CS 268: Multicast Transport

Ion Stoica April 10-15, 2003

The Goal

Transport protocol for multicast Reliability

- Apps: file distribution, non-interactive streaming
- Low delay
- Apps: conferencing, distributed gaming
- Congestion control for multicast flows
 Critical for all applications

istoica@cs.berkeley.edu

2

Reliability: The Problems

- Assume reliability through retransmission
 Even with FEC, may still have to deal with
 retransmission (why?)
- Sender can not keep state about each receiver
 - E.g., what receivers have received, RTTNumber of receivers unknown and possibly very large
- Sender can not retransmit every lost packet
 - Even if only one receiver misses packet, sender must retransmit, lowering throughput
- N(ACK) implosion
 Described next

istoica@cs.berkeley.edu









Inter-node Latency Estimation **Repair Request Timer Randomization** - Chosen from the uniform distribution on Every node estimates latency to every other node $2^{i}[C_{1}d_{S,A}, (C_{1}+C_{2})d_{S,A}]$ (A)(B) - Uses session reports (< 5% of - A - node that lost the packet bandwidth) - S-source Assume symmetric latency $\begin{array}{l} C_{i}, C_{2} - \text{algorithm parameters} \\ - d_{S,A} - \text{latency between S and A} \\ - i - \text{iteration of repair request tries seen} \end{array}$ What happens when group d becomes very large? Algorithm Detect loss \rightarrow set timer - Receive request for same data \rightarrow cancel timer, set new timer, possibly with new iteration $d_{A,B} = (t_2 - t_1 - d)/2$ - Timer expires \rightarrow send repair request istoica@cs.berkeley.edu istoica@cs.berkeley.edu 7 8











Adaptive Timers

 ${\it C}$ and ${\it D}$ parameters depends on topology and congestion \rightarrow choose adaptively

After sending a request:

- Decrease start of request timer interval

- Before each new request timer is set:
- If requests sent in previous rounds, and any dup requests were from further away: · Decrease request timer interval
 - Else if average dup requests high:
 - · Increase request timer interval
 - Else if average dup requests low and average request delay too high: Decrease request timer interval

istoica@cs.berkeley.edu

13

15

Local Recovery

- Some groups are very large with low loss correlation between nodes Multicasting requests and repairs to entire group wastes
- bandwidth
- Separate recovery multicast groups
 - e.g. hash sequence number to multicast group address
- only nodes experiencing loss join group - recovery delay sensitive to join latency
- TTL-based scoping
 - send request/repair with a limited TTL
 - how to set TTL to get to a host that can retransmit
 - how to make sure retransmission reaches every host that heard request

istoica@cs.berkeley.edu

14

16

Application Layer Framing (ALF)

- [Clark and Tennenhouse 90]
- Application should define Application Data Unit
- (ADU) to lower layers
- ADU is unit of error recovery
 - app can recover from whole ADU loss
 - app treats partial ADU loss/corruption as whole loss
- App names ADUs
- App can process ADUs out of order
- Small ADUs (e.g., a packet): low delay, keep app busy
- Large ADUs (e.g., a file): more efficient use of bw and cycles
- Lower layers can minimize delay by passing ADUs to apps out of order

istoica@cs.berkeley.edu

Multicast Congestion Control Problem

- Unicast congestion control:
- send at rate not exceeding smallest fair share of all links along a path
- Multicast congestion control:
 - send at minimum of unicast fair shares across all receivers
 - · problem: what if receivers have very different bandwidths?
 - segregate receivers into multicast groups according to current available bandwidth

istoica@cs.berkeley.edu



















Resilient Multicast: STORM [Rex et al '97]

- Targeted applications: continuous-media applications
 - E.g., video and audio distribution
- Resilience
 - Receivers don't need 100% of dataPackets must arrive in time for repairs
 - Data is continuous, large volume
 - Old data is discarded

istoica@cs.berkeley.edu

istoica@cs.berkeley.edu

29

Design Implications

- Recovery must be fast
 SRM not appropriate (why?)
- Protocol overhead should be small
- No ACK collection or group management

istoica@cs.berkeley.edu

30

32

Solution	Details
 Build an application recovery structure Directed acyclic graph that span the set of receiver Does not include routers! Typically, a receiver has multiple parents Structure is built and maintained dsitributedly Properties Responsive to changing conditions Achieve faster recovery Reduced overhead 	 Use multicast (expanding ring search) to find parents When there is a gap in sequence number send a NACK Note: unlike SRM in which requests and repairs are multicast, with STORM NACKs and repairs are unicast Each node maintain List of parent nodes A quality estimator for each parent node A delay histogram for all packets received A list of timers for NACKs sent to the parent A list of NACKs note served yet Note: excepting the list of NACKs shared by parent-child all other info is local

31

istoica@cs.berkeley.edu





