Multicast technology: past, present, future

C. Pham

University Claude Bernard Lyon 1 LIP/INRIA RESO

Congduc.Pham@ens-lyon.fr www710.univ-lyon1.fr/~cpham

Sunday, February 6th 2005, DFMA 05

multicast! THE MULTICAST How multicast can change the way people use the Internet?



Purpose of this tutorial

- Multicast on the Internet has been introduced in 1986
- Provides a comprehensive overview of multicast technologies that span almost 20 years of communication networks
- Present the current technologies and concerns
- Show future directions in multicasting

Outline

 \square « the past »

- Only protocols
- One-size-fits-all philosophy

\Box « the present »

- Simpler communication models
- New concerns for congestion control and fairness

« the future »

- Interoperability with new network technologies
- New approaches, not necessarily based on IP multicast



This tutorial will...

- explain how multicast can change the way
 people use the Internet
- present the main technologies behind multicast, both at the routing and transport level
- show what are the main problems and how they can be solved
- state on the current deployment of multicast technologies and the problems encountered for large scale deployment

From unicast...

Sending same data to many receivers via unicast is inefficient Popular WWW sites become serious bottlenecks



...to multicast on the Internet.

- Not n-unicast from the sender perspective
- Efficient one to many data distribution
- Towards low latence, high bandwidth



New applications for the Internet Think about...

high-speed www video-conferencing video-on-demand interactive TV programs remote archival systems tele-medecine, white board high-performance computing, grids virtual reality, immersion systems



distributed interactive simulations/gaming...

A whole new world for multicast...



Introduction

The delivery models (1)

model 1: streaming (e.g. for audio/video)
 multimedia data requires efficiency due to its size
 requires real-time, semi-reliable delivery



The delivery models (2)

model 2: push delivery

- synchronous model where delivery is started at t0
- usually requires a fully reliable delivery, limited number of receivers
- Ex: synchronous updates of software



The delivery models (3)

□ model 3: *on-demand* delivery

- popular content (video clip, software, update, etc.) is continuously distributed in multicast
- users arrive at any time, download, and leave
- possibility of millions of users, no real-time constraint





A very simple example in figures File replication (PUSH) with ftp 10MBytes file I source, n receivers (replication sites) 512KBits/s upstream access □ n=100 $- T_x = 4.55$ hours □ n=1000 - $T_x = 1$ day 21 hours 30 mins!

A real example: LHC (DataGrid)



Reliable multicast: a big win for grids



Wide-area interactive simulations



The challenges of multicast

SCALABILITY - SECURITY - TCP Friendliness - MANAGEMENT

SCALABILITY





Multicast BONE at the ENS Lyon



Basics IP multicast

MBone tools - RAT

The Robust Audio Tool (RAT) is a an opensource audio conferencing and streaming application that allows users to participate in audio conferences over the internet. These can be between two participants directly, or between a group of participants on a common multicast group.



MBone tools - VIC

VIC is a video conferencing application developed by the Network Research Group at the LBNL in collaboration with the University of California, Berkeley.



MBone tools - WBD

WBD is a shared whiteboard

compatible with the LBL whiteboard, WB. It was originally written by Julian Highfield at Loughborough University and has since been modified by Kristian Hasler at UCL.



MBone - Advertising sessions

□ SDR is a session directory tool designed to allow the advertisement and joining of multicast conferences on the Mbone. It was originally modelled on sd written by Van Jacobson at LBNL.

| ₹ sd | r:rbennett @ | Prat.cs. | ucl.ac.u | |
|--|---------------------|----------|----------|------|
| New | Calendar | Prefs | Help | Quit |
| Public Sessions | | | | |
| 🐏 IMJ Channel 2 | | | | |
| 🐏 IP Multicast Summit 1999 – Busine | | | | |
| 👫 IP Multicast Summit 1999 – Deploy | | | | |
| 🐏 IP Multicast Summit 1999 – Keynot | | | | |
| 🐏 IP Multicast Summit 99 – Technolo 🚽 | | | | |
| ∠ JMRC | | | | |
| ? Lectures and Seminars | | | | |
| ? Low-Bandwidth Sessions | | | | |
| 🐏 Lund University: FilosoficirkeIn der | | | | |
| MECCANO Project meeting | | | | |
| (t) 2 | | | | |
| Private Sessions | | | | |
| | | | | |
| | | | | |
| Enter passphrase to view encounted sessions: | | | | |
| | | | | |
| Multicast Session Directory v2.7 | | | | |

A look back in history of multicast

History

- Long history of usage on shared medium networks
- Resource discovery: ARP, Bootp.



The Internet group model

multicast/group communications means...

 $\Box 1 \rightarrow n \quad \text{as well as} \quad n \rightarrow m$

- a group is identified by a class D IP address (224.0.0.0 to 239.255.255.255)
 - abstract notion that does not identify any host!



The group model is an open model

anybody can belong to a multicast group

- no authorization is required

a host can belong to many different groups

- no restriction

a source can send to a group, no matter whether it belongs to the group or not

- membership not required

the group is dynamic, a host can subscribe to or leave at any time

a host (source/receiver) does not know the number/identity of members of the group

Example: video-conferencing

The user's perspective





IP multicast TODO list

- Receivers must be able to subscribe to groups, need group management facilities
- A communication tree must be built from the source to the receivers
- Branching points in the tree must keep multicast state information
- Inter-domain routing must be reconsidered for multicast traffic
- Need to consider non-multicast clouds good luck...

unicast island









incremental deployment groups management session advertising tree construction address allocation duplication engine forwarding state routing

multicast island



Basics IP multicast



The two sides of IP multicast

Iocal-area multicast

- use the potential diffusion capabilities of the physical layer (e.g. Ethernet)
- efficient and straightforward

wide-area multicast

- requires to go through multicast routers, use IGMP/multicast routing/...(e.g. DVMRP, PIM-DM, PIM-SM, PIM-SSM, MSDP, MBGP, BGMP, MOSPF, etc.)
- routing in the same administrative domain is simple and efficient
- inter-domain routing is <u>complex</u>, not fully operational

IP Multicast Architecture



Part I « The past »



Basic of IP multicast model

Early group management

IP multicast routing

First steps in reliability

Early group mngt

Internet Group Management Protocol (RFC 1112)

- IGMP: "signaling" protocol to establish, maintain, remove groups on a subnet.
- Objective: keep router upto-date with group membership of entire LAN
 - Routers need not know who all the members are, only that members exist
- Each host keeps track of which mcast groups are subscribed to


IGMP: subscribe to a group (1)



224.0.0.1 reach all multicast host on the subnet Early group mngt



IGMP: subscribe to a group (3)



Data distribution example





| TCAAD | G 🖄 🖈 Sdr: Session Information 🗢 □ | | | | | | | |
|--|--|-----------------------------|----------|---------|-------------|---------------|--|--|
| LGWL | Encryption: none (NOENC) Places all over the world Authentication: none (NOAUTH) | | | | | | | |
| Join | Low bandwidth video (10-25 kb/s) with views from all over the world. | | | | | | | |
| | Contact Details | | | | | | | |
| | TTL: 127 Key: | | | | | | | |
| | Heard rom 128.253.115.224 at 21 Mar 2003 15:57 CET | | | | | | | |
| | Join | Invite | Record | | Dism | iss | | |
| ©≈≈ | | <c .pture=""> - Etherea</c> | | | | | | |
| <u>F</u> ile <u>E</u> dit <u>C</u> apture <u>D</u> isplay <u>T</u> ools | | | | | | <u>H</u> elp | | |
| No. Time Source | Destina | tion | Protocol | Info | | | | |
| 2520 284.239720 venezia.c | <u>ri2000.ens-1 224.</u> | 2.172.238 | | Host re | esponse (v2 |) Doctinat | | |
| 2522 284.271785 koralli.c | sc.fi 224. | 2.172.238 | UDP | Source | port: 1035 | Destinat | | |
| 2523 284.271902 stinky.dc | .luth.se 224. | 2.172.238 | UDP | Source | port: 1161 | Destinat | | |
| 2524 284.291741 IPTVserve | r.ldc.lu.se 224. | 2.172.238 | UDP | Source | port: 1062 | Destinat | | |
| 2525 284.296050 IPIVserve | r.ldc.lu.se 224. | 2.172.238 | UDP | Source | port: 1062 | Destinat | | |
| 2520 284.290525 IPTVSerVe | r ldc lu so 224. | 2.172.238 | | Source | port: 1062 | Destinat | | |
| 2528 284, 299740 IPTVserve | r.ldc.lu.se 224. | 2.172.238 | UDP | Source | port: 1062 | Destinat | | |
| 2529 284.300016 IPIVserve | r.ldc.lu.se 224. | 2.1/2.238 | UDP | Source | port: 1062 | Destinat | | |
| 2530 284.300283 IPTVserve | r.ldc.lu.se 224. | 2.172.238 | UDP | Source | port: 1062 | Destinat | | |
| 2531 284.300561 IPTVserve | r.ldc.lu.se 224. | 2.172.238 | UDP | Source | port: 1062 | Destinat | | |
| RI | | | | | | | | |
| ■ Frame 2520 (46 on wire, ■ Ethernet II ■ Internet Protocol □ Internet Group Manageme Version: 1 Tvpe: 6 (Host respons Unused: 0x00 Checksum: 0x5d0e Group address: 224.2. | 46 captured) nt Protocol e (v2)) 172.238 (224.2.17 | 72.238) | | | | | | |

IGMP: leave a group (1)



224.0.0.2 reach the multicast enabled router in the subnet Early group mngt

IGMP: leave a group (2)





IGMP: leave a group (4)





IGMP: leave a group (5)



| TCAAD | <u> ৩ % ४</u> | Sdr: Session Information | | | | | | | | |
|--|---|--------------------------|---------------------|----------|-------------------|---|-------------------------------|--|--|--|
| LGWL | Tencryption: none (NOENC) Places all over the world Authentication: none (NOAUTH) | | | | | | | | | |
| | Low bandwidth video (10-25 kb/s) with views from all over the world. | | | | | | | | | |
| Leave | | | | | | | _ | | | |
| | Sentact Details | | | | | | | | | |
| | Format: H.261 Proto: RTP Addr: 224.2.172.238 Port: 51482 TTL: 127 Kev: | | | | | | | | | |
| | Heard from 128.253.115.224 at 21 Mar 2003 15:57 CET | | | | | | | | | |
| | Join | | Invite | Rec | ord | Dism | niss | | | |
| | Distant Distantishing - Artic | | | | |) and each state of a state of state of state | Februari Contractul II. anti- | | | |
| | | < | car ture> - Etherea | | | | 0 | | | |
| <u>File Edit Capture Display Tools</u> | | | | | | | <u>H</u> e | | | |
| No. Time Source | | Destination | | Protocol | Info | | | | | |
| 757 19.217944 D-128-208 | -20-224.dhcp | 224.2.172. | 138 | UDP | Source p | ort: 4062 | Destinati | | | |
| 758 19.219191 venezia.c | ri2000.ens-1 | 224.2.172 | 238 | UDP | Source p | ort: 1060 | Destinati | | | |
| 759 19.220437 venezia.c | ri2000.ens-1 | ALL-ROUTER | S.MCAST.NET | IGMP | Leave gr | oup (v2) | | | | |
| 760 19.220531 venezia.c | <u>ri2000.ens-1</u> | ALL-ROUTER | S.MCAST.NET | IGMP | <u>Leave gr</u> | oup (v2) | | | | |
| 761 19.220835 switch-gi | ga.ens–1yon. | ALL-SYSTEM | IS.MCAST.NET | IGMP | Router q | uery | | | | |
| 762 19.221013 switch gi | ga.ens ly o n. | ALL SYSTEM | IS.MCAST.NET | IGMP | R o uter q | uery | | | | |
| 763 19.238200 stinky.dc | .luth.se | 224.2.172. | 238 | UDP | S o urce p | ort: 1161 | Destinati | | | |
| 764 19.265099 D-128-208 | –20–224.dhcp | 224.2.172. | 238 | UDP | Source p | ort: 4062 | Destinati | | | |
| 765 19.265401 D-128-208 | –20–224.dhcp | 224.2.172. | 238 | UDP | Source p | ort: 4062 | Destinati | | | |
| 766 19.265425 D-128-208 | -20-224.dhcp | 224.2.172. | 238 | UDP | Source p | ort: 4062 | Destinati | | | |
| 767 19.265582 leejh.kr. | apan.net | 224.2.172. | 238 | UDP | Source p | ort: 1253 | Destinati | | | |
| 768 19.270547 rgt-451-p | cO2.wmin.ac. | SAP.MCAST. | NET | SAP | Announce | ment (v1) | | | | |
| 4 | | | | | | | | | | |

■ Frame 760 (46 on wire, 46 captured)
 ■ Ethernet II
 ■ Internet Protocol
 □ Internet Group Management Protocol
 ∨ersion: 1
 Tvpe: 7 (Leave group (v2))
 Unused: 0x00
 Checksum: 0x5c0e
 Group address: 224.2.172.238 (224.2.172.238)

IGMP: leave a group (5)



OK, now I can express local interest, so what?



Does all paths lead to Roma?



Before going further...

Multicast on Ethernet LAN

- How can a end-host get link-layer (MAC) packets?
- Review of Ethernet filtering
 - By default, the Ethernet device listen on
 - its (Ethernet) MAC address fixed in a PROM
 - The broacast MAC address FF:FF:FF:FF:FF:FF
 - Other Ethernet addresses must be explicitly programmed into the driver
 - For multicast, one must listen at:
 - the Ethernet-equivalent of 224.0.0.1 (all multicast host in the LAN)
 - The Ethernet-equivalent address on which multicast sessions are advertised

Mapping of IP multicast address A MAC address is built from a mapping of IP multicast addr (Deering88)



Part I « The past »

Basic of IP multicast model

Early group management

IP multicast routing

First steps in reliability



IP multicast routing

- Find a tree (dedicated, shared) between the source(s) and the receivers
- Dense Mode
 - Assumes that there are many many receivers willing to get multicast traffic
 - Uses the « push » model: every body can receive

Sparse Mode

- Assumes that the number of receivers is small
- Uses the « pull » model: requires an explicite query from the receivers.

Dense mode protocols, DVMRP

The Ancestor: DVMRP (Distance Vector Multicast Routing)

Based on Reverse Path Forwarding (RPF)

A multicast router forwards packets received from a link which is on the shortest path to the source, and drops other packets



DVMRP... (cont')

□resulting multicast distribution tree



□different sources lead to diff. trees

 \Rightarrow improves load distribution on the links



Creates a spanning tree...

DVMRP... (cont')

add "flood and prune" algorithm to dynamically update the tree

step 1: flood the Internet (only limited by the packet's TTL)



DVMRP... (cont')

flooding/pruning is done periodically to update the tree

 required to discover new receivers and remove branches to receivers who left the session

limitations:

- creates signaling load (PRUNE message)
- periodically creates important traffic (flooding)
- all routers keep some state for all the multicast groups in use in the Internet

DVMRP deployment

 large scale deployment of DVMRP in the MBONE (multicast backbone) since 1992
 tunnels are set up to link "multicast islands" through unicast areas



Multicast tunnelling illustrated



The early MBone with tunnels



source K. Almeroth's paper. IEEE Networks Magazine, Vol.14(1)

Mixing tunnels and native multicast



source K. Almeroth's paper. IEEE Networks Magazine, Vol.14(1)

DVMRP on Linux: the mrouted daemon

kwad ~ vifs_with_neighbors = 1 [This host is a leaf] Virtual Interface Table Name Local-Address Vif M Thr Rate Flags eth0 193.253.175.161 subnet: 193.253.175.128/26 1 Û 1 Û leaf group host (time left): 239.2.11.73 193.253.175.135 (0:03:47)239.2.11.72 193.253.175.142 0:03:49239.2.11.71 193.253.175.134 0:03:49224.0.0.4 193.253.175.161 0:03:49224.0.0.2 193.253.175.161 (0:03:41) IGMP querier: 193.253.175.129 up 50:21:15 last heard 0:00:40 ago Nbr bitmaps: 0x0000000000000000 pkts/bytes in : 772010/38687700 pkts/bytes out: 0/0 eth1 193.253.175.249 subnet: 193.253.175.248/30 1 querier leaf 1 0 1 193.253.175.249 (group host (time left): 224.0.0.4 0:02:29)193.253.175.249 (0:02:31)224.0.0.2 IGMP querier: 193,253,175,249 (this system) Nbr bitmaps: 0x0000000000000000 pkts/bytes in : 0/0 pkts/bytes out: 7936/10780820 2 eth2 193,253,175,253 subnet: 193,253,175,252/30 1 querier leaf 1 0 group host (time left): 224.0.0.4 193.253.175.253 (0:02:33) 224.0.0.2 193.253.175.253 (0:02:28)IGMP querier: 193.253.175.253 (this system) Nbr bitmaps: 0x0000000000000000 -More--(62%)

DVMRP summary

□ it works but... this is far from perfect

- periodical flooding creates a heavy load on routers/links
- each multicast router must keep some forwarding state for each group
- tunneling quickly became anarchic
- this is a flat architecture (the same protocol is used everywhere)

conclusion: "dense mode protocols" like DVMRP are not scalable enough for WAN multicast routing

- dense mode assumes a dense distribution of receivers, wrong in practice!

DVMRP uses Source-based Trees





Moving to a Shared Tree





Shared vs. Source-Based Trees

Source-based trees

Shortest path trees - low delay, better load distribution

More state at routers (per-source state)

Efficient in dense-area multicast

Shared trees

 Higher delay (bounded by factor of 2), traffic concentration

Choice of core/RP affects efficiency

Per-group state at routers

Efficient for sparse-area multicast

Sparse mode protocols

□ The newcomers: PIM-SM/MSDP/MBGP

- **PIM-SM** : Protocol Independent Multicast Sparse Mode
- MSDP: Multicast Source Discovery Protocol
- MBGP: Multi-protocol Border Gateway Protocol
- □ domain ≈ site, or ISP network
 - similar to "autonomous systems" of unicast routing
- intra-domain mcast routing uses PIM-SM
- inter-domain mcast routing requires MBGP
- the discovery of sources in other domains requires MSDP

PIM-SM Protocol Overview

Basic protocol steps

- Shared trees are unidirectional
- Routers with local members Join toward Rendezvous Point (RP) to join shared tree
- Routers with local sources encapsulate data in Register messages to RP
- Routers with local members may initiate data-driven switch to source-specific shortest path trees
- □PIM v.2 Specification (RFC 2362)

PIM-SM: Build Shared Tree



Data Encapsulated in Register



RP Send Join to High Rate Source


Build Source-Specific Distribution Tree



IP multicast routing

PIM-SM... (cont')

moving to a per-source tree is efficient for bulk data transfer, but has a higher cost in case of multiple sources

one tree per source versus a single shared



PIM-SM on Internet routers

PIM-SM is implemented on all major Internet routers (CISCO, JUNIPER, Alcatel AVICI, PROCKET...)

A linux package exists, see
<u>http://netweb.usc.edu/pim/</u> (I haven't tried it yet)

Example: PIM-SM on VTHD



IP multicast routing

Source doc VTHD

Configuration on CISCO routers
Enabling PIM

Declaring the RP IP addr of the RP

Part I « The past »



Basic of IP multicast model Early group management Early IP multicast routing First step in reliability



Reliability Models

- Reliability => requires redundancy to recover from uncertain loss or other failure modes.
- Two types of redundancy:
 - Spatial redundancy: independent backup copies
 - Forward error correction (FEC) codes
 - Problem: requires huge overhead, since the FEC is also part of the packet(s) it cannot recover from erasure of all packets
 - Temporal redundancy: retransmit if packets lost/error
 - Lazy: trades off response time for reliability
 - Design of status reports and retransmission optimization important

Reliability

Temporal Redundancy Model



End-to-end reliability models

- Sender-reliable
 - Sender detects packet losses by gap in ACK sequence
 - Easy resource management
- Receiver-reliable
 - Receiver detect the packet losses and send NACK towards the source



Challenge: scalability (1)

□ many problems arise with 10,000 receivers...

Problem 1: scalable control traffic

- ACK every 2 packets (à la TCP)...oops, 10000ACKs / 2 pkt!
- NAK (negative ack) only if failure... oops, if pkt is lost close to the source, 10000 NAKs!



Challenge: scalability (2)

problem 2: scalable repairs/exposure

 receivers may receive several time the same packet



A piece of the solutions (1)

□ solutions to problem 1: scalable control traffic

solution 1: feedback suppression at the receivers

- each node picks a random backoff timer
- send the NAK at timeout if loss not corrected
- solution 2: proactive FEC (forward error correction)
 - send data plus additional FEC packets
 - any FEC packet can replace any lost data packet

solution 3: use a tree of intelligent routers/servers

- use a tree for ACK aggregation and/or NAK suppression
- PGM, ARM, DyRAM

Reliability

A piece of the solutions (2)

□ solutions to problem 2: scalable repairs

- solution 1: use TTL-scoped retransmissions
 - repair packets have limited scope

solution 2: use proactive/reactive FEC

- proactive: always send data + FEC
- reactive: in case of retransmission, send FEC

solution 3: use a tree of retransmission servers

- a receiver can be a retransmission server if he has the requested data

Scalable Reliable Multicast Floyd et al., 1995

- Receiver-reliable, NACK-based
- NACK local suppression
 - Delay before sending
 - Based on RTT estimation
 - Deterministic + Stochastic
- Every member may multicast NACK or retransmission
- Periodic session messages
 - Sequence number: detection of loss
 - Estimation of distance matrix among members

Reliability









































Deterministic Suppression



Simple TTL-scoped of repairs

use the TTL field of IP packets to limit the scope of the repair packet



Summary: reliability problems

- What is the problem of loss recovery?
 - feedback (ACK or NACK) implosion
 - ACK/NACK aggregation based on timers are approximative!
 - replies/repairs duplications
 - TTL-scoped retransmissions are approximative!
 - Heterogeneity of receivers (crying baby, congestion control)
 - difficult adaptability to dynamic membership changes

- Design goals
 - reduce the feedback
 traffic
 - reduce recovery latencies
 - improve recovery isolation



