# Distributed Algorithms

**Benoît Garbinato**
distributed object programming lab

---

# Distributed systems

"As long as there were no ~~machines~~ *networks*, ~~distributed~~ programming was no problem at all; when we had a few weak ~~computers~~ *networks*, ~~distributed~~ programming became a mild problem and now that we have gigantic ~~computers~~ *networks*, *distributed* programming has become an equally gigantic problem. In this sense the electronic industry has not solved a single problem, it has only created them - it has created the problem of using its products."

Edgster Dijkstra, The Humbel Programmer.
Communication of the ACM, vol. 15, no. 10.
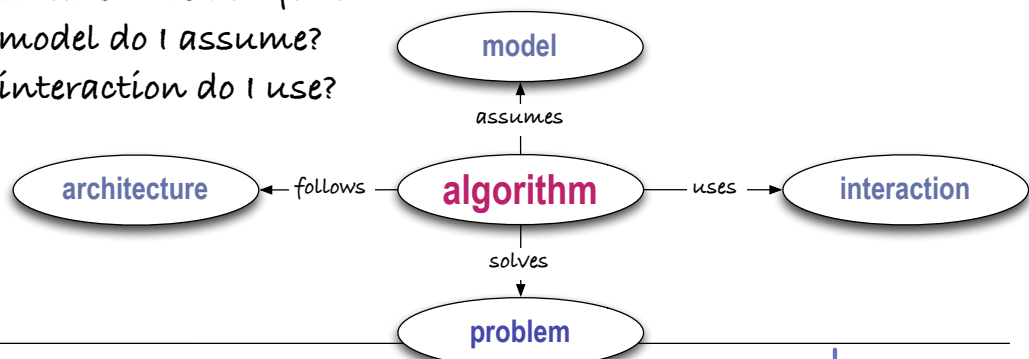October 1972. Turing Award Lecture.

# Our approach

- The practitioner needs the theoretical perspective to understand the implicit assumptions hidden in the technologies, and their consequences

- The theoretician needs the practical perspective to validate that theoretical models, problems & solutions work in accordance to existing technologies

- To achieve this, we approached distributed systems through four complementary views:
  - The model view
  - The interaction view
  - The architecture view
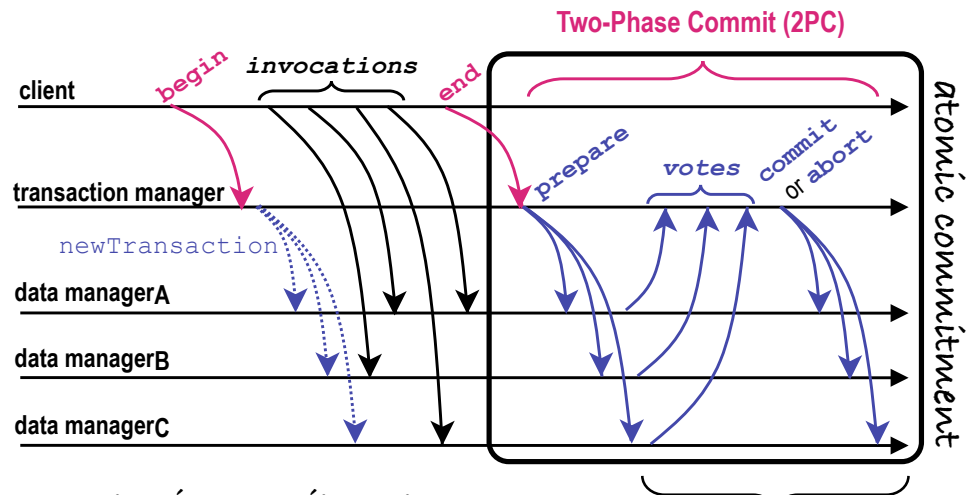  - The algorithm view

dop
l a b

---

# The big picture

When implementing a distributed program, you will always end up writing some algorithm. In doing so, you will have to answer the following questions:

- What problem am I trying to solve?
- What architecture do I follow?
- What model do I assume?
- What interaction do I use?

dop
l a b

# Atomic commitment

**Two-Phase Commit (2PC)**



Problem: atomic commitment
Interaction: reliable message passing
Model: synchronous crash-recovery
Algorithm: 2-phase commit protocol

---

# A few observations

☐ Most atomic commitment protocols guarantee that safety will always hold, but not necessarily liveness

☐ Liveness is compromised when failures prevent the Termination property from holding; in such a case, we say that the protocol is blocking

☐ In the crash-recovery model, a blocking protocol cannot terminate until crashed processes have recovered

☐ Upon recovery, a failed process reads it log file from stable storage and acts according to its last operation

☐ In atomic commitment terms, this implies that the recovering process should be able to decide commit or abort from what it finds in its log file

# Agreement problems

☐ The atomic commitment is an instance of a more general agreement problem, also known as the consensus problem

☐ There exists many variants of the consensus problem, which are not necessarily equivalent to each other

dop : : :
l a b

---

# Problem specification

The atomic commitment problem corresponds to the following consensus variant, with the transaction manager and data managers being processes, value 1 corresponding to commit and value 0 corresponding to abort

<u>Agreement</u>                                                    *(safety property)*
No two processes decide on different values

<u>Validity</u>                                                      *(safety property)*
• If any process starts with 0, then 0 is the only possible decision
• If all processes start with 1 and there are no failures, then 1 is the only possible decision

<u>Termination</u>                                                  *(liveness property)*
  <u>Weak</u>:   if there are no failures, then all processes eventually decide
  <u>Strong</u>: all non faulty processes eventually decide

dop : : :
l a b

# Two-phase commit (2PC)

Premises:
- synchronous model, reliable channels
- crash-recovery failures of data managers $D_i$
- transaction manager T acts as coordinator but also votes

Phase 1:
- each $D_i$ process sends its initial value to process T
- any process $D_i$ whole initial value is 0 decides 0
- if process T times out waiting for some initial value, it decides 0; otherwise it decides for the minimum of all values

Phase 2:
- process T broadcasts its decision to all $D_i$ processes
- any process that has not yet decided adopts this decision

## What <u>Termination</u> property is ensured?

dop
l a b

---

# Upon recovery (2PC)

Premises:

- operations are logged onto stable storage before execution
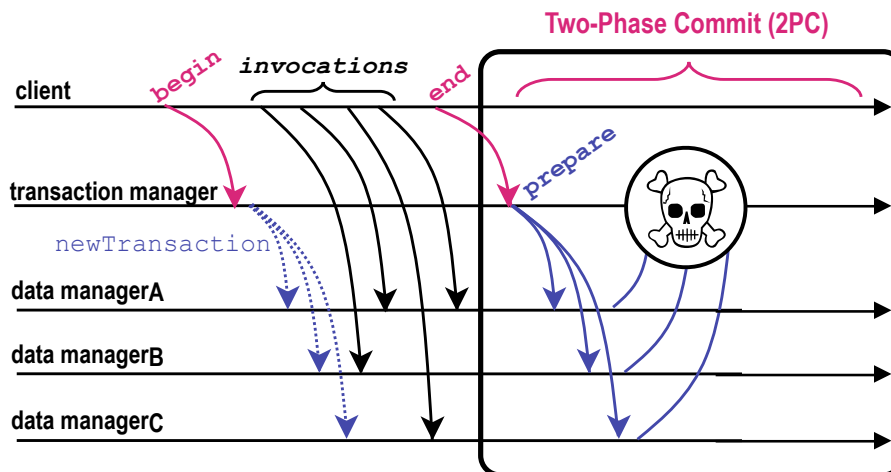- the logging of an operation and its execution are atomic

Recovery of a $D_i$ process:

$D_i$ reads its log file from stable storage

- ▸ if it voted 0 or if it crashed before sending its vote to T, it aborts

- ▸ otherwise, it asks T for the outcome of the transaction and acts accordingly

dop
l a b

# Limits of 2PC

Questions:   what happens if the transaction manager crashes
             before sending the final commit or abort message?

---

# If process T crashed...

Case A:   some process has decided 0

⇒ it knows that T has either not decided or that it has decided 0

⇒ it can inform all other $D_i$ process that it is safe to decide 0

Case B:   all $D_i$ processes have voted 1

⇒ no decision is possible (blocking) :

1.  T might have decided 0, so deciding 1 violates Agreement

2.  T might have decided 1, so deciding 0 violates Agreement

# Three-phase commit (3PC)

Premises:
- synchronous model, reliable channels
- crash-recovery failures of any process
- transaction manager T acts as coordinator but also votes

Phase 1:
- each $D_i$ process sends its initial value to process T
- any process $D_i$ whose initial value is 0 decides 0
- if process T times out waiting for some initial value or receives 0 from some process, it decides 0; otherwise it goes to ready state

Phase 2:
- if process T decided 0, it broadcasts its decision to all $D_i$ processes, so any process that has not yet decided adopts this decision
- if process T is ready state, it broadcasts a pre-commit message, so all processes go to ready state and send an ack message to T
- if process T crashes, the other processes time out and decide 0

Phase 3:
- if process T receives ack messages from all processes, it decides 1 and broadcast its decision, so all processes decide 1 as well
- if process T time out waiting for some ack message, it decides 0 and broadcast its decision, so all processes decide 0 as well
- if process T crashes, the other processes time out and decide 1

# Upon recovery (3PC)

Premises: same as 2PC

Recovery of a $D_i$ process:

$D_i$ reads its log file from stable storage
- ▸ if it voted 0 or if it crashed before acknowledging the pre-commit message, it aborts
- ▸ otherwise, it asks T for the outcome of the transaction and acts accordingly

Recovery of T:

T reads its log file from stable storage:
- ▸ if it crashed before sending pre-commit , it aborts
- ▸ otherwise, it commits

# Limits of 3PC

If $T$ fail in Phase 3, no other process is allowed to fail

<u>Problematic scenario in Phase 3</u>:

1. some $D_i$ crashes before acknowledging pre-commit message
2. $T$ decides 0 but crashes before broadcasting its decision
3. all other $D_i$ time out waiting for the decision and decide 1

$\Rightarrow$ Agreement is violated!

### Why not have all other $D_i$ decide 0 then?

dop
l a b

---
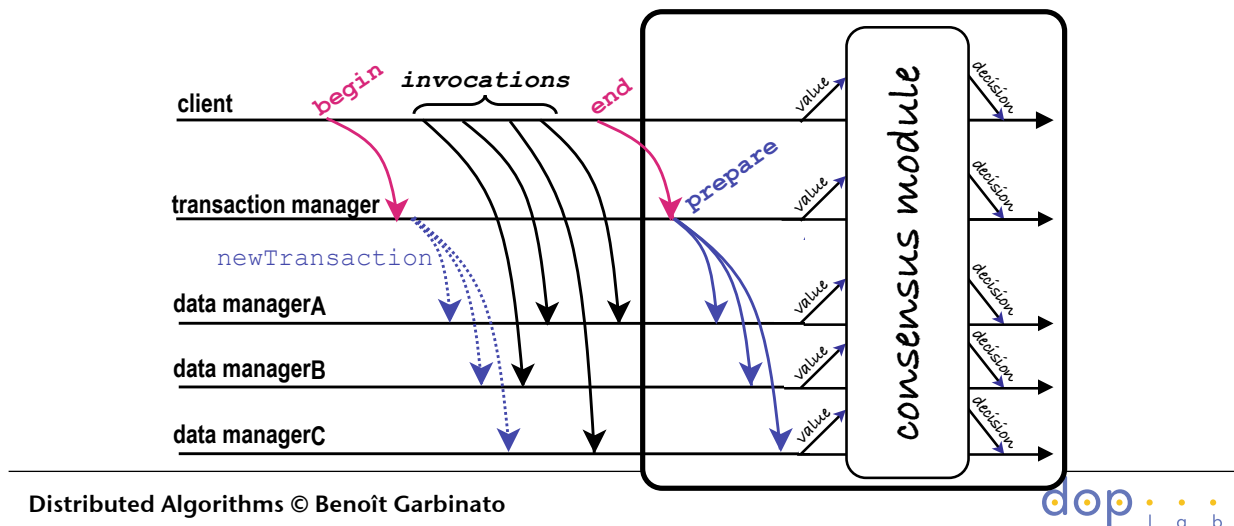
# Further problems

<u>Unrealistic assumptions</u>: synchronous processes and network, no network partitions
$\Rightarrow$ reliable failure detection

<u>Drastic limitation on failures</u>: see previous slide
$\Rightarrow$ $T$ is a single point of failure/vulnerability

<u>Hidden assumptions</u>: logging an action & executing it must be atomic, deciding & broadcasting the decision must also be atomic
$\Rightarrow$ strong underlying atomic mechanisms

dop
l a b

# Back to consensus

☐ If we express the atomic commitment protocol in terms of some consensus module, we can benefit from all the algorithmic work done on the subject

---

# Consensus & asynchrony

☐ Consensus cannot be solved in asynchronous systems; this is the famous Fisher-Lynch-Paterson (FLP) impossibility result

☐ For atomic commitment, the FLP result implies that we cannot answer the question "how long should we wait before aborting?"

 ‣ if we do not wait long enough, safety is at stake
 ‣ if we wait forever, liveness is at stake

☐ Real distributed systems are partially synchronous, i.e., they are mostly synchronous but they experience asynchronous periods every now and then. So, if we can solve a given problem during a synchronous period, that's all we need.
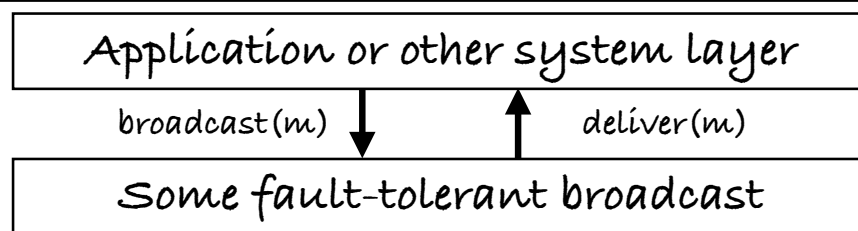
# Failure detectors

☐ A failure detector is a module that provides each process with hints about possible crashes of other processes

☐ A failure detector encapsulates time assumptions and turns them into logical properties: completeness & accuracy. For example, the eventually strong failure detector ($\diamond S$) ensures:

> <u>Strong Completeness</u>. Eventually, every process that crashes is permanently suspected by every correct process.
> <u>Eventual Weak Accuracy</u>. Eventually, there exists a correct process that is never suspected by any correct process

☐ The actual implementability of a given failure detector depends on the underlying timing assumption

---

# Failure detectors & consensus

☐ The $\diamond S$ failure detector was proven to be the weakest failure detector to solve consensus, provided that there are less than half incorrect processes

☐ The algorithm relies on the rotating coordinator paradigm, where a different process has the opportunity to become the next coordinator each time the current coordinator is suspected to have crashed

☐ The <u>Strong Completeness</u> of $\diamond S$ ensures that no process will wait forever for the decision of a crashed coordinator

☐ The <u>Eventual Weak Accuracy</u> of $\diamond S$ ensures that at least one of the coordinators will be able to decide

# Fault-tolerant broadcasts

□ The ability to <u>broadcast messages</u> with some <u>dependable guarantees</u> is a key issue when building fault-tolerant distributed systems

□ Besides the <u>reliable delivery</u> of messages, their <u>ordering</u> is another aspect of this issue

□ For example, if messages represent <u>updates</u> sent to the <u>replicas of a database</u>, reliable delivery and total ordering are necessary

---

| Application or other system layer |
|---|

broadcast(m) ↓    ↑ deliver(m)

| Some fault-tolerant broadcast |
|---|

dop
l a b

---

# Reliable broadcast (basis)

In the following, we assume that each message m includes (1) the identity of the sender, written <u>sender(m)</u> , and (2) a sequence number, denoted <u>seq#(m)</u>. These two fields are what makes each message unique.

<u>Validity</u>
If a correct process broadcasts a message m, then it eventually delivers m

<u>Agreement</u>

<u>Standard</u>: If a correct process delivers a message m, then all correct processes eventually deliver m

<u>Uniform</u>:  If a process delivers a message m, then all correct processes eventually deliver m

<u>Integrity</u>
For any message m, every correct process delivers m at most once, and only if m was previously broadcasted by sender(m)

dop
l a b

# Fifo broadcast

To obtain the specification of fifo broadcast, we simply add the following fifo order property to the aforementioned validity, agreement and integrity properties. That is,
<u>fifo broadcast ⇔ reliable broadcast + fifo order</u>

<u>Fifo order</u>

If a process broadcasts a message m before it broadcasts a message m', then no correct process delivers m' unless it has previously delivered m

dop
l a b

---

# Atomic broadcast (total order)

To obtain the specification of atomic broadcast, we simply add the following total order property to the aforementioned validity, agreement and integrity properties. That is,
<u>atomic broadcast ⇔ reliable broadcast + total order</u>

<u>Total order</u>

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

<u>Question</u>: does this imply Fifo Order ?

dop
l a b

# Causal broadcast

- Very often, perfectly <u>synchronized clocks</u> are <u>not available</u>, due to drifts, impreciseness, etc.

- However, physical time of not necessarily what we need: only <u>causality relationships between events</u> often need to be preserved

- In this context, an <u>event</u> is typically the <u>sending or the reception</u> of some message
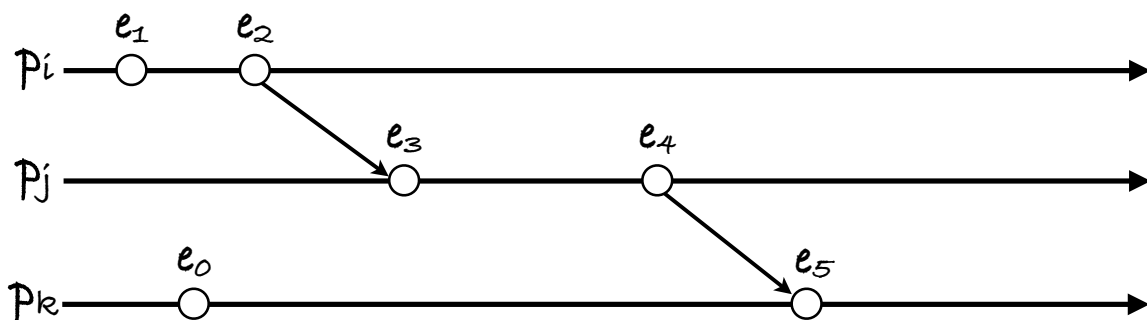
# Causality relationship (1)

- In order to specify the causal broadcast, we must first introduce a <u>partial order relationship</u>

- Let $\to_\ell$ be a partial order on the set of events expressing direct dependencies such that:
  - Let $e_1$ and $e_2$ be two events occurring at the same process p: $e_1 \to_\ell e_2$ if and only if $e_1$ happened before $e_2$ at process p
  - In particular, we have that for each message m: $send(m) \to_\ell receive(m)$

# Causality relationship (2)

- ☐ We now define the causal ordering relationship, noted $\rightarrow_C$ , as the transitive closure of $\rightarrow_\ell$

- ☐ Note that $\rightarrow_C$ also defines a partial order and is sometimes called the happened-before relationship

- ☐ Let $e_1$ and $e_2$ be two events occurring anywhere in the system, i.e., possibly at two distinct processes, we say that $e_1$ causally precedes $e_2$ if and only if we have $e_1 \rightarrow_C e_2$

---

# Illustration of causality



- ☐ Here we have $e_1 \rightarrow_C e_5$ , via $e_2$ , $e_3$ and $e_4$

- ☐ However, $e_1$ and $e_0$ are concurrent, i.e., they are not ordered (hence $\rightarrow_C$ is a partial order)

# Causal broadcast (partial order)

We now specify causal broadcast by simply adding the causal order property given hereafter (based on the happened-before partial order) to the reliable broadcast properties

> ### Causal order
>
> If the broadcast of a message m <u>causally precedes</u> the broadcast of a message m', then no correct process delivers m' unless it has previously delivered m

So: <u>causal broadcast</u> ⇔ <u>reliable broadcast + causal order</u>

---

# Causal broadcast (alternative)

We can also see causal order as a <u>generalization of fifo order</u>. In this case, we define causal broadcast by adding the <u>local order</u> property given hereafter to the fifo broadcast properties

> ### Local order
>
> If a process broadcasts a message m and a process delivers m before broadcasting m', then no correct process delivers m' unless it has previously delivered m.

So: <u>causal broadcast</u> ⇔ <u>fifo broadcast + local order</u>

# Relationship among broadcasts

```
┌─────────────┐   total order   ┌─────────────┐
│  Reliable   │ ──────────────► │   Atomic    │
│  Broadcast  │                 │  Broadcast  │
└─────────────┘                 └─────────────┘
       │ fifo order                   │ fifo order
       ▼                              ▼
┌─────────────┐   total order   ┌─────────────┐
│    Fifo     │ ──────────────► │ Fifo Atomic │
│  Broadcast  │                 │  Broadcast  │
└─────────────┘                 └─────────────┘
       │ local order                  │ local order
       ▼                              ▼
┌─────────────┐   total order   ┌──────────────┐
│   Causal    │ ──────────────► │Causal Atomic │
│  Broadcast  │                 │  Broadcast   │
└─────────────┘                 └──────────────┘
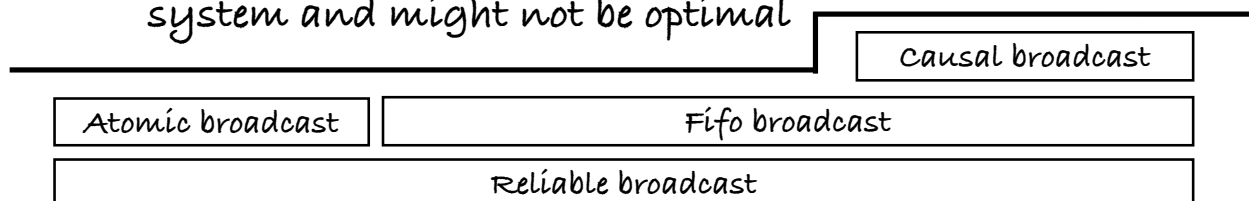```

causal order (left)   causal order (right)

# Implementing broadcasts

☐ There exists numerous algorithms solving the various broadcast primitives we presented

☐ The algorithms we are presenting hereafter are taken from two major papers:

> **[Hadzilacos93]**  Hadzilacos, V. and Toueg, S. 1993. *Fault-tolerant broadcasts and related problems*. In Distributed Systems (2nd Ed.), S. Mullender, Ed. Acm Press Frontier Series. ACM Press/Addison-Wesley Publishing Co., New York, NY, 97-145.

> **[Chandra96]**  Chandra, T. D. and Toueg, S. 1996. *Unreliable failure detectors for reliable distributed systems*. J. ACM 43, 2 (Mar. 1996), 225-267.

☐ These algorithms all assume a partially synchronous system and might not be optimal

| Causal broadcast |
| Atomic broadcast | Fifo broadcast |
| Reliable broadcast |

# Reliable broadcast

Algorithm for process $p$:
To execute broadcast(R, $m$):
    send($m$) to $p$

deliver(R, $m$) occurs as follows:
    **upon** receive($m$) **do**
        **if** $p$ has not previously executed deliver(R, $m$)
        **then**
            send($m$) to all neighbors
            deliver(R, $m$)

**[Hadzilacos93]**

Every process $p$ executes the following:

To execute R-broadcast($m$):
    send $m$ to all (including $p$)

R-deliver($m$) occurs as follows:
    **when** receive $m$ for the first time
        **if** $sender(m) \neq p$ **then** send $m$ to all
        R-deliver($m$)

**[Chandra96]**

Comment: This is typically a flooding algorithm

dop · · ·
    l  a  b

---

# Fifo broadcast

Algorithm for process $p$:
Initialization:
    $msgSet := \emptyset$
    $next[s] := 1,$   for each process $s$

To execute broadcast(F, $m$):
    broadcast(R, $m$)

deliver(F, $-$) occurs as follows:
    **upon** deliver(R, $m'$) **do**
        $s := sender(m')$
        **if** $next[s] = seq\#(m')$
        **then**
            deliver(F, $m'$)
            $next[s] := next[s] + 1$
            **while** ($\exists m \in msgSet : sender(m) = s$
                **and** $next[s] = seq\#(m)$) **do**
            deliver(F, $m$)
            $next[s] := next[s] + 1$
        **else**
            $msgSet := msgSet \cup \{m'\}$

**[Hadzilacos93]**

dop · · ·
    l  a  b

# Causal broadcast

*Algorithm for process p:*
*Initialization:*
$rcntDlvrs := \bot$

*To execute* $\mathtt{broadcast}(\mathtt{C}, m)$:
$\mathtt{broadcast}(\mathtt{F}, \langle rcntDlvrs \,\|\, m \rangle)$
$rcntDlvrs := \bot$

$\mathtt{deliver}(\mathtt{C}, -)$ *occurs as follows:*
    **upon** $\mathtt{deliver}(\mathtt{F}, \langle m_1, m_2, \ldots, m_l \rangle)$ for some $l$ **do**
        **for** $i := 1..l$ **do**
            **if** $p$ has not previously executed $\mathtt{deliver}(\mathtt{C}, m_i)$
            **then**
                $\mathtt{deliver}(\mathtt{C}, m_i)$
                $rcntDlvrs := rcntDlvrs \,\|\, m_i$

**[Hadzilacos93]**

## Comments:

- rcntDelvrs is the sequence of messages that p delivered since it previous causal broadcast
- || is the concatenation operator on sequences of messages

---

# Back to consensus…

The atomic broadcast can be reduced to the consensus problem. Note however that this version of consensus is different from the version we used when discussing atomic commitment. This second version is defined in terms of two primitives, <u>propose(v)</u> and <u>decide(v)</u>, with v some value. When some process executes propose(v), we say that it proposes value v, and when it executes decides(v), we say it decides value v.

<u>Termination.</u> Every correct process eventually decides on some value.

<u>Uniform integrity.</u> Every process decides at most once.

<u>Agreement</u>. No two correct processes decide differently.

<u>Uniform validity</u>. If a process decides v, then v was proposed by some process.

# Atomic broadcast

Initialization:
$R\_delivered := \emptyset$
$A\_delivered := \emptyset$
$k := 0$

To execute $\texttt{broadcast}(\texttt{A}, m)$:

$\texttt{broadcast}(\texttt{R}, m)$

$\texttt{deliver}(\texttt{A}, -)$ occurs as follows:

upon $\texttt{deliver}(\texttt{R}, m)$ do
$R\_delivered := R\_delivered \cup \{m\}$

do forever
$A\_undelivered := R\_delivered - A\_delivered$
if $A\_undelivered \neq \emptyset$ then
$k := k + 1$
$\texttt{propose}(k, A\_undelivered)$
wait for $\texttt{decide}(k, msgSet)$
$batch(k) := msgSet - A\_delivered$
A-deliver all messages in $batch(k)$ in some deterministic order
$A\_delivered := A\_delivered \cup batch(k)$

**[Hadzilacos93]**

Initialisation:

$R\_delivered \leftarrow \emptyset$
$A\_delivered \leftarrow \emptyset$
$k \leftarrow 0$

To execute $A\text{-}broadcast(m)$:

$R\text{-}broadcast(m)$

$A\text{-}deliver(-)$ occurs as follows:

when $R\text{-}deliver(m)$
$R\_delivered \leftarrow R\_delivered \cup \{m\}$

when $R\_delivered - A\_delivered \neq \emptyset$
$k \leftarrow k + 1$
$A\_undelivered \leftarrow R\_delivered - A\_delivered$
$propose(k, A\_undelivered)$
**wait until** $decide(k, msgSet^k)$
$A\_deliver^k \leftarrow msgSet^k - A\_delivered$
atomically deliver all messages in $A\_deliver^k$ in some deterministic order
$A\_delivered \leftarrow A\_delivered \cup A\_deliver^k$

**[Chandra96]**

Comment: consensus execution are numbered and ordered (k)

dop
l a b