

Advanced Course

Distributed Systems

Reliable Broadcast



COURSE TOPICS



- ▶ Intro to Distributed Systems
- ▶ Basic Abstractions and Failure Detectors
- ▶ Reliable and Causal Order Broadcast
- ▶ Distributed Shared Memory
- ▶ Consensus (Paxos, Raft, etc.)
- ▶ Dynamic Reconfiguration
- ▶ Time Abstractions and Interval Clocks (Spanner etc.)
- ▶ Consistent Snapshotting (Stream Data Management)
- ▶ Distributed ACID Transactions (Cloud DBs)

Quorums

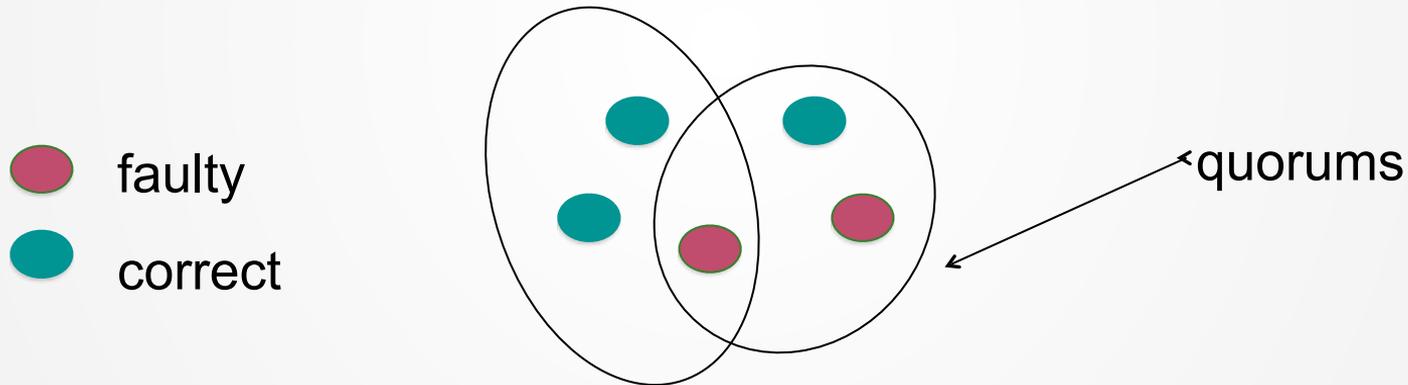
QUORUMS

- For N crash-stop processes
- Quorum is **any set of majority of processes**
- i.e., a set with at least $\lfloor N/2 \rfloor + 1$ processes

- The algorithms will rely on a majority of processes will not fail
 - $f < N/2$ (f is the max number of faulty processes)
- f is the **resilience** of the algorithm

QUORUMS CRASH-STOP/RECOVERY MODEL - $F < N/2$

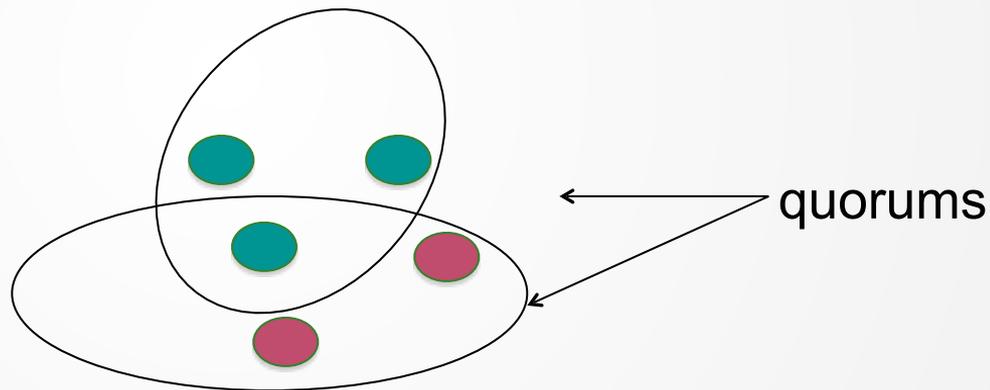
Two quorums always intersect in at least ONE process



QUORUMS CRASH-STOP/RECOVERY MODEL - $F < N/2$

There is at least ONE quorum with only correct processes

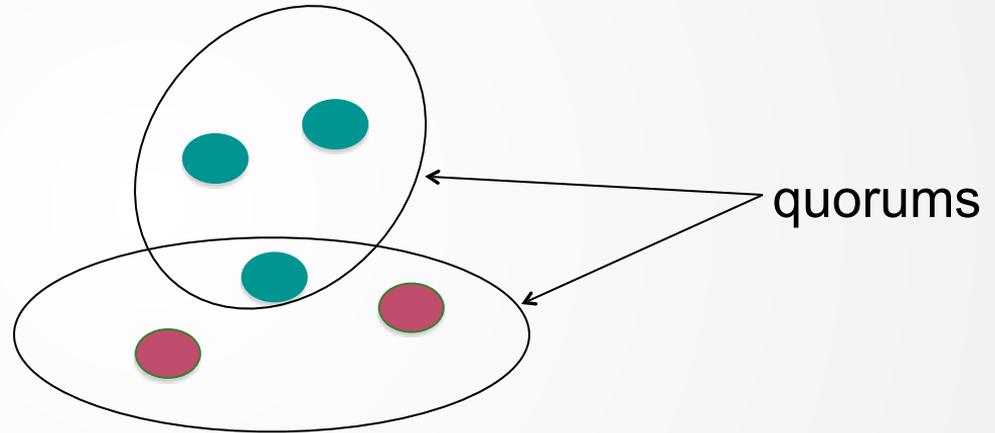
-  faulty
-  correct



QUORUMS CRASH-STOP/RECOVERY MODEL - $F < N/2$

There is at least ONE correct process in each quorum

-  faulty
-  correct

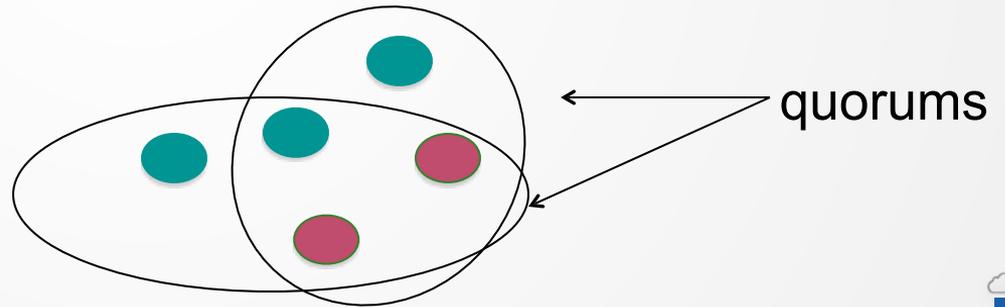


QUORUMS

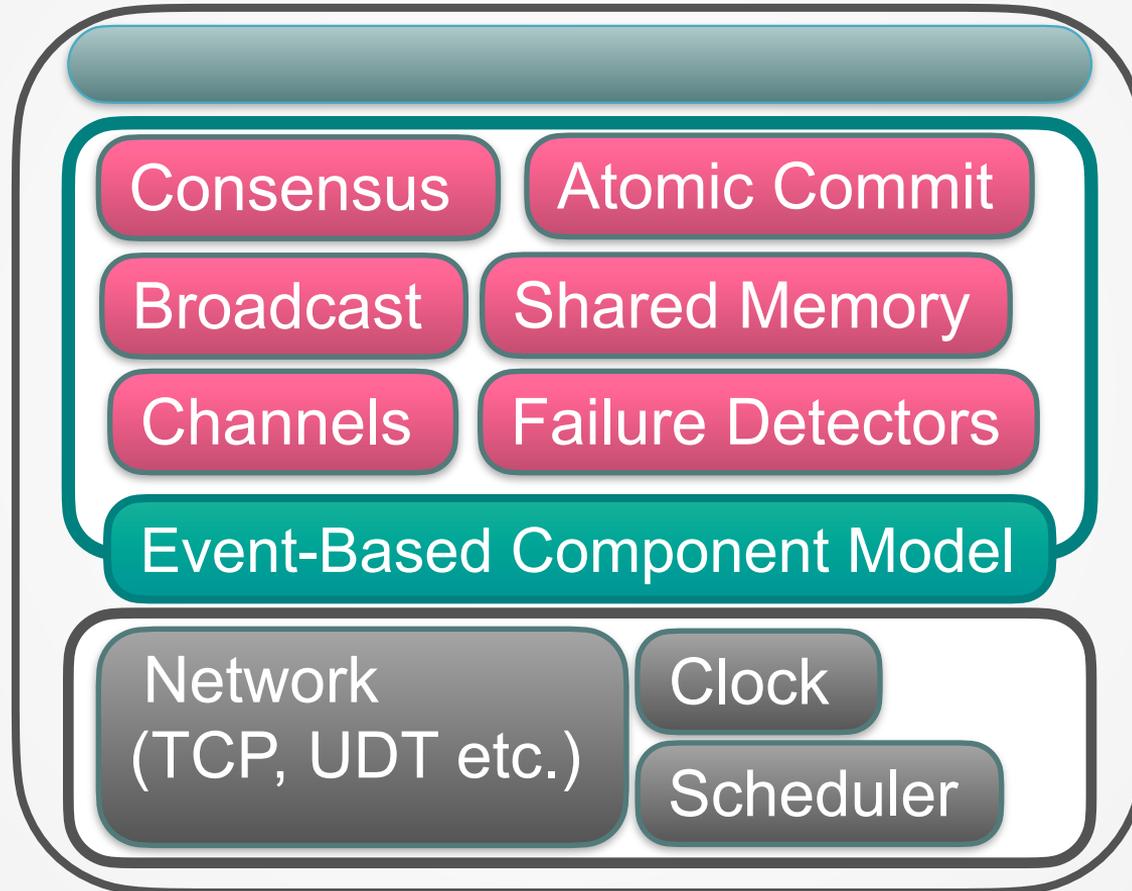
Quorums used in Fail-Silent and Fail-Noisy algorithms

A process never waits for messages from more than $\lfloor N/2 \rfloor + 1$ (different) processes

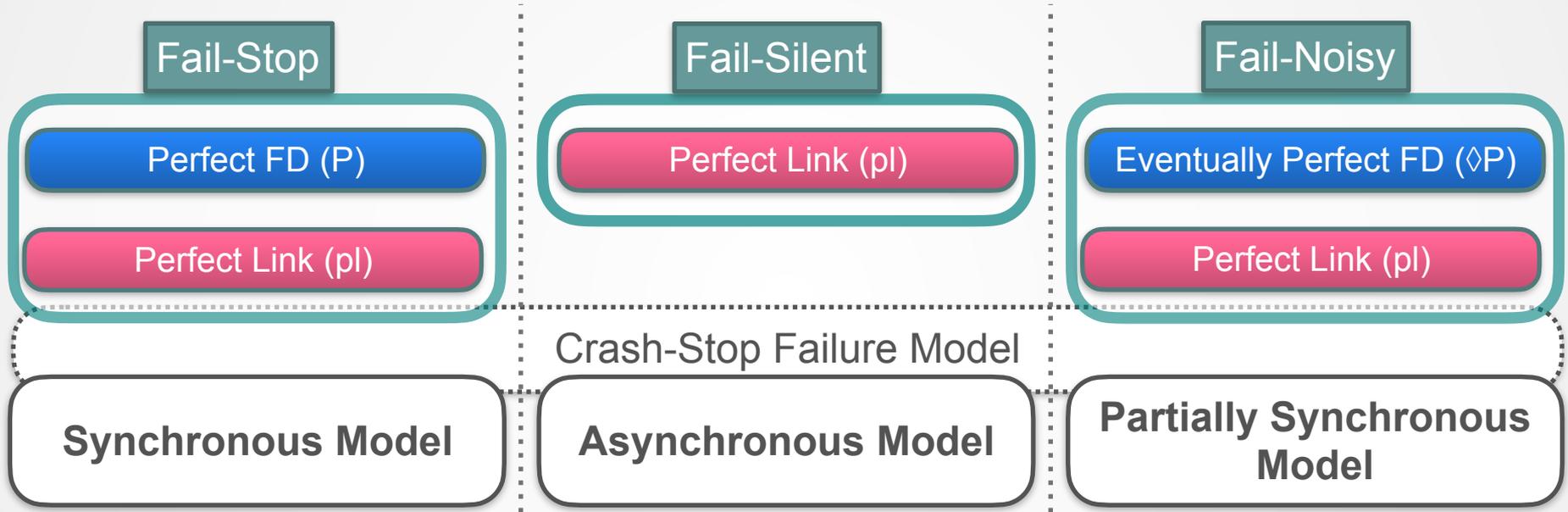
-  faulty
-  correct



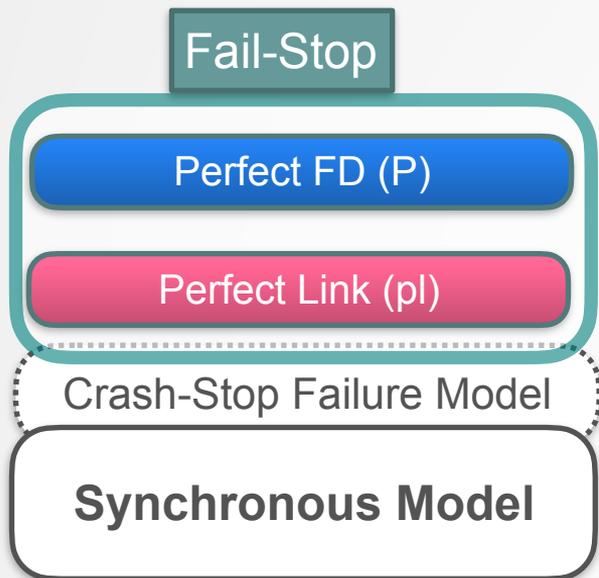
LET'S DEFINE OUR BUNDLES



LET'S MAKE SOME BUNDLES



THE FAIL-STOP



- **How we work with it**
- Local algorithms can **track the set of correct processes at any time.**
- Without violating liveness properties: use
 - Request/Reply protocols.
 - Wait for **correct** processes to reply.

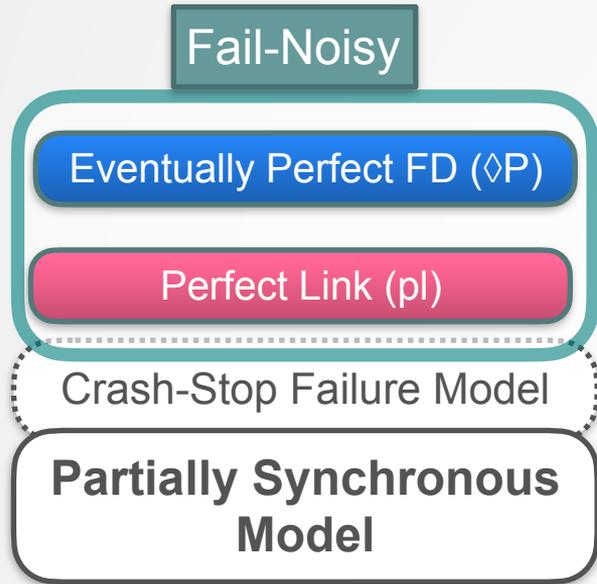
THE FAIL-SILENT



- **How we work with it**

- Failure detection is impossible.
- **Correctness assumptions:** a majority of processes are always correct.
- Protocols work with **majority quorums**.
 - Expect at least $\lceil n/2 \rceil + 1$ responses.

THE FAIL-NOISY

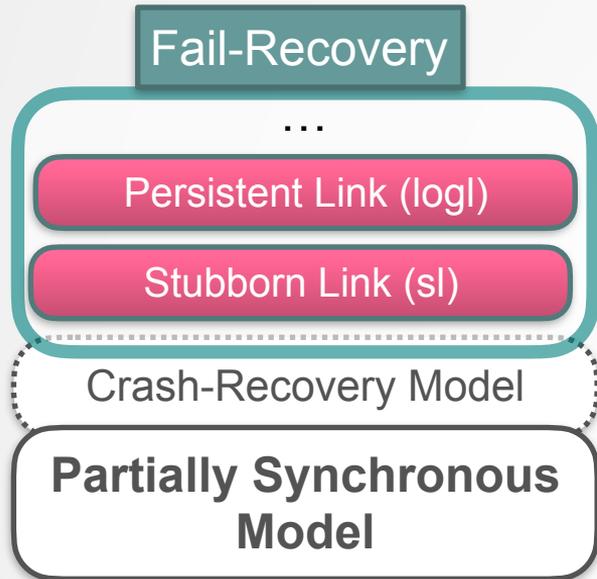


- **Key ideas:**

- To guarantee safety properties any algorithm has to assume the failure detector can be **inaccurate**.
- Eventual strong accuracy is only used to guarantee **liveness**.

Quorum-based ideas also apply here.

A FAIL-RECOVERY BUNDLE



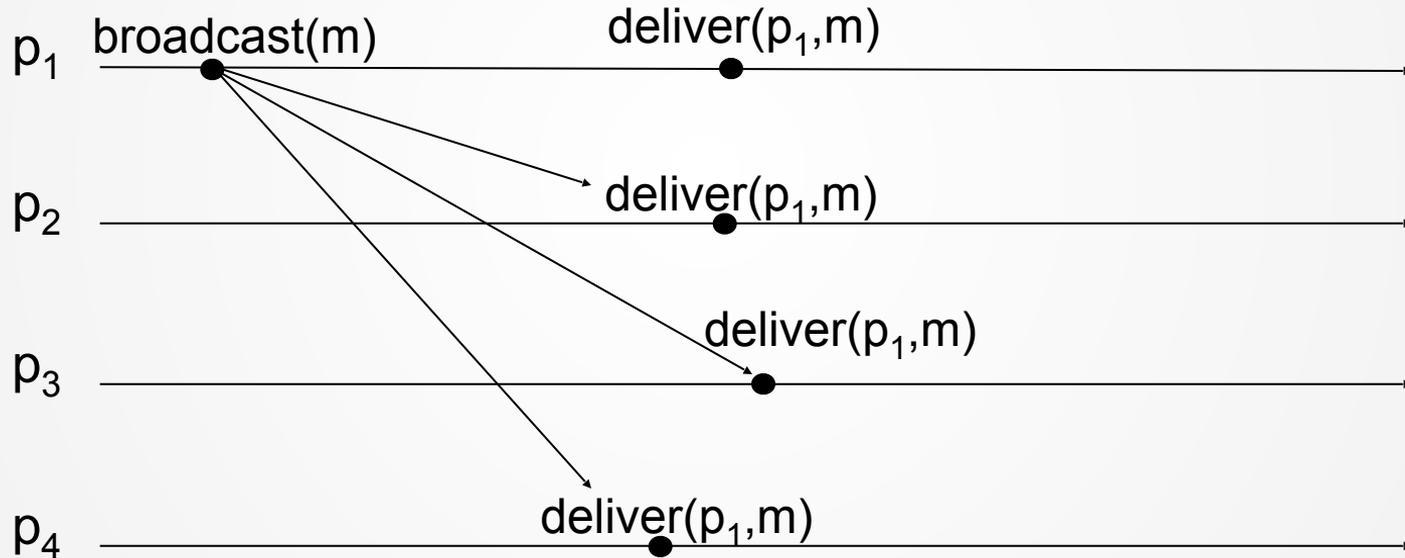
- **Key ideas:**

- Relies often on a **persistent memory** to store and retrieve critical information
- After recovery a process may contact other process to retrieve up to date state information
- Some algorithms **relax** the reliability conditions on channels allowing message loss/duplication/reordering

Broadcast Abstractions

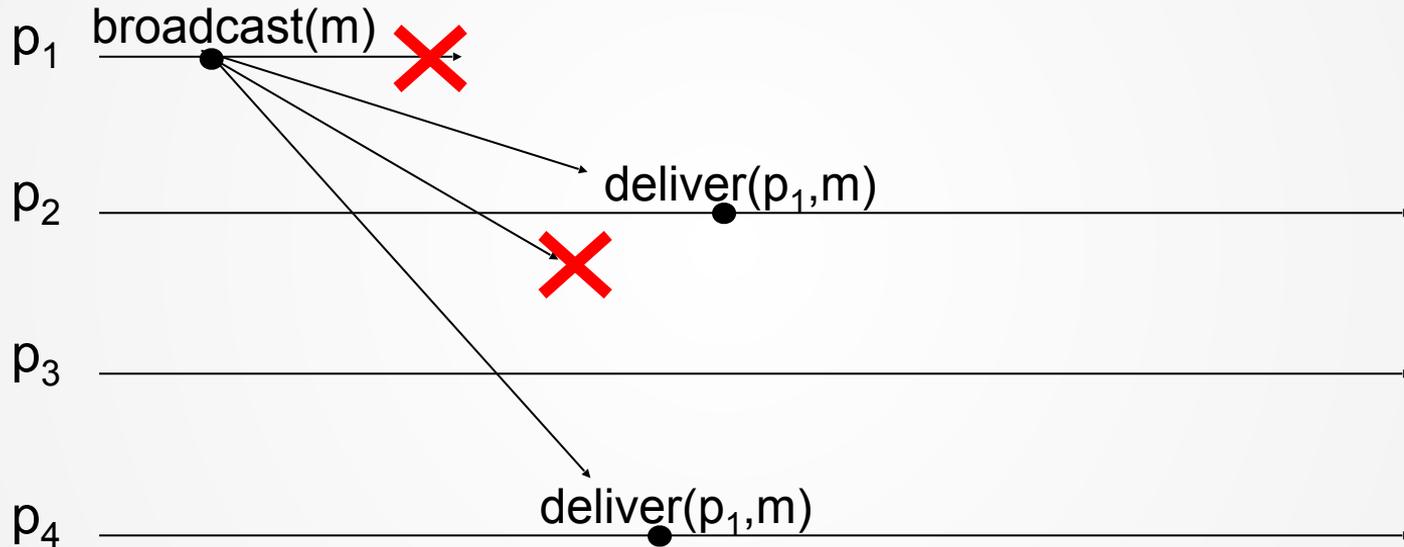
BROADCAST SERVICES

Send a message to a group of processes



UNRELIABLE BROADCAST

✗ crash event



RELIABLE BROADCAST ABSTRACTIONS

- **Best-effort broadcast**
 - Guarantees reliability **only if sender is correct**
- **Reliable broadcast**
 - Guarantees reliability **independent of whether sender is correct**
- **Uniform reliable broadcast**
 - Also **considers behaviour of failed nodes**
- **FIFO reliable broadcast**
 - Reliable broadcast with **FIFO delivery order**
- **Causal reliable broadcast**
 - Reliable broadcast with **causal delivery order**

RELIABLE BROADCAST ABSTRACTIONS

- **Probabilistic reliable broadcast**
 - Guarantees reliability **with high probability**
 - Scales to large number of nodes
- **Total order (atomic) reliable broadcast**
 - Guarantees reliability **and** same order of delivery

Specification of Broadcast Abstractions

BEST-EFFORT BROADCAST (BEB)

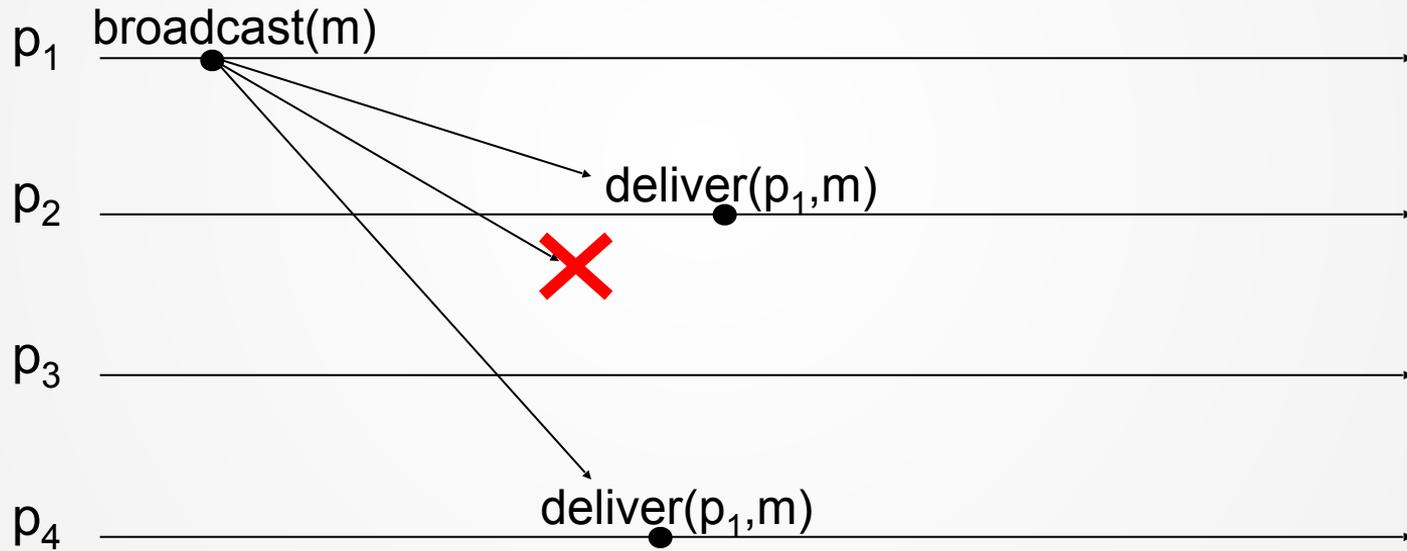
- *Instance beb*
- *Events*
 - Request: $\langle \text{beb Broadcast} \mid m \rangle$
 - Indication: $\langle \text{beb Deliver} \mid \text{src}, m \rangle$
- *Properties: BEB1, BEB2, BEB3*

BEST-EFFORT BROADCAST (BEB)

- *Intuitively*: everything perfect unless sender crash
- *Properties*
 - **BEB1. Best-effort-Validity**: If p_i and p_j are **correct**, then any broadcast by p_i is eventually delivered by p_j
 - **BEB2. No duplication**: No message delivered more than once
 - **BEB3. No creation**: No message delivered unless broadcast

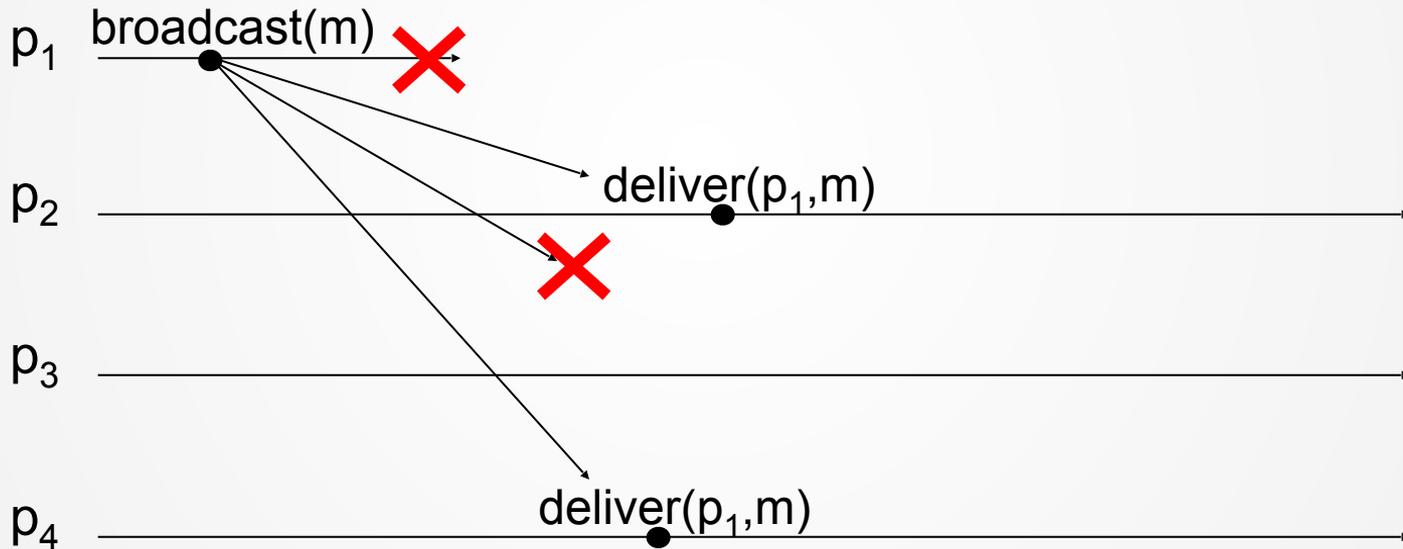
BEB EXAMPLE

Is this allowed? **No**



BEB EXAMPLE (2)

Is this allowed? **Yes**



RELIABLE BROADCAST

- BEB gives no guarantees if **sender crashes**
 - Strengthen to give guarantees if sender crashes
- Reliable Broadcast Intuition
 - Same as BEB, plus
 - If sender crashes:
 - ensure *all or none* of the correct nodes get msg

RELIABLE BROADCAST (RB)

Instance rb

Events

Request: $\langle \text{rb Broadcast} \mid m \rangle$

Indication: $\langle \text{rb Deliver} \mid \text{src}, m \rangle$

Properties: RB1, RB2, RB3, RB4

RELIABLE BROADCAST PROPERTIES

- **Properties**

- **RB1 = BEB1. Validity**

- **RB2 = BEB2. No duplication**

- **RB3 = BEB3. No creation**

- **RB4. Agreement.**

- If a **correct process delivers** m , then every correct process delivers m

REFINING CORRECTNESS

Can weaken RB1 without any effect

Old Validity

←equivalent with→

New Validity

RB1 = BEB1 Validity

- If p_i and p_j are **correct**, then any broadcast by p_i is eventually delivered by p_j

RB2 = BEB2. No duplication

RB3 = BEB3. No creation

RB4. Agreement.

- If a **correct node delivers** m , then every correct node delivers m

RB1 Validity.

- If **correct** p_i broadcasts m , p_i itself eventually delivers m

RB2 = BEB2. No duplication

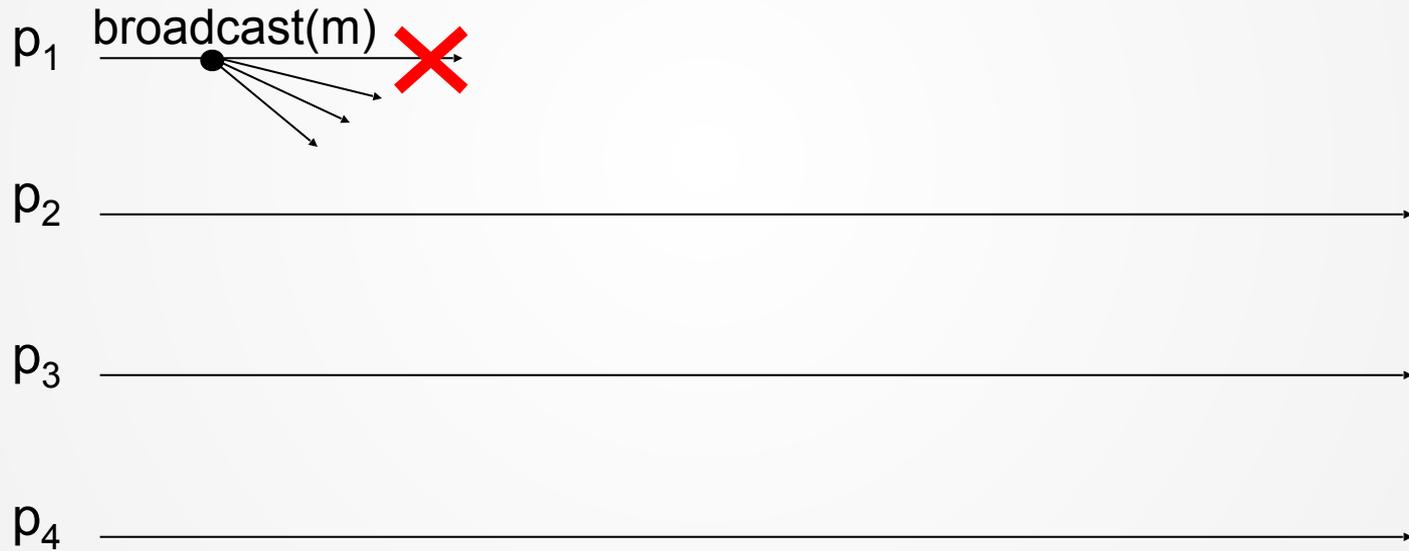
RB3 = BEB3. No creation

RB4. Agreement.

- If a **correct node delivers** m , then every correct process delivers m

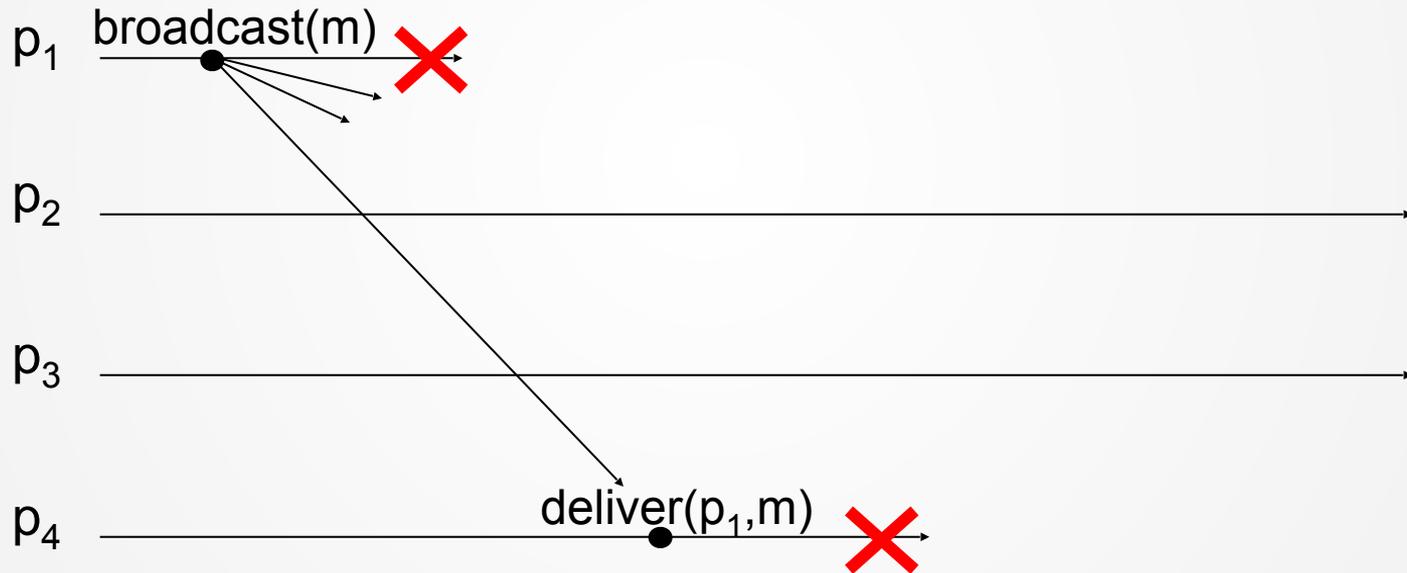
RB EXAMPLE

Is this allowed? **Yes**



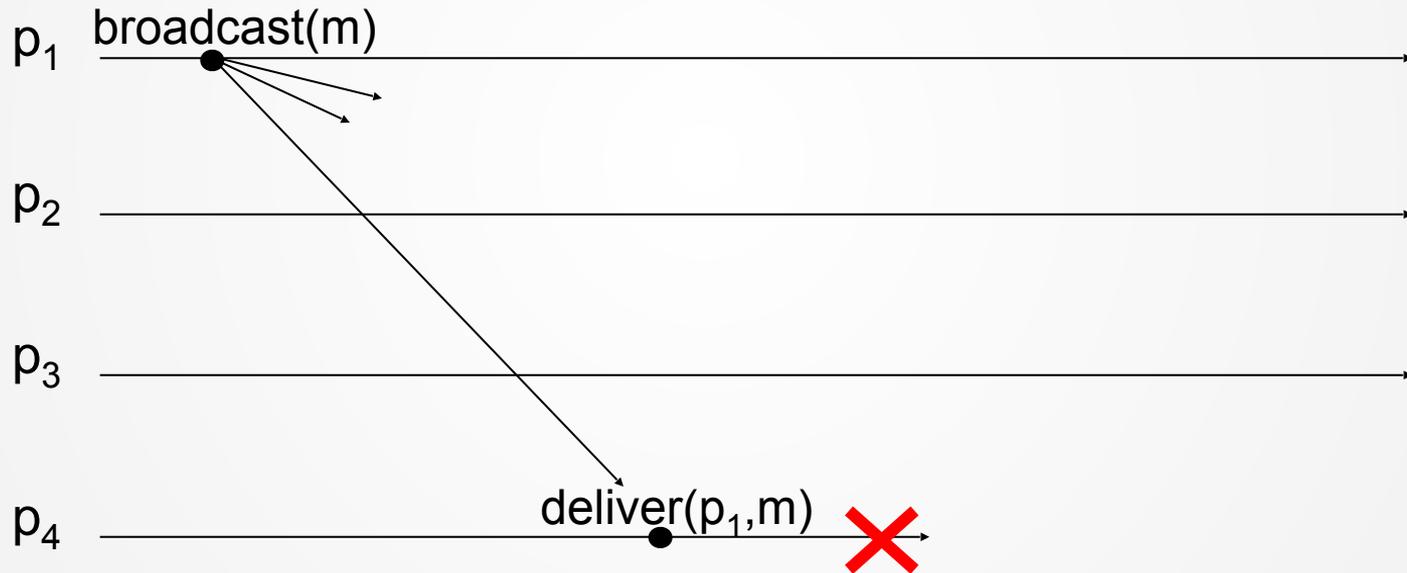
RB EXAMPLE

Is this allowed? **Yes**



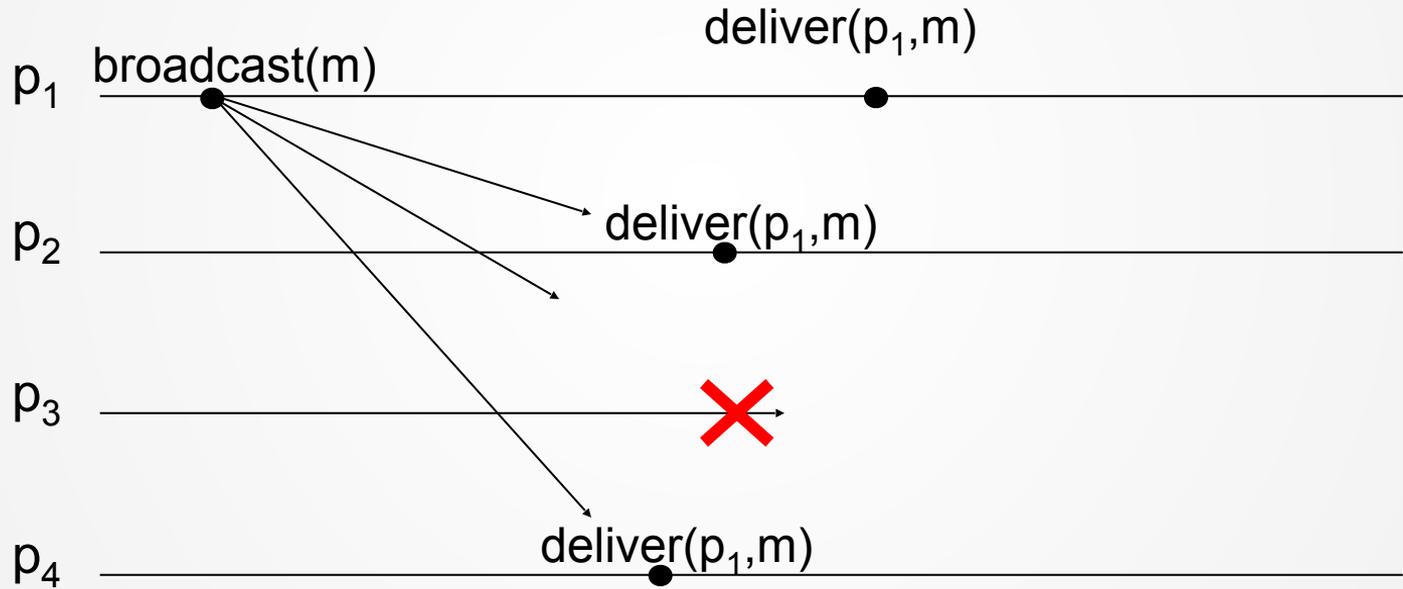
RB EXAMPLE

Is this allowed? **No**



RB EXAMPLE

Is this allowed? **Yes**



UNIFORM RELIABLE BROADCAST

- Assume sender broadcasts message
 - Sender fails
 - No correct process delivers message
 - Some failed processes deliver message
- Assume the broadcast enforces
 - Printing a message on paper
 - Withdrawing money from account
- **Uniform** reliable broadcast intuition
 - If a failed node delivers, everyone must deliver...
 - At least correct nodes, we cannot revive the dead...

UNIFORM BROADCAST (URB)

Events

Request: $\langle \text{urb Broadcast} \mid m \rangle$

Indication: $\langle \text{urb Deliver} \mid \text{src}, m \rangle$

Properties:

URB1

URB2

URB3

URB4

UNIFORM BROADCAST PROPERTIES

Properties

URB1 = RB1.

URB2 = RB2.

URB3 = RB3.

*URB4. **Uniform Agreement:*** For any message m , if **a process delivers** m , then every correct process delivers m

Wanted: Dead & Alive!

Broadcast Abstractions

Implementation of Broadcast Abstractions

IMPLEMENTING BEB

- Use Perfect channel abstraction
 - Upon $\langle \text{beb Broadcast} \mid m \rangle$ send message m to all processes (for-loop)
- Correctness
 - If sender doesn't crash, every other correct process receives message by perfect channels (**Validity**)
 - **No creation** & **No duplication** already guaranteed by perfect channels

Fail-Stop

Lazy Reliable Broadcast

FAIL-STOP: LAZY RELIABLE BROADCAST

- Requires perfect failure detector (**P**)
- To broadcast m :
 - **best-effort broadcast** m
 - When get **beb** Deliver
 - Save message, and
 - **rb** Deliver message
- If **sender** s crash, detect & relay msgs from s to all
 - **case 1**: get m from s , detect crash s , redistribute m
 - **case 2**: detect crash s , get m from s , redistribute m
- Filter duplicate messages before delivery

FAIL-STOP: LAZY RELIABLE BROADCAST

If **sender s** crashes, detect & relay msgs from s to all

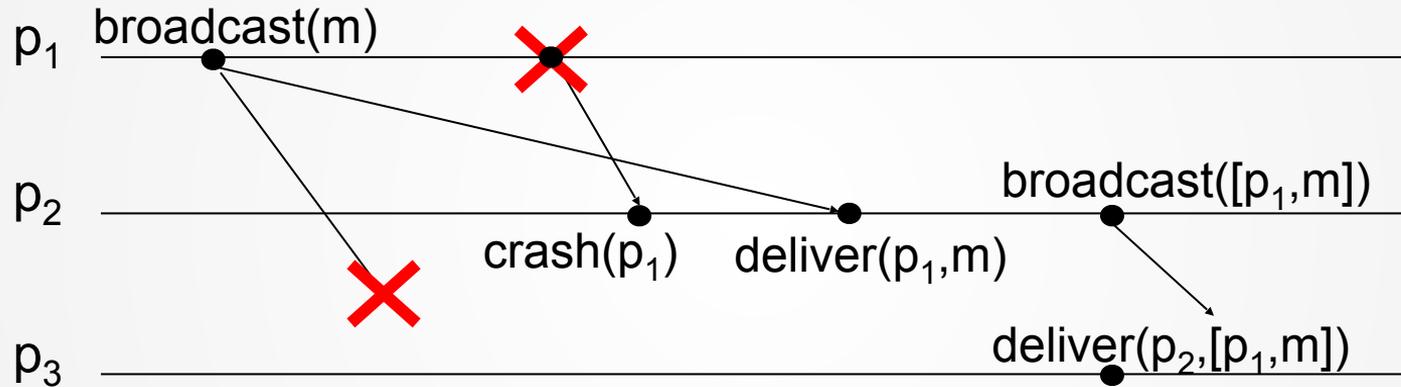
case 1: get m from s, detect crash s, redistribute m

case 2: detect crash s, get m from s, redistribute m

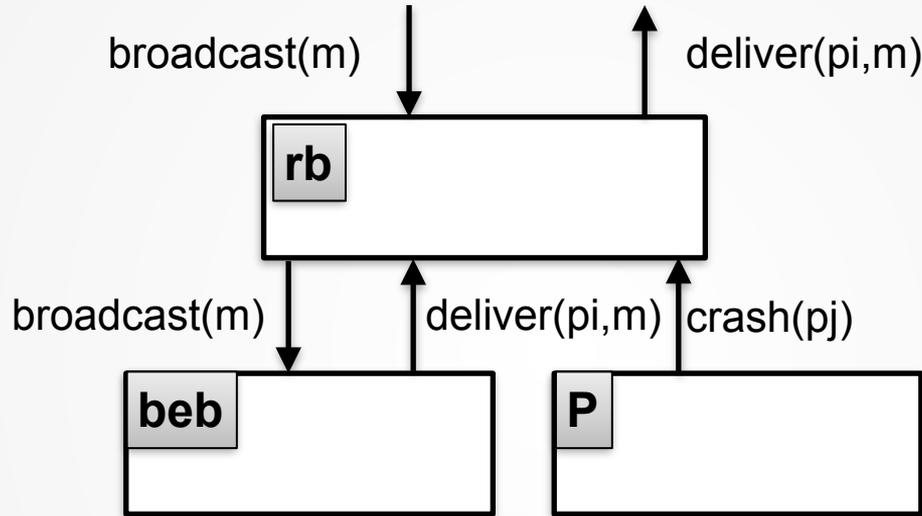
Why case 2? [d]

LAZY RELIABLE BROADCAST

Case 2



FAIL-STOP LAZY RELIABLE BROADCAST



LAZY RELIABLE BROADCAST

Implements: ReliableBroadcast (rb)

Uses:

BestEffortBroadcast (beb)
PerfectFailureDetector (P)

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

delivered := \emptyset

correct := Π

forall $p_i \in \Pi$ **do** from $[p_i]$:= \emptyset

upon event $\langle \text{rb Broadcast} \mid m \rangle$ **do**

trigger $\langle \text{beb Broadcast} \mid (\text{DATA}, \text{self}, m) \rangle$

**for filtering
duplicates**

**storage for saved
messages**

LAZY RELIABLE BROADCAST (2)

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

correct := correct $\setminus \{p_i\}$

forall $(s_m, m) \in \text{from}[p_i]$ **do**

trigger $\langle \text{beb Broadcast} \mid (\text{DATA}, s_m, m) \rangle$

Case 1: redistribute anything we have from failed node

upon event $\langle \text{beb Deliver} \mid p_i, (\text{DATA}, s_m, m) \rangle$ **do**

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

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

from[p_i] := from[p_i] $\cup \{(s_m, m)\}$

trigger $\langle \text{rb Deliver} \mid s_m, m \rangle$

if $p_i \notin \text{correct}$ **then**

trigger $\langle \text{beb Broadcast} \mid (\text{DATA}, s_m, m) \rangle$

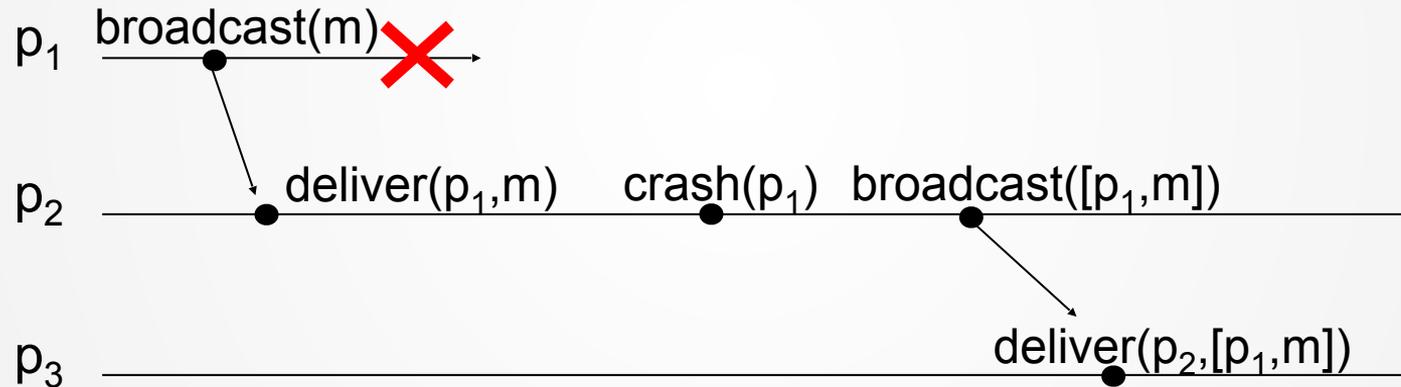
Avoid duplicates

Store for future

Case 2: redistribute

RB EXAMPLE

Which case? **Case 1**



CORRECTNESS OF LAZY RB

- ***RB1-RB3*** satisfied by BEB
- Need to prove ***RB4***
 - If a **correct node delivers** m , then every correct node delivers m
- Assume Correct p_k delivers message bcast by p_i
 - If p_i is correct, BEB ensures correct delivery
 - If p_i crashes,
 - p_k detects this (completeness)
 - p_k uses BEB to ensure (BEB1) every correct node gets it

Measuring Performance

MESSAGE COMPLEXITY

- The number of messages required to terminate an operation of an abstraction
- Lazy reliable broadcast
 - The number of messages initiated by broadcast(m)
 - Until a deliver(src, m) event is issued at each process
- Bit complexity
 - Number of bits sent, if messages can vary in size

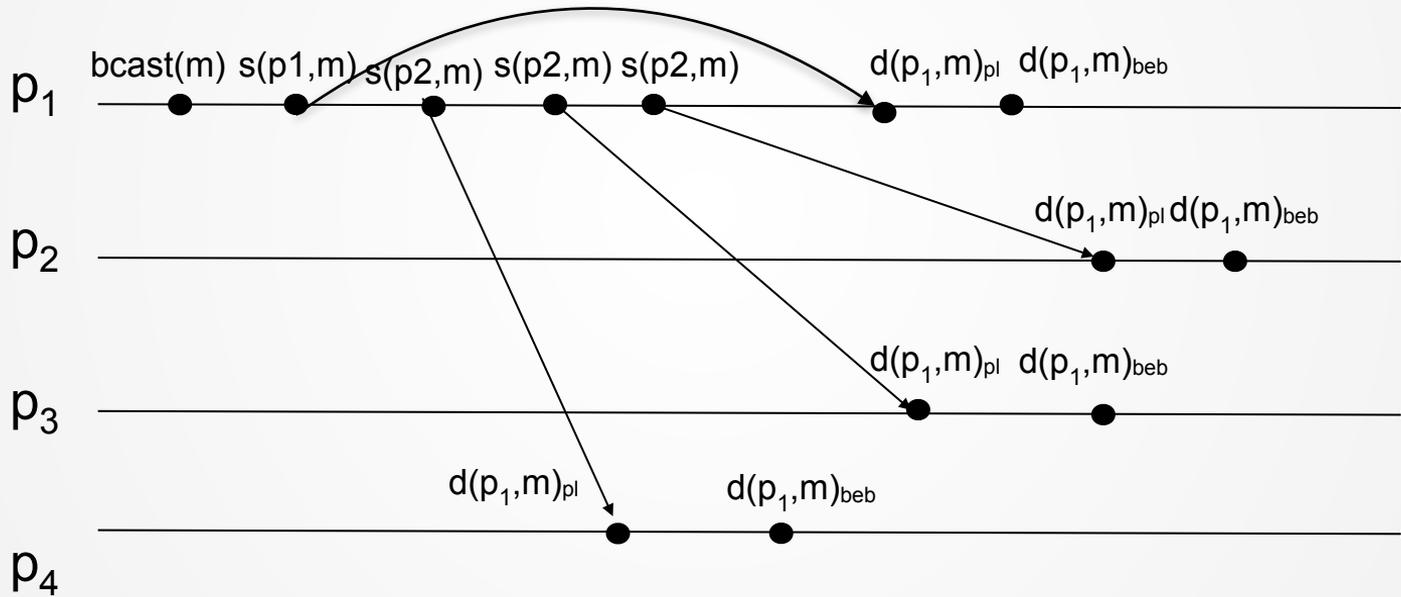
TIME COMPLEXITY ~ #ROUNDS

- **One time unit** in an Execution E is the **longest** message delay in E
- **Time Complexity is** Maximum time taken by any execution of the algorithm under the assumptions
 - A process can execute any finite number of actions (events) in **zero** time
 - The time between $\text{send}(m)_{i,j}$ and $\text{deliver}(m)_{i,j}$ is **at most one** time unit
- In most algorithms we study we assume all communication steps takes one time unit. We also call this a **round or step**.

BEST EFFORT BROADCAST

Takes **one time unit** from $\text{broadcast}(m)_p$ to last deliver(p,m)

We also call it one **communication step / round**.



COMPLEXITY OF LAZY RELIABLE BROADCAST

- Assume N processes
- Message complexity
 - Best case: $O(N)$ messages
 - Worst case: $O(N^2)$ messages
- Time complexity
 - Best case: 1 time unit
 - Worst case: 2 time units

Fail-Silent

Eager Reliable Broadcast

EAGER RELIABLE BROADCAST

What happens if we replace P with $\diamond P$? [d]

- Only affects performance
- Only affects correctness
- No effect
- Affects performance and correctness

EAGER RELIABLE BROADCAST

Can we modify Lazy RB to not use P? [d]

Just assume all processes failed

BEB Broadcast as soon as you get a msg

EAGER RELIABLE BROADCAST

Uses: BestEffortBroadcast (beb)

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

delivered := \emptyset

upon event $\langle \text{rb Broadcast} \mid m \rangle$ **do**

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

trigger $\langle \text{rb Deliver} \mid \text{self}, m \rangle$

trigger $\langle \text{beb Broadcast} \mid (\text{DATA}, \text{self}, m) \rangle$

Immediately deliver

**Immediately BEB
broadcast**

upon event $\langle \text{beb Deliver} \mid p_i, (\text{DATA}, s_m, m) \rangle$ **do**

if $m \notin$ delivered **then**

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

trigger $\langle \text{rb Deliver} \mid s_m, m \rangle$

trigger $\langle \text{beb Broadcast} \mid (\text{DATA}, s_m, m) \rangle$

Immediately deliver

**Immediately BEB
broadcast**

CORRECTNESS OF EAGER RB

- ***RB1-RB3*** satisfied by BEB
- Need to prove ***RB4***
 - If a **correct process delivers** m , then every correct node delivers m
- Assume correct p_k delivers message bcast by p_i
 - p_k uses BEB to ensure (BEB1) every correct process gets it

Uniform Reliable Broadcast

UNIFORMITY

- Is the proposed algorithm also **uniform**? [d]
- Uniformity necessitates
 - If a **failed process** delivers a message m
 - then every correct node delivers m

UNIFORMITY

- No.
 - Sender p immediately RB delivers and crashes
 - Only p delivered message
- **upon event** $\langle \text{rb Broadcast} \mid m \rangle$ **do**
 - $\text{delivered} := \text{delivered} \cup \{m\}$
 - **trigger** $\langle \text{rb Deliver} \mid \text{self}, m \rangle$
 - **trigger** $\langle \text{beb Broadcast} \mid (\text{DATA}, \text{self}, m) \rangle$

UNIFORM EAGER RB

- Necessary condition for uniform RB delivery
 - All correct processes will get the msg
 - How do we know the correct processes got msg? [d]
- Messages are **pending** until all correct processes get it
 - Collect acks from processes that got msg
- Deliver once all correct processes acked
 - Use **perfect FD**
 - **function canDeliver(m):**
 - **return** correct \subseteq ack[m]

← Use vector **ack[m]** at p_i : the set of processes that acked m

UNIFORM EAGER RB IMPLEMENTATION

- **upon event** $\langle \text{urb Broadcast} \mid m \rangle$ **do**

- $\text{pending} := \text{pending} \cup \{(\text{self}, m)\}$
- **trigger** $\langle \text{beb Broadcast} \mid (\text{DATA}, \text{self}, m) \rangle$

remember sent messages

- **upon event** $\langle \text{beb Deliver} \mid p_i, (\text{DATA}, s_m, m) \rangle$ **do**

- $\text{ack}[m] := \text{ack}[m] \cup \{p_i\}$
- **if** $(s_m, m) \notin \text{pending}$ **then**
 - $\text{pending} := \text{pending} \cup (s_m, m)$
 - **trigger** $\langle \text{beb Broadcast} \mid (\text{DATA}, s_m, m) \rangle$

p_i obviously got m

avoid resending

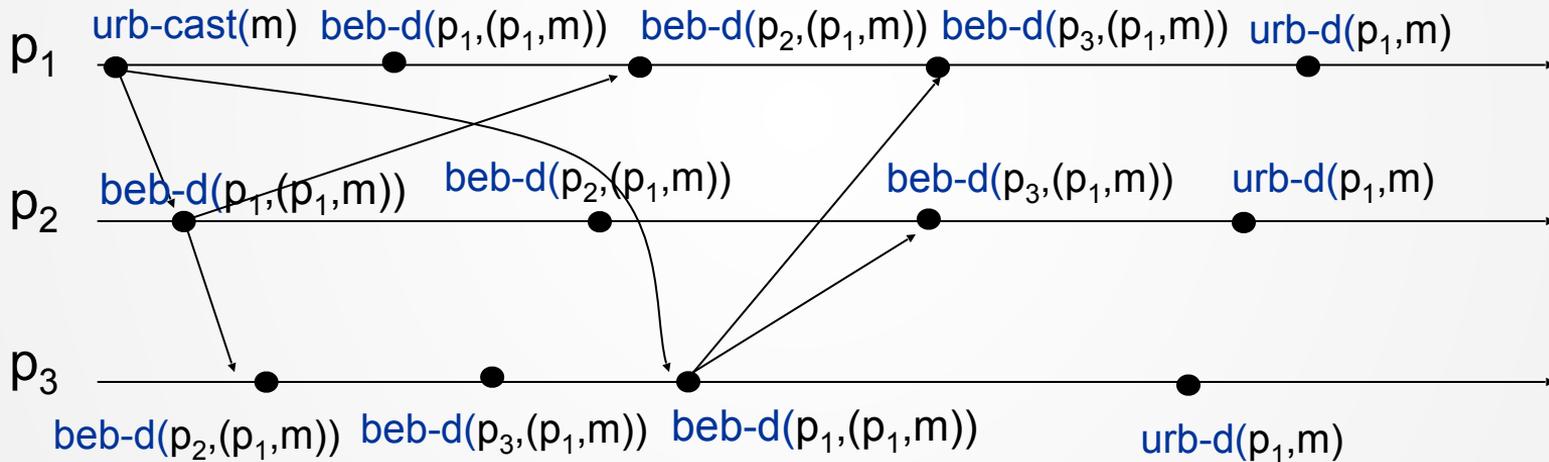
- **Upon exists** $(s_m, m) \in \text{pending}$ **s.t.**

- **canDeliver(m) and** $m \notin \text{delivered}$ **do**

- $\text{delivered} := \text{delivered} \cup \{m\}$
- **trigger** $\langle \text{urb Deliver} \mid s_m, m \rangle$

deliver when all correct nodes have acked

URB EAGER ALGORITHM EXAMPLE



CORRECTNESS OF UNIFORM RB

- **No creation** from BEB
- **No duplication** by using ***delivered*** set
- **Lemma**
 - If a **correct** process p_i bebDelivers m , then p_i eventually urbDelivers m
- **Proof**
 - Correct process p_i bebBroadcasts m as soon as it gets m
 - By BEB1 every correct process gets m and bebBroadcasts m
 - p_i gets bebDeliver(m) from every correct process by BEB1
 - By completeness of **P**, it will not wait for dead nodes forever
 - **canDeliver(m)** becomes true and p_i delivers m

CORRECTNESS OF UNIFORM RB

Validity

If sender s is correct, it'll by **validity** (BEB1)
bebDeliver m

By the **lemma**, it will eventually urbDeliver(m)

CORRECTNESS OF UNIFORM RB

- **Uniform agreement**
 - Assume some process (possibly failed) URB delivers m
 - Then $\text{canDeliver}(m)$ was true,
by **accuracy** of P **every** correct process has BEB delivered m
 - By **lemma** each of the nodes that BEB delivered m will URB deliver m

Uniform Broadcast

Fail-Silent

HOW USEFUL IS THE UNIFORM ALGORITHM?

- Strong failure detectors necessary for URB?
 - No, we'll provide RB for **fail-silent** model
- Assume a **majority** of correct nodes
 - Majority = $\lfloor n/2 \rfloor + 1$, i.e. 6 of 11, 7 of 12...
- Every node eagerly BEB broadcast m
 - URB deliver m when received m from a majority

MAJORITY-ACK UNIFORM RB

- Same algorithm as uniform eager RB
 - Replace one function
 - **function canDeliver(m)**
 - **return** $|ack[m]| > n/2$ ← **majority has acknowledged m**
- Agreement (main idea)
 - If a process URB delivers, it got ack from majority
 - In that majority, one node, p , must be correct
 - p will ensure all correct processes BEB deliver m
 - The correct processes (majority) will ack and URB deliver

MAJORITY-ACK UNIFORM RB

Validity

If correct sender sends m

All correct nodes BEB deliver m

All correct nodes BEB broadcast

Sender receives a majority of acks

Sender URB delivers m

RESILIENCE

- The maximum number of faulty processes an algorithm can handle
- The Fail-Silence algorithm
 - Has resilience less than $N/2$
- The Fail-Stop algorithm
 - Has resilience = $N - 1$