Multicast Communication

Brian Nielsen bnielsen@cs.aau.dk

Communication modes in DS

- Uni-cast
 - Messages are sent from exactly <u>one</u> process to <u>one</u> process
- Broad-cast
 - Messages are sent from exactly <u>one</u> process to <u>all</u> processes on the network.
- Multi-cast
 - Messages are sent from exactly <u>one</u> process to <u>several</u> processes on the network (named group).
- Any-cast
 - Message is sent to <u>one</u> (eg "best" or "nearest") of a set of possible receivers
- Geo-cast:
 - Message sent to geographically close neighbors

 $g = \{p_1, p_2, p_3\}$

Multicast applications

- Shared whiteboards
- Chat applications
- Audio/Video conferencing
- Communication with server group
- Data replication
- Event notification
- Discovery of services
- Large scale interactive gaming on the Internet
- Distributed databases
- Large scale news distribution
- ... etc.



LAN Multicast

• Hardware support = 1 message is sent





WAN Multicast



Unicast to multiple receivers



Unicast

- With 4 receivers, sender must replicate the stream 4 times.
- Consider good quality audio/video streams are about 1.5Mb/s (a T1)
- Each additional receiver requires another 1.5Mb/s of capacity on the sender network
- Multiple duplicate streams over expensive WAN links

IP - Multicast



IP-Multicast Efficiency

- IP-multicast more **Efficient** than n sends!
 - Source transmits one stream of data for n receivers
 - Replication happens inside routers and switches
 - WAN links only need one copy of the data, not n copies.
- IP datagram multicast:
 - Hosts join/leave on a class D address
 - IGMP constructs and maintains multicast tree

IP-Multicast Failures

- HW- and IP-multicast Failure model ~ UDP
 - Omission failures
 - Delivery to none
 - Delivery to some
 - No ordering guarentees
 - Consequetive multicasts may be received in different order
- However, ordering and reliability are required by many applications
- Reliable & Ordered multicast requires "fancy"
 algorithms









FIFO-ORDERING



FIFO-ORDERING



TOTAL ORDERING



Multicast-API

- X-multicast(g,m)
- X-deliver(m)
- X is one of
 - B: Basic,
 - R: Reliable
 - FO: FIFO,
 - CO: Causal,
 - TO: Total
- Incoming messages (Receive)



The Hold-back queue



Basic Multicast

- A **basic multicast** primitive guarantees
 - All correct process eventually delivers the message, as long as the sender (multicasting process) does not crash
 - A "correct" process = a process that exhibits no failures at any execution point under consideration
 - NB: *NOT* satisfied by HW (IP) multicast
- A straightforward way to implement B-multicast is to use a reliable one-to-one send operation:
 - *B*-multicast(g,m): for each process p in g, send (p,m).
 - receive(m) at p: *B-deliver(m)*.

B-Multicast



•If P_n crashes, message not delivered in p_4 and p_5 •Hence, Unreliable

Reliable Uni-cast

- Integrity: A correct process p delivers a message m at most once. Furthermore, m is unmodified and was destined for p.
- Validity: If *m* was sent and the receiver is correct, it eventually delivers *m*.

Reliable multicast

- Integrity: A correct process p delivers a message m at most once. Furthermore, p∈group(m) and m was supplied to a multicast operation by sender(m).
- *Validity:* If a correct process multicasts message *m*, then it will eventually deliver m.
- Agreement: If a correct process delivers m, then all other correct processes in group(m) will eventually deliver m.
- *Liveness*=Validity+agreement

Reliable multicast

Algorithm 1 with B-multicast

On initialization

Received := $\{\};$

For process p to R-multicast message m to group g

B-multicast(g, m); // $p \in g$ is included as a destination

On B-deliver(m) at process q with g = group(m)

if $(m \notin Received)$ then

> Received := Received $\cup \{m\}$; if $(q \neq p)$ then B-multicast(g, m); end if R-deliver m;

end if

Each R-multicast message is sent |g| times, ie O(N²).

Reliable multicast

- Correct?
 - Integrity
 - Validity
 - Agreement
- Efficient?
 - NO: each message transmitted |g| times

- Each process maintains sequence numbers
 - $-S_{g}^{p}$ next message to be sent
 - R^q_g (for all q∈g) latest message delivered from q
- On *R-multicast* of *m* to group *g*, attach S_{g}^{p} and all pairs $\langle q, R_{g}^{q} \rangle$
- *R-deliver* in process *q* happens iff
 S^p_q=R^p_q+1
 - if $S^{p}_{g} < R^{p}_{g} + 1$, process *q* has seen the message before,
 - if $S^{p}_{g} > R^{p}_{g} + 1$ or if $R > R^{p}_{g}$ for some pair $\langle q, R \rangle$ in message a message has been lost



Data structures at process p:

 S_q^{p} : sending sequence number

 R_{g}^{q} : sequence number of the latest msg p delivered from q (for each q) On initialization:

 $S_g^p = 0$, $R_g^q = -1$, for all $q \in g$

For process p to R-multicast message m to group g

IP-multicast (g, <m, S_g^p , <**R**_g>>)

S_g^p ++

On IP-deliver (<m, S, <**R**>>) at q from p

(continued)

On IP-deliver (<m, S, <**R**>>) at q from p

save m

if S =
$$R_{g}^{p} + 1$$

then R-deliver (m)

R^p_g ++

check hold-back queue

else if $S > R_g^p + 1$

then store m in hold-back queue

request missing messages endif

endif

if $\exists p. r_g^p \in \mathbf{R}$ and $r_g^p > R_g^p$ then request missing messages endif

- 3 processes in group: P, Q, R
- State of process:
 - S: Next sequence number
 - R_a: Already delivered from Q
 - Set of Stored messages!
- Presentation:

P: 2 R: 5 Q: 3 <>

• Initial state:





• Arrival multicast by P at Q:



• New state:

• Multicast by Q:

Q: m_{q0}, 0, <P:0, R:-1>







- R detects missing message!
- When to delete stored messages?



$$\begin{array}{c|c} \mathbf{Q}: & 1 \\ P: 0 & R: -1 \\ < m_{p0}, , m_{q0} > \end{array}$$

- Correct?
 - Integrity:
 - seq numbers (duplicate detection) + checksums in IP multicast
 - Validity:
 - Self delivery assumed for IP
 - Agreement:
 - if missing messages are detected
 - \Rightarrow Correct processes multicasts indefinitely
 - if copy of message remains available
 - IMPROVE IT!

Ordered multicast

- **FIFO:** If a correct process issues $multicast(g, m_1)$, and then $multicast(g, m_2)$, then every correct process that delivers m_2 will deliver m_1 before m_2
- Causal: If multicast(g, m_1) \rightarrow multicast(g, m_2), where \rightarrow is the happed-before relation for messages send to g, then any correct process that delivers m_2 will deliver m_1 before m_2
- **Total:** If a correct process delivers message m_1 before it delivers m_2 then any other correct process that delivers m_2 will deliver m_1 before m_2

Total, FIFO and causal ordering

Notice the consistent ordering of totally ordered messages T_1 and T_2 , the FIFO-related messages F_1 and F_2 and the causally related messages C_1 and C_3 – and the otherwise arbitrary delivery ordering of messages.



FIFO multicast

- Analyse our algorithm for reliable multicast on top of IP-multicast.
- A process q delivers all messages from p in p sending order (S^p_g) by comparing to local expected sequence number R^p_g

(Unreliable) TO-multicast

- Basic approach as FIFO:
 - Uses globally unique IDs instead of per process unique IDs (as FIFO)
 - Receiver: deliver as for FIFO ordering
- Alg. 1: use a (single) sequencer process
- Alg. 2: participants collectively agree on the assignment of sequence numbers

TO-multicast: sequencer



 $\frac{On B-deliver(\langle m, i \rangle) with g = group(m)}{B-multicast(g, \langle \text{``order''}, i, s_g \rangle)};$

 $s_{\sigma} := s_{\sigma} + 1;$

(Unreliable) TO-multicast: ISIS

- Approach:
 - Sender:
 - B-multicasts message
 - Receivers:
 - Propose sequence numbers to sender
 - Sender:
 - uses returned sequence numbers to generate agreed sequence number

The ISIS algorithm for total ordering



The ISIS algorithm

- Process q maintains sequence numbers
 - $-A^{q}_{g}$ the largest agreed seq nr q has observed for g
 - $-Pq_{g}q$'s own largest proposed sequence number q
- Process *p* performs *B*-multicast($\langle m, i \rangle, g$), where *i* as a unique identifier for message *m*.
- Each process *q* replies *p* with a proposed sequence number P_q^q :=max(A_q^q , P_q^q)+1.
- Process *p* collects proposed sequence numbers and chooses the largest, let's call it *a*. Then *p* performs *B-multicast((i,a),g)*.
- Each process q in g sets A^q_g:=max(A^q_g,a) and attach sequence number a to message m

TO-multicast: ISIS alg.

- Correct?
 - Processes will agree on sequence number for a message
 - Sequence numbers are monotonically increasing
 - No process can prematurely deliver a message
- Performance
 - 3 serial messages!

CO-multicast

- Each process p_i maintains vector clock
 - Vⁱ_g[j] is the number of messages from each process
 Pⁱ_j that happened-before next message to be multicast
- To CO-multicast(m): P_i increments V_gⁱ [i] and B-multicasts(g,< V_gⁱ,m>)
- $P_i CO-delivers(m)$ from P_i iff
 - a) It has delivered any earlier message send by P_j $V_g^{j}[j] = V_g^{i}[j] + 1$, and
 - b) It has delived any message that P_j had delivered at the time it multicast the message: V_g^j[k] ≤ V_gⁱ[k] +1,k≠j
- **E.g.** message: V^{j} =[3,6,2] Receiver V^{i} =[2,5,2]

Summary

- So you thought multi-cast was simple??!!
- Applications have different semantic ordering, reliability and cost requirements
 - Unreliable / reliable multicast
 - FiFo, Causal, Causal-Fifo, Total, ...
 - FiFo+Total (Exercise)
- Many algorithms available with different cost / ordering tradeoff
- Did you see an algorithm for totally ordered reliable multicasting ????

