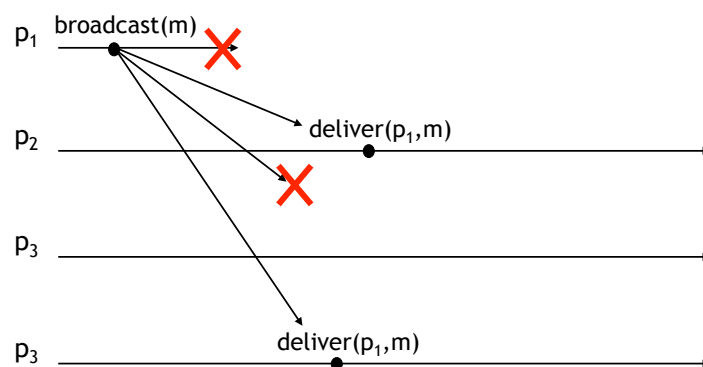


Reliable Broadcast

Seif Haridi - Royal Institute of Technology
Peter Van Roy - Université catholique de Louvain

haridi(at)kth.se
peter.vanroy(at)uclouvain.be

Unreliable Broadcast



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 behavior of failed nodes**
- **FIFO reliable broadcast**
 - Reliable broadcast with **FIFO delivery order**
- **Causal reliable broadcast**
 - Reliable broadcast with **causal delivery order**

19/3/2012

Seif Haridi

3

Specification of Broadcast Abstractions

Best-effort broadcast (beb)

■ **Events**

- Request: $\langle \text{bebBroadcast} \mid m \rangle$
- Indication: $\langle \text{bebDeliver} \mid \text{src}, m \rangle$

■ **Properties: BEB1, BEB2, BEB3**

Best-effort broadcast (beb)

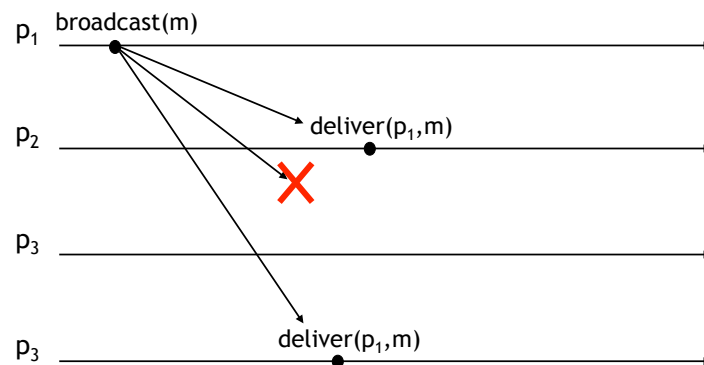
■ **Intuitively:** everything perfect unless sender crash

■ **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**



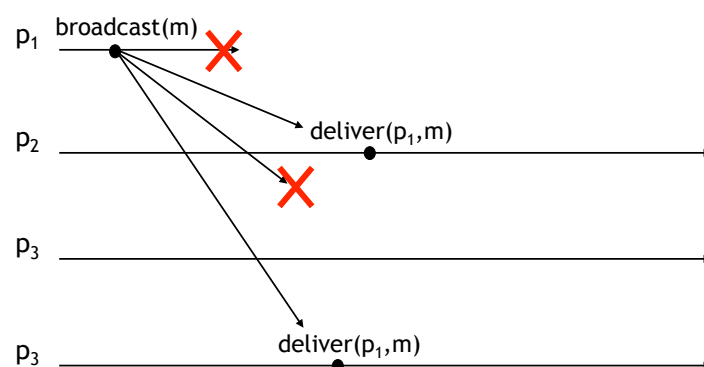
19/3/2012

Seif Haridi

7

BEB Example (2)

■ Is this allowed? **Yes**



19/3/2012

Seif Haridi

8

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)

- **Events**
 - Request: $\langle \text{rbBroadcast} \mid m \rangle$
 - Indication: $\langle \text{rbDeliver} \mid \text{src}, m \rangle$
- **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 node delivers** m , then every correct node delivers m

19/3/2012

Seif Haridi

11

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 node delivers m

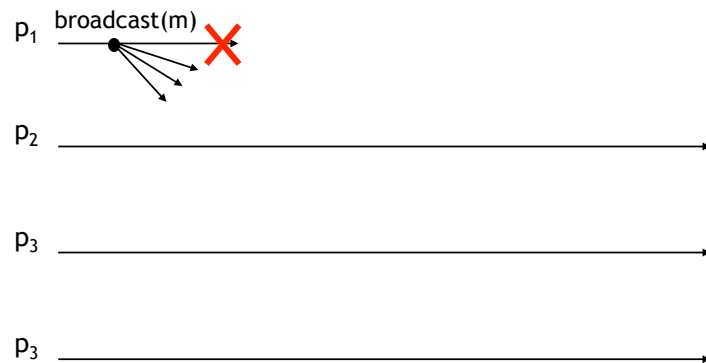
19/3/2012

Seif Haridi

12

RB Example

- Is this allowed? **Yes**



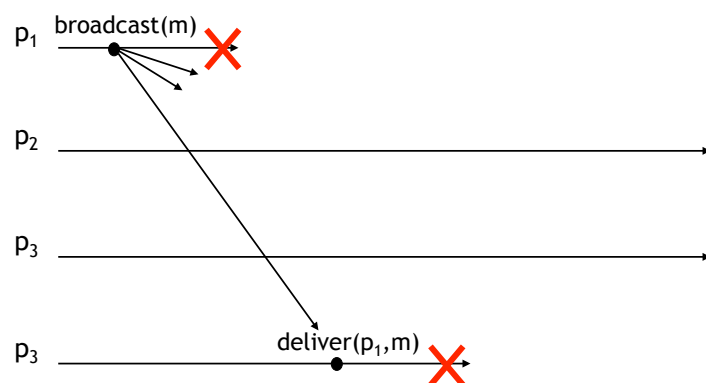
19/3/2012

Seif Haridi

13

RB Example

- Is this allowed? **Yes**



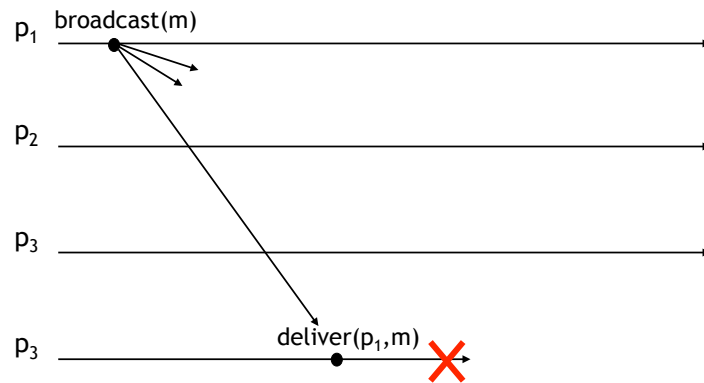
19/3/2012

Seif Haridi

14

RB Example

■ Is this allowed? **No**



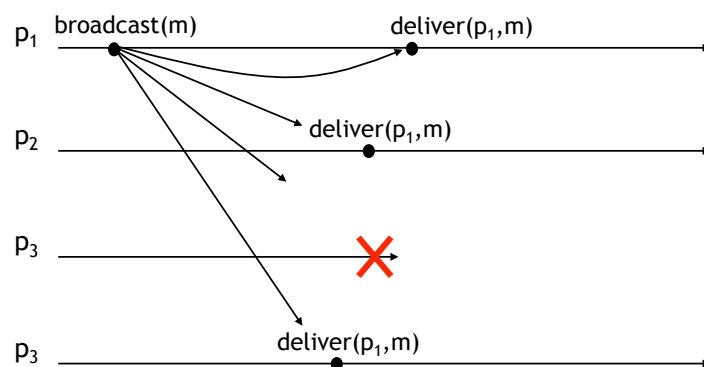
19/3/2012

Seif Haridi

15

RB Example

■ Is this allowed? **Yes**



19/3/2012

Seif Haridi

16

Uniform Reliable Broadcast

- Assume the broadcast enforces some real-world action
 - Printing a message on paper
 - Withdrawing money from account in variable
 - Launching a missile
- Assume sender broadcasts message
 - Sender fails
 - No correct node delivers message
 - *Failed nodes* might or might not deliver message, is it ok?
- **Uniform** reliable broadcast intuition
 - If a *failed node* delivers, *everyone* must deliver...
(At least the correct nodes, we cannot revive the dead...)

19/3/2012

Seif Haridi

17

Uniform Reliable Broadcast (urb)

- **Events**
 - Request: $\langle \text{urbBroadcast} \mid m \rangle$
 - Indication: $\langle \text{urbDeliver} \mid \text{src}, m \rangle$
- **Properties:**
 - **URB1**
 - **URB2**
 - **URB3**
 - **URB4**

19/3/2012

Seif Haridi

18

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 AND
Alive!

Implementation of Broadcast Abstractions

Implementing BEB

- Use Perfect channel abstraction
 - Upon `<bebBroadcast|m>` send message to all nodes (for-loop)
- Correctness
 - If sender doesn't crash, every other correct node receives message by perfect channels
 - No creation & no duplication already guaranteed by perfect channels

19/3/2012

Seif Haridi

21

Fail-Stop: Lazy Reliable Broadcast

- Requires perfect failure detector (P)
- To broadcast `m`:
 - **best-effort broadcast** `m`
 - When get `bebDeliver`
 - Save message, and
 - `rbDeliver` message
- 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]**
- Filter duplicate messages

19/3/2012

Seif Haridi

22

Lazy Reliable Broadcast

- **Implements:** ReliableBroadcast (rb)

- **Uses:**

- BestEffortBroadcast (beb)
- PerfectFailureDetector (P)

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

- $\text{delivered} := \emptyset$
- $\text{correct} := \Pi$
- **forall** $p_i \in \Pi$ **do** $\text{from}[p_i] := \emptyset$

for filtering
duplicates

storage for saved
messages

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

- **trigger** $\langle \text{bebBroadcast} \mid (\text{DATA}, \text{self}, m) \rangle$

19/3/2012

Seif Haridi

23

Lazy Reliable Broadcast (2)

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

- $\text{correct} := \text{correct} \setminus \{p_i\}$
- **forall** $(s_m, m) \in \text{from}[p_i]$ **do**
 $\text{trigger} \langle \text{bebBroadcast} \mid (\text{DATA}, s_m, m) \rangle$

Case 1: redistribute
anything we have
from failed node

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

- **if** $m \notin \text{delivered}$ **then**
- $\text{delivered} := \text{delivered} \cup \{m\}$
- $\text{from}[p_i] := \text{from}[p_i] \cup \{(s_m, m)\}$
- **trigger** $\langle \text{rbDeliver} \mid s_m, m \rangle$
- **if** $p_i \notin \text{correct}$ **then**
 $\text{trigger} \langle \text{bebBroadcast} \mid (\text{DATA}, s_m, m) \rangle$

Avoid duplicates

Store for future

Case 2: redistribute

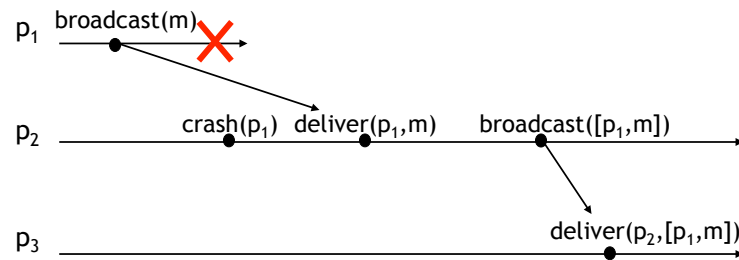
19/3/2012

Seif Haridi

24

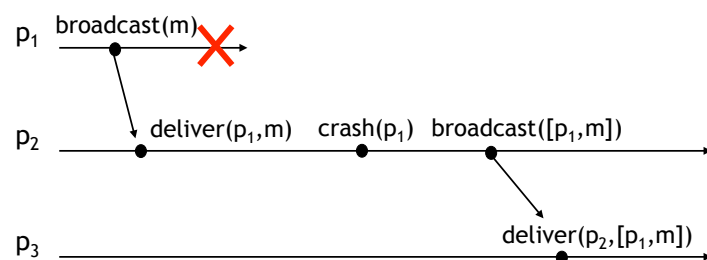
RB Example

- Which case? **Case 2**



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
- If correct p_j delivers msg broadcast by p_i
 - If p_i is correct, BEB ensures correct delivery
 - If p_i crashes,
 - p_j detects this (completeness)
 - p_j uses BEB to ensure (BEB1) every correct node gets it
 - This is a proof by induction. Why? **[d]**

19/3/2012

Seif Haridi

27

Eager Reliable Broadcast

- What happens if we replace P with $\Diamond P$? **[d]**
 - Only affects performance, not correctness.
- Can we modify Lazy RB to not use P ? **[d]**
 - Just assume all nodes failed
 - BEB Broadcast as soon as you get a msg

19/3/2012

Seif Haridi

28

Eager Reliable Broadcast

- **Uses:** BestEffortBroadcast (beb)
- **upon event** $\langle \text{Init} \rangle$ **do**
 - $\text{delivered} := \emptyset$
- **upon event** $\langle \text{rbBroadcast} \mid m \rangle$ **do**
 - $\text{delivered} := \text{delivered} \cup \{m\}$
 - **trigger** $\langle \text{rbDeliver} \mid \text{self}, m \rangle$ Immediately deliver
 - **trigger** $\langle \text{bebBroadcast} \mid (\text{DATA}, \text{self}, m) \rangle$ Immediately BEB broadcast
- **upon event** $\langle \text{bebDeliver} \mid p_i, (\text{DATA}, s_m, m) \rangle$ **do**
 - **if** $m \notin \text{delivered}$ **then**
 - $\text{delivered} := \text{delivered} \cup \{m\}$
 - **trigger** $\langle \text{rbDeliver} \mid s_m, m \rangle$ Immediately deliver
 - **trigger** $\langle \text{bebBroadcast} \mid (\text{DATA}, s_m, m) \rangle$ Immediately BEB broadcast

19/3/2012

Seif Haridi

29

Correctness of Eager RB

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

19/3/2012

Seif Haridi

30

Uniformity

- Is the proposed algorithm also **uniform**? [d]
- No.
 - Sender p immediately RB delivers and crashes
 - Only p delivered message
- Uniformity necessitates
 - If a **failed** node delivers a message m then **every correct** node delivers m

19/3/2012

Seif Haridi

31

Uniform Eager RB

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

Use vector **ack[m]** at p_i :
Set of nodes who
acked m

19/3/2012

Seif Haridi

32

Uniform Eager RB implementation

- **upon event** $\langle \text{urbBroadcast} \mid m \rangle$ **do**
 - $\text{pending} := \text{pending} \cup \{(\text{self}, m)\}$ ← remember sent messages
 - **trigger** $\langle \text{bebBroadcast} \mid (\text{DATA}, \text{self}, m) \rangle$

- **upon event** $\langle \text{bebDeliver} \mid p_i, (\text{DATA}, s_m, m) \rangle$ **do**
 - $\text{ack}[m] := \text{ack}[m] \cup \{p_i\}$ ← p_i obviously got m
 - **if** $(s_m, m) \notin \text{pending}$ **then** ← avoid resending
 - $\text{pending} := \text{pending} \cup (s_m, m)$
 - **trigger** $\langle \text{bebBroadcast} \mid (\text{DATA}, s_m, m) \rangle$

- **Upon exists** $(s_m, m) \in \text{pending}$ **s.t.** $\text{canDeliver}(m)$ **and** $m \notin \text{delivered}$ **do**
 - $\text{delivered} := \text{delivered} \cup \{m\}$
 - **trigger** $\langle \text{urbDeliver} \mid s_m, m \rangle$ ← deliver when all correct nodes have acked

19/3/2012

Seif Haridi

33

Correctness of Uniform RB

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

19/3/2012

Seif Haridi

34

Correctness of Uniform RB

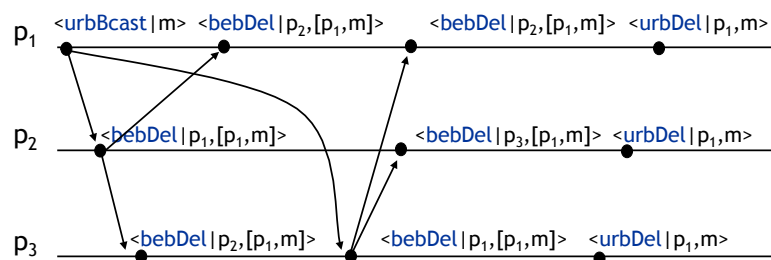
■ Validity

- If sender s is correct, it'll by **validity** (BEB1) $\text{bebDeliver } m$
- By the **lemma**, it will eventually $\text{urbDeliver}(m)$

■ Uniform agreement

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

URB Eager Algorithm Example



How difficult is URB?

- Strong failure detectors necessary for URB?
 - No, we'll provide URB 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

19/3/2012

Seif Haridi

37

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 node URB delivers, it got ack from majority
 - In that majority, one node, p , must be correct
 - p will ensure all correct nodes BEB deliver m
 - The correct nodes (majority) will ack and URB deliver

19/3/2012

Seif Haridi

38

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