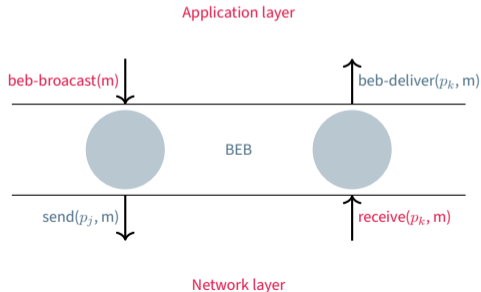
Upon beb-broadcast(m) **do:**

forall $p_j \in \Pi$:

trigger send(p_j, m)

Upon receive(p_k, m) **do:**

trigger beb-deliver(p_k, m)



Best-effort Broadcast: Correctness

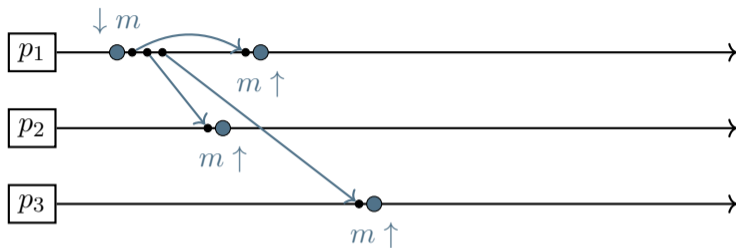
Why does it work?

- BEB-Validity holds because Perfect-Link model guarantees no creation and there is no other way for messages to appear, only through beb-broadcast
- BEB-Integrity holds because Perfect-Link model guarantees no duplication
- BEB-Termination holds because Perfect-Link model guarantees reliable delivery

Perfect-Link Model

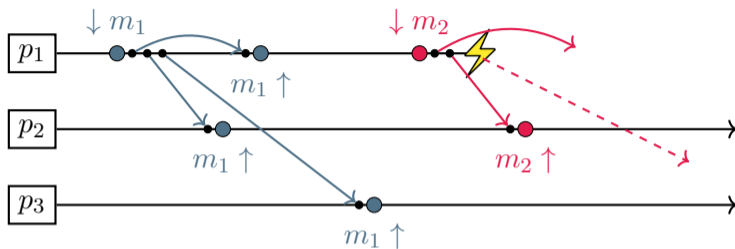
- **Reliable Delivery:** Considering two correct processes i and j ; if i sends a message m to j , then j eventually delivers m .
- **No Duplication:** No message is delivered by a process more than once.
- **No Creation:** If a correct process j delivers a message m , then m was sent to j by some process i .

Visualizing Executions with Space-Time Diagrams



- $\downarrow m$ = broadcast message m
- $\uparrow m$ = deliver message m

Best-effort Broadcast: Sender crashes



Limitations of Best-effort Broadcast

What happens if a process fails while sending a message?

- If the sender crashes before being able to send the message to all processes, some process will not deliver the message.
- Even in the absence of communication failures!

Limitations of Best-effort Broadcast

What happens if a process fails while sending a message?

- If the sender crashes before being able to send the message to all processes, some process will not deliver the message.
- Even in the absence of communication failures!

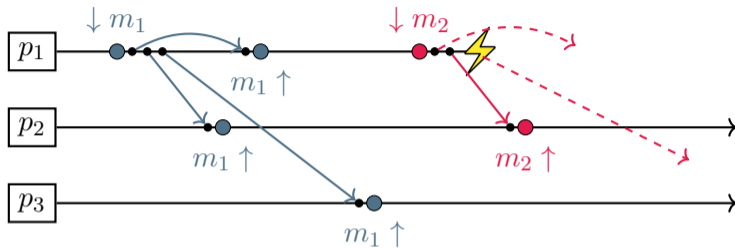
Let's try for a reliable version of broadcast!

- Guarantees that all or none of the correct nodes gets the message
- Even if sender crashes!

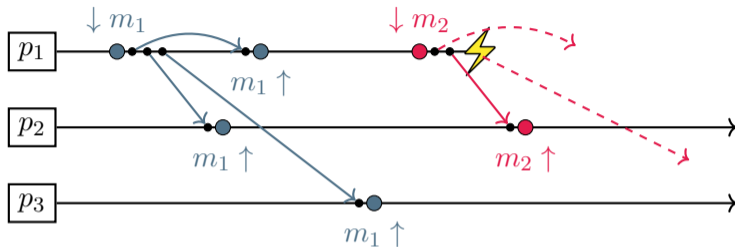
Reliable Broadcast (RB): Specification

- *RB-Validity*: If a correct process p_i rb-delivers a message m , then m has been previously rb-broadcast.
- *RB-Integrity*: A process rb-delivers a message m at most once.
- *RB-Termination-1*: If a correct process p_i rb-broadcasts message m , then p_i rb-delivers the message m .
- *RB-Termination-2*: If a correct process p_i rb-delivers a message m , then each correct process rb-delivers m .

Reliable Broadcast: Scenario 1



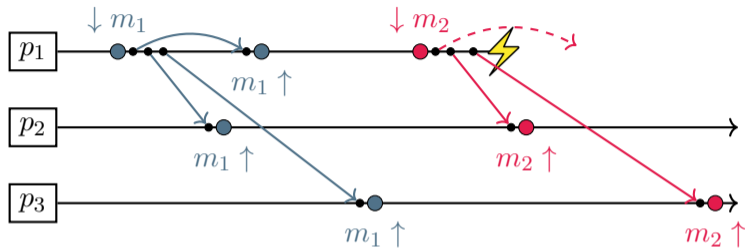
Reliable Broadcast: Scenario 1



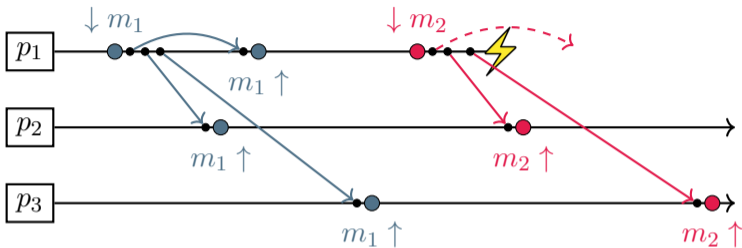
Not possible under Reliable Broadcast: RB-Termination-2 is violated!

If correct process p_2 delivers m , then correct process p_3 must also rb-deliver m .

Reliable Broadcast: Scenario 2

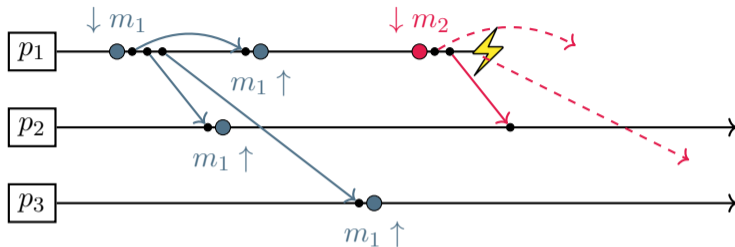


Reliable Broadcast: Scenario 2

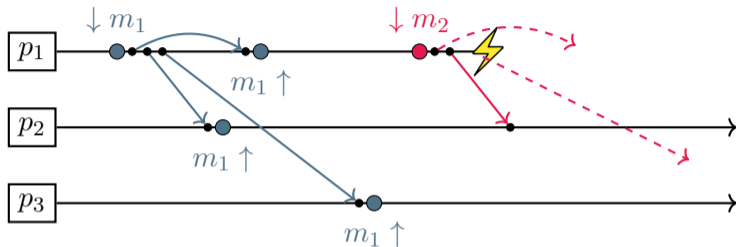


The fact that process p_1 does not deliver m_2 is not a problem, because only correct processes are required to deliver their own messages.

Reliable Broadcast: Scenario 3

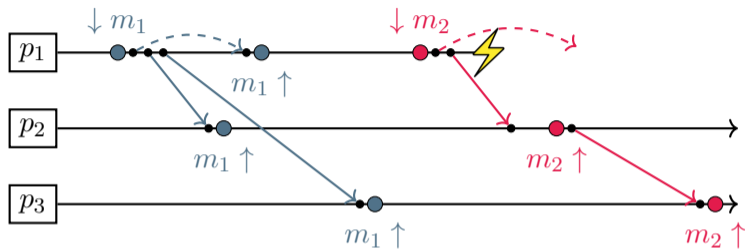


Reliable Broadcast: Scenario 3



The fact that no process delivers m_2 is not a problem (even though p_2 received it), because process p_1 has crashed and no process delivers m_2 .

Reliable Broadcast: Idea!



Reliable Broadcast: Algorithm

State:

delivered

Upon Init do:

delivered $\leftarrow \emptyset$

Upon rb-broadcast(m) do

$m_{id} \leftarrow \text{generateUniqueID}(m)$

trigger beb-broadcast($[m_{id}, m]$)

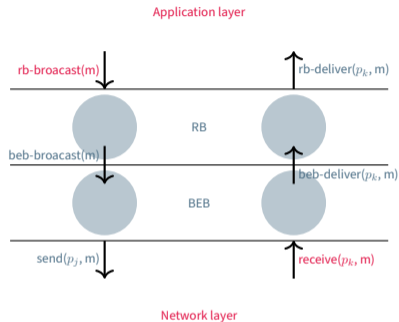
Upon beb-deliver($p_k, [m_{id}, m]$) do

if ($m_{id} \notin \text{delivered}$) **then**

delivered $\leftarrow \text{delivered} \cup \{m_{id}\}$

trigger rb-deliver(p_k, m)

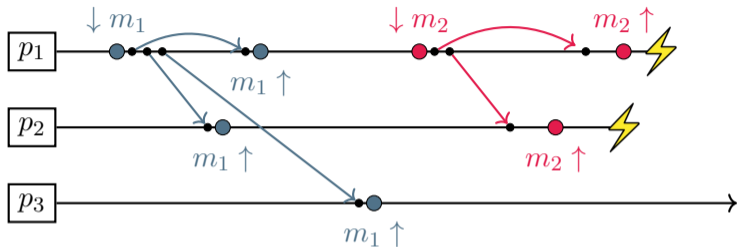
trigger beb-broadcast($[m_{id}, m]$)



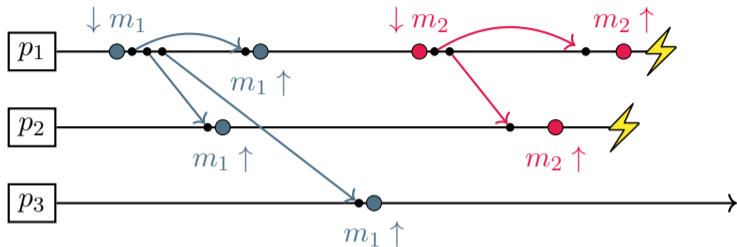
Reliable Broadcast: Correctness

- *RB-Validity*: If a correct process p_i rb-delivers a message m , then m has previously been rb-broadcast.
 - By BEB-Validity.
- *RB-Integrity*: A process rb-delivers a message m at most once.
 - By BEB-Integrity and handling the set of delivered messages.
- *RB-Termination-1*: If a correct process p_i broadcasts message m , then p_i eventually rb-delivers m .
 - By BEB-Termination and handling of the set of delivered messages.
- *RB-Termination-2*: If a correct process p_i rb-delivers a message m , then each correct process rb-delivers m .
 - After rb-delivering m , a correct process forwards m to all processes. By BEB-Termination and p_i being correct, all correct processes will eventually beb-deliver m and hence rb-deliver it.

Reliable Broadcast: Scenario 4



Reliable Broadcast: Scenario 4



The fact that m_2 has been delivered by faulty p_1 and p_2 does not imply that m_2 has to be delivered by p_3 as well under rb-broadcast. Yet, this situation is not desirable, because two processes deliver something and another one does not.

⇒ Interaction with external world!

Uniform Reliable Broadcast (URB): Specification

- *URB-Validity*: If a correct process p_i urb-delivers a message m , then m was urb-broadcast to p_i by some process p_j .
- *URB-Integrity*: A process p_i urb-delivers a message m at most once.
- *URB-Termination-1*: If a correct process p_i urb-broadcasts a message m , then p_i eventually urb-delivers m .
- *URB-Termination-2*: If a process p_i urb-delivers a message m , then each correct process p_j eventually urb-delivers m .

An Impossibility Result

- n : total number of processes
- t : upper bound on the number of processes that can fail
- Fail-silent system model: crash-stop + perfect point-to-point links

Theorem

There is no algorithm implementing URB under the fail-silent system model if a majority of processes can fail, i.e. if $t \geq n/2$.

Proof sketch

By contradiction.

- Assume there exists algorithm A that implements URB under the fail-silent model for $t \geq n/2$.
- Partition $\Pi = P_1 \cup P_2$ such that
 - $P_1 \cap P_2 = \emptyset$
 - $|P_1| = \lceil n/2 \rceil$ and $|P_2| = \lfloor n/2 \rfloor$ ($|P_1| \geq |P_2|$)
- Consider two executions E_1 and E_2
- Execution E_1 :
 - All $p_i \in P_2$ crash initially, all processes in P_1 are correct.
 - If $p_x \in P_1$ issues $\text{urb-broadcast}(m)$ using algorithm A , then every process in P_1 must eventually $\text{urb-deliver } m$ (assuming A correctly implements URB)
 - Under fail-silent model, the decision to deliver must be independent of the status of the processes in P_2

Proof sketch (2)

- Execution E_2 :
 - All $p_i \in P_2$ are correct, and initially all processes in P_1 are well-behaving
 - If $p_x \in P_1$ issues $\text{urb-broadcast}(m)$ using algorithm A , then every process in P_1 must eventually $\text{urb-deliver } m$ (assuming A correctly implements URB)
 - The decision to $\text{urb-deliver } m$ is made by the same algorithm A as before, i.e. it is independent of the status of the processes in P_2
 - Assume that after the delivery, all processes in P_1 crash
 - If m has not reached any process in P_2 , yet, it cannot be urb-delivered by processes in P_2 anymore, because the perfect-link model requires sender and receiver to be correct for reliable delivery.
- Contradiction to the assumption that the algorithm implements URB

Uniform Reliable Broadcast for $t < n/2$: Algorithm

State:

```
delivered //set of message ids that have been delivered
pending // set of messages to be delivered
ack // map  $m_{id}$  to received acknowledgments
```

Upon Init do:

```
delivered  $\leftarrow \emptyset$ 
pending  $\leftarrow \emptyset$ 
 $\forall m_{id}$ : ack[ $m_{id}$ ]  $\leftarrow \emptyset$ 
```

Upon urb-broadcast(m) do

```
 $m_{id} \leftarrow \text{generateUniqueID}(m)$ 
pending  $\leftarrow$  pending  $\cup$  {[self,  $m_{id}$ , m]}
trigger beb-broadcast([self,  $m_{id}$ , m])
```

Uniform Reliable Broadcast for $t < n/2$: Algorithm (2)

```
Upon beb-deliver( $p_k$ , [ $p_j$ ,  $m_{id}$ ,  $m$ ]) do  
  ack[ $m_{id}$ ]  $\leftarrow$  ack[ $m_{id}$ ]  $\cup$  { $k$ }  
  if ( ( $p_j$ ,  $m_{id}$ ,  $m$ )  $\notin$  pending ) then  
    pending  $\leftarrow$  pending  $\cup$  ( $p_j$ ,  $m_{id}$ ,  $m$ )  
    trigger beb-broadcast([ $p_j$ ,  $m_{id}$ ,  $m$ ])  
  
Upon exists ( $p_j$ ,  $m_{id}$ ,  $m$ )  $\in$  pending  
  with |ack[ $m_{id}$ ]|  $>$   $n/2$  and  $m_{id} \notin$  delivered  
  delivered  $\leftarrow$  delivered  $\cup$   $m_{id}$   
  trigger urb-deliver( $p_j$ ,  $m$ )
```

Uniform Reliable Broadcast: Correctness (Sketch)

- Assume majority of correct processes ($t < n/2$)
- If a process urb-delivers a message, it has received acknowledgement from majority
- In this majority, at least one process p must be correct
- p ensures that all correct processes beb-deliver m by forwarding the message

Resilience

- Defined by maximum number of faulty processes that an algorithm can handle
- Algorithm for URB under fail-silent model has resilience $< n/2$

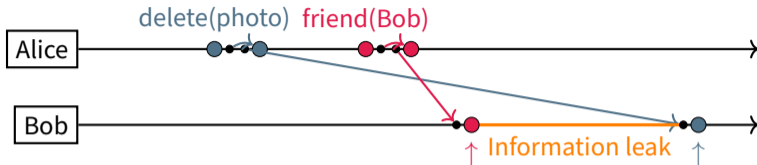
Problem: Message ordering

- Given the asynchronous nature of distributed systems, messages may be delivered in *any* order.
- Some services, such as replication, need messages to be delivered in a consistent manner, otherwise replicas may diverge.

Breakout: I know what you have seen

- Think about email threads between multiple persons exchanging information
- How can you determine when two answers are given concurrently?
- How can you reconstruct what information a person who answers in the thread has seen?

FIFO Order

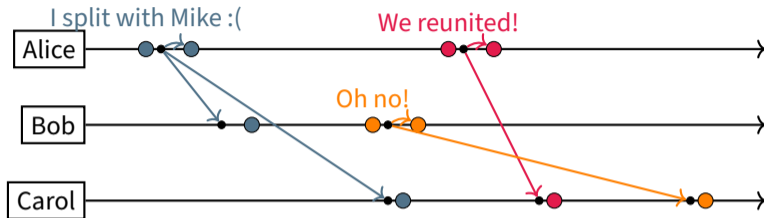


FIFO Property

If a process p broadcasts a message m before the same process broadcasts another message m' , then no correct process q delivers m' unless it has previously delivered m .

$$\text{broadcast}_p(m) \rightarrow \text{broadcast}_p(m') \Rightarrow \text{deliver}_q(m) \rightarrow \text{deliver}_q(m')$$

Causal Order



Causality Property

If the broadcast of a message m happens-before the broadcast of some message m' , then no correct process delivers m' unless it has previously delivered m .

$$\text{broadcast}_p(m) \rightarrow \text{broadcast}_q(m') \Rightarrow \text{deliver}_r(m) \rightarrow \text{deliver}_r(m')$$

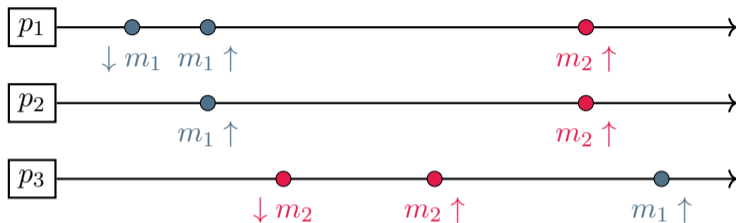
Total Order

Total Order Property

If correct processes p and q both deliver messages m, m' , then p delivers m before m' if and only if q delivers m before m' .

$$\text{deliver}_p(m) \rightarrow \text{deliver}_p(m') \Rightarrow \text{deliver}_q(m) \rightarrow \text{deliver}_q(m')$$

Message ordering: Quizzzzzz



Is this allowed under FIFO Order, Causal Order, Total Order?

(Reliable) FIFO Broadcast (FIFO): Specification

- All properties from reliable broadcast
- *FIFO delivery*: If a process fifo-broadcasts m and later m' , then no process fifo-delivers m' unless it already delivered m .

FIFO-Broadcast: Algorithm

State:

```
next    // array mapping process id to seq numer
seq     // sequence numbers for broadcast messages
pending // messages to be delivered
```

Upon Init do:

```
next  $\leftarrow$  [0, ..., 0]; seq  $\leftarrow$  0; pending  $\leftarrow$   $\emptyset$ 
```

Upon fifo-broadcast(m) do

```
 $m_{id} \leftarrow$  seq++ // generate next seq number
```

```
trigger rb-broadcast([ $m_{id}$ , m])
```

Upon rb-deliver(p_k , [m_{id} , m]) do

```
if  $m_{id} =$  next[ $p_k$ ] then
```

```
    trigger fifo-deliver( $p_k$ , m)
```

```
    next[ $p_k$ ]++
```

```
    while exists ( $p_k, n_{id}, n$ )  $\in$  pending with  $n_{id} =$  next[ $p_k$ ] do
```

```
        trigger fifo-deliver( $p_k$ , n)
```

```
        next[ $p_k$ ]++
```

```
        pending  $\leftarrow$  pending  $\setminus$  {( $p_k, n_{id}, n$ )}
```

```
else pending  $\leftarrow$  pending  $\cup$  {( $p_k, m_{id}, m$ )}
```

(Reliable) Causal Broadcast (RCO): Specification

- All properties from reliable broadcast
- *Causal delivery*: No process p_i delivers a message m' unless p_i has already delivered every message m such that $m \rightarrow m'$.

Idea

- Each messages carries past_m , an ordered list of messages that causally precede m
- When a process rb-delivers m ,
 - it co-delivers first all causally preceding messages in past_m
 - it co-delivers m
 - avoiding duplicates using `delivered`

Causal Broadcast (RCO): Algorithm 1 (No-waiting)

State:

```
delivered //set of msg ids that have been rco-delivered  
past      // ordered list
```

Upon Init do:

```
delivered <- ∅  
past      <- []
```

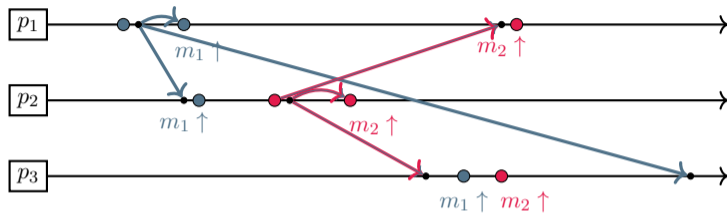
Upon rco-broadcast(m) do

```
 $m_{id}$  <- generateUniqueID(m)  
trigger rb-broadcast([ $m_{id}$  , past, m])  
past <- past ++ [(self,  $m_{id}$ , m)] // append at the end
```

Causal Broadcast (RCO): Algorithm 1 (No-waiting) - Continued

```
Upon rb-deliver( $p_k$ , [ $m_{id}$ ,  $past_m$ ,  $m$ ]) do
  if (  $m_{id} \notin delivered$  ) then
    for ( $p_j$ ,  $n_{id}$ ,  $n$ ) :  $past_m$  do // from old to recent
      if ( $n_{id} \notin delivered$  ) then
        trigger rco-deliver( $p_j$ ,  $n$ )
        delivered  $\leftarrow$  delivered  $\cup$  { $n_{id}$ }
        if ( $p_j$ ,  $n_{id}$ ,  $n$ )  $\notin past$  then
          past  $\leftarrow$  past ++ [( $p_j$ ,  $n_{id}$ ,  $n$ )]
      trigger rco-deliver( $p_k$ ,  $m$ )
    delivered  $\leftarrow$  delivered  $\cup$  { $m_{id}$ }
  if ( $p_k$ ,  $m_{id}$ ,  $m$ )  $\notin past$  then
    past  $\leftarrow$  past ++ [( $p_k$ ,  $m_{id}$ ,  $m$ )]
```

Causal Broadcast: Scenario 1



Causal Broadcast - Algorithm 1: Correctness

- Validity follows directly from rb-broadcast
- Integrity follows from rb-broadcast and the check before rco-delivering messages from past_m
- Termination follows directly from rb-broadcast and the fact that no waiting occurs
 - Every message is rco-delivered once rb-delivered
- Causal delivery
 - Each message m carries its causal past
 - Causal past is in order delivered before m
 - Proof by induction on trace prefix
 - Initial state
 - For every delivery

Remarks

- Message from causal past of m are delivered before message m (*causal delivery*)
- Message id's could be reused for rb-broadcast
- Size of messages grows linearly with every message that is broadcast since it includes the complete causal past

Idea: Garbage collect the causal past

- If we know when a process fails (i.e., under the fail-stop model), we can remove messages from the causal past.
 - When a process rb-delivers a message m , it rb-broadcasts an acknowledgement message to all other processes.
 - When an acknowledgement for message m has been rb-delivered by all correct processes, m is removed from *past*
 - N^2 additional ack messages for each application message
 - Typically, acknowledgements are grouped and processed in batch mode
- ⇒ Requires still unbounded messages sizes

Efficient representation of causal past: Vector clocks

- A vector clock [6] is a mapping from processes to natural numbers $p_i \mapsto \mathbb{N}$.
 - Example: $[p_1 \mapsto 3, p_2 \mapsto 4, p_3 \mapsto 1]$
 - If processes are numbered $1, \dots, n$, this mapping can be represented as a vector, e.g., $[3, 4, 1]$
 - Intuitively: $p_1 \mapsto 3$ means “observed 3 events (here: messages broadcast) from process p_1 ”

Vector time

Each process p_i stores current causal past as a vector clock VC .

- Initially, $VC[k] := 0$ for all k
- On each local event, process p_i increments its own entry as follows:
 $VC[i] := VC[i] + 1$
- On sending a message m , p_i attaches VC to m and increments VC for itself afterwards
- On receiving message m with vector time VC_m , the messages are processed in causal order, potentially waiting for missing updates, and increment the local vector accordingly

Relating vector times

Let u, v denote time vectors.

- $u \leq v$ iff $u[k] \leq v[k]$ for $k = 1, \dots, n$
- $u < v$ iff $u \leq v$ and $u \neq v$
- $u \parallel v$ iff $u \not\leq v$ and $v \not\leq u$

Causal Broadcast (RCO): Algorithm 2 [8]

State:

pending //set of messages that cannot be delivered yet
VC // vector clock

Upon Init do:

pending $\leftarrow \emptyset$
forall $p_i \in \Pi$ **do:** VC[p_i] $\leftarrow 0$

Upon rco-broadcast(m) do

trigger rco-deliver(self, m)
trigger rb-broadcast(VC, m)
VC[self] \leftarrow VC[self] + 1

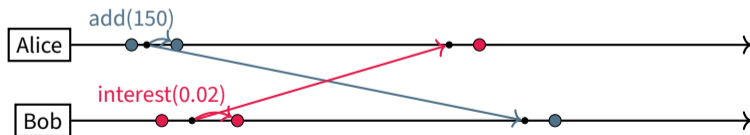
Upon rb-deliver(p_k , VC $_m$, m) do

if ($p_k \neq$ self) **then**
 pending \leftarrow pending \cup {(p_k , VC $_m$, m)}
 while exists (q, VC $_{m_q}$, m_q) \in pending **with** VC \geq VC $_{m_q}$ **do**
 pending \leftarrow pending \setminus {(q, VC $_{m_q}$, m_q)}
 trigger rco-Deliver(q, m_q)
 VC[q] \leftarrow VC[q] + 1

Limitations of Causal Broadcast

Example: Replicated database handling bank accounts

- Initially, account A holds 1000 Euro.
- User deposits 150 Euro, triggers broadcast of message
 $m_1 = \text{'add 150 Euro to A'}$
- Concurrently, bank initiates broadcast of message
 $m_2 = \text{'add 2% interest to A'}$
- Diverging state because processes can observe messages in different order



Total-order broadcast (aka Atomic Broadcast)

- All processes deliver their messages in the same order
- Typical use case: for Replicated State Machines (RSM) to implement

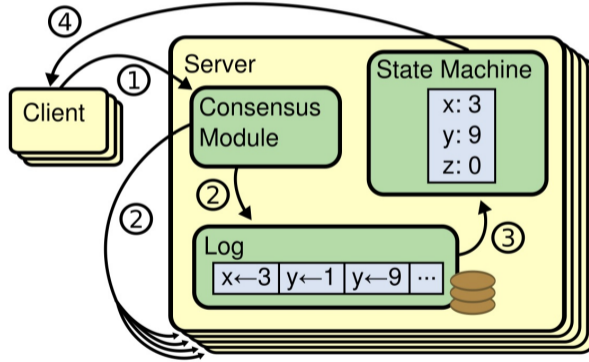


Figure: [7]

The FLP Theorem [2]

There is no deterministic protocol that solves consensus in an asynchronous system in which a single process may fail by crashing.

- 2001 Dijkstra prize for the most influential paper in distributed computing
- Proof Strategy
 - Assume that there **is** a (deterministic) protocol to solve the problem
 - Reason about the properties of **any** such protocol
 - Derive a contradiction \Rightarrow Done :)

Core Ideas

- Replicated log \Rightarrow State-machine replication
 - Each server stores a log containing a sequence of state-machine commands.
 - All servers execute the same commands in the same order.
 - Once one of the state machines finishes execution, the result is returned to the client.
- Consensus module ensures correct log replication
 - Receives commands from clients and adds them to the log
 - Communicates with consensus modules on other servers such that every log eventually contains same commands in same order
- *Failure model*: Nodes may crash, recover and rejoin, delayed/lost messages

Classification

- *Leader-less (symmetric)*
 - All servers are operating equally
 - Clients can contact any server
- *Leader-based (asymmetric)*
 - One server (called leader) is in charge
 - Other server follow the leader's decisions
 - Clients interact with the leader, i.e. all requests are forwarded to the leader
 - If leader crashes, a new leader needs to be (s)elected
 - Quorum for choosing leader in next epoch (i.e. until the leader is suspected to have crashed)
 - Then, overlapping quorum decides on proposed value \Rightarrow Only accepted if no node has knowledge about higher epoch number

Classic approaches I

- Paxos[4]
 - The original consensus algorithm for reaching agreement on a **single value**
 - Leader-based
 - Two-phase process: Promise and Commit
 - Clients have to wait 2 RTTs
 - Majority agreement: The system works as long as a majority of nodes are up
 - Monotonically increasing version numbers
 - Guarantees safety, but not liveness

Classic approaches II

- Multi-Paxos
 - Extends Paxos for a stream of agreement problems (i.e. total-order broadcast)
 - The promise (Phase 1) is not specific to the request and can be done before the request arrives and can be reused
 - Client only has to wait 1 RTT
- View-stamped replication (revisited)[5]
 - Variant of SMR + Multi-Paxos
 - Round-robin leader election
 - Dynamic membership

The Problem with Paxos

[...] I got tired of everyone saying how difficult it was to understand the Paxos algorithm.[...] The current version is 13 pages long, and contains no formula more complicated than $n1 > n2$. [3]

Still significant gaps between the description of the Paxos algorithm and the needs of a real-world system

- Disk failure and corruption
- Limited storage capacity
- Effective handling of read-only requests
- Dynamic membership and reconfiguration

In Search of an Understandable Consensus Algorithm: Raft[7]

- Yet another variant of SMR with Multi-Paxos
- Became very popular because of its understandable description

In a nutshell

- Strong leadership with all other nodes being passive
- Dynamic membership and log compaction

Summary

- Different variants of solution to the Broadcast Problem
 - **Best-effort broadcast:** Reliable only if sender is correct
 - **Reliable broadcast:** Reliable independent of whether sender is correct
 - **Uniform reliable broadcast:** Considers also behavior of failed nodes
 - **FIFO broadcast:** Reliable broadcast with FIFO delivery order
 - **Causal broadcast:** Reliable broadcast with causal delivery order
 - **Total-order broadcast:** Reliable and same order of delivery at all nodes
- Correctness proofs based on properties of underlying level + algorithmic properties

Further reading I

- [1] Christian Cachin, Rachid Guerraoui, and Luis Rodrigues. *Introduction to Reliable and Secure Distributed Programming (2. ed.)* Springer, 2011. ISBN: 978-3-642-15259-7. DOI: 10.1007/978-3-642-15260-3. URL: <https://doi.org/10.1007/978-3-642-15260-3>.
- [2] Michael J. Fischer, Nancy A. Lynch, and Mike Paterson. “Impossibility of Distributed Consensus with One Faulty Process”. In: *J. ACM* 32.2 (1985), pp. 374–382. DOI: 10.1145/3149.214121. URL: <http://doi.acm.org/10.1145/3149.214121>.
- [3] Leslie Lamport. “Paxos Made Simple”. In: *SIGACT News* 32.4 (Dec. 2001), pp. 51–58. ISSN: 0163-5700. DOI: 10.1145/568425.568433. URL: <http://research.microsoft.com/users/lamport/pubs/paxos-simple.pdf>.

Further reading II

- [4] Leslie Lamport. “The Part-Time Parliament”. In: *ACM Trans. Comput. Syst.* 16.2 (1998), pp. 133–169. DOI: 10.1145/279227.279229. URL: <http://doi.acm.org/10.1145/279227.279229>.
- [5] Barbara Liskov and James Cowling. *Viewstamped Replication Revisited (Technical Report)*. MIT-CSAIL-TR-2012-021. MIT, July 2012.
- [6] Friedemann Mattern. “Virtual Time and Global States of Distributed Systems”. In: *Parallel and Distributed Algorithms*. North-Holland, 1988, pp. 215–226.
- [7] Diego Ongaro and John K. Ousterhout. “In Search of an Understandable Consensus Algorithm”. In: *2014 USENIX Annual Technical Conference, USENIX ATC '14, Philadelphia, PA, USA, June 19-20, 2014*. Ed. by Garth Gibson and Nickolai Zeldovich. USENIX Association, 2014, pp. 305–319. URL: <https://www.usenix.org/conference/atc14/technical-sessions/presentation/ongaro>.

Further reading III

- [8] Michel Raynal, André Schiper, and Sam Toueg. “The Causal Ordering Abstraction and a Simple Way to Implement it”. In: *Inf. Process. Lett.* 39.6 (1991), pp. 343–350. doi: 10.1016/0020-0190(91)90008-6. URL: [https://doi.org/10.1016/0020-0190\(91\)90008-6](https://doi.org/10.1016/0020-0190(91)90008-6).
- [9] Peter Zeller, Annette Bieniusa, and Carla Ferreira. “Teaching practical realistic verification of distributed algorithms in Erlang with TLA+”. In: *Proceedings of the 19th ACM SIGPLAN International Workshop on Erlang, Erlang@ICFP 2020, Virtual Event, USA, August 23, 2020*. Ed. by Annette Bieniusa and Viktória Fördós. ACM, 2020, pp. 14–23. doi: 10.1145/3406085.3409009. URL: <https://doi.org/10.1145/3406085.3409009>.