# Distributed Call Signaling (DCS) & Dynamic Quality of Service (DQoS) Architectures

W. Marshall, K. K. Ramakrishnan, E. Miller, G. Russell,  B. Beser,
M. Mannette, K. Steinbrenner, D. Oran, Bill Guckel, J. Pickens,
P. Lalwaney, J. Fellows, D. Evans, K. Kelly, F. Andreasen

AT&T, CableLabs, 3Com, Cisco, Com21, General Instrument,
Lucent Cable, NetSpeak, Telcordia

Nov. 1999
IETF Presentation

1

# Agenda

Presentation

- ◆ Introduction of DCS and DQoS Architectural Framework
- ◆ Walk through a Basic Call Flow: highlight DCS enhancements


- ◆ Integration of Resource Management
- ◆ Call Authorization; example DQoS flow setup using RSVP & COPS


- ◆ Proxy-Proxy Communication: additional info' to be passed
- ◆ Privacy: motivation and suggested enhancements
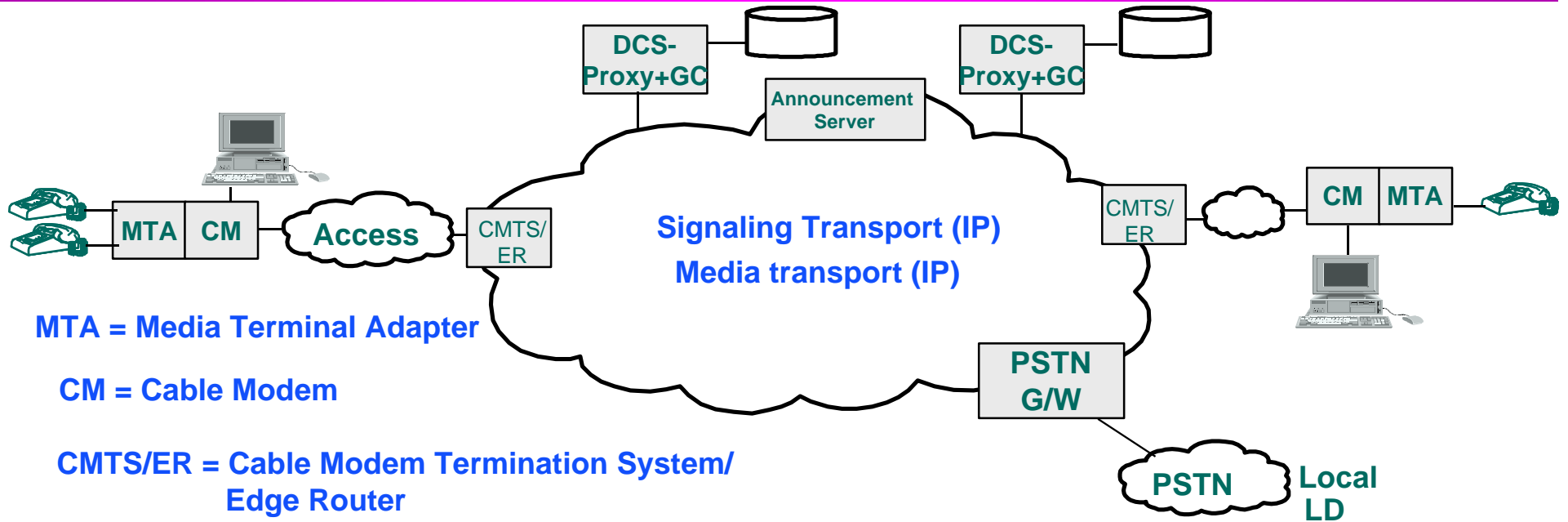- ◆ Communication of Call State Information to untrusted entities

Discussion

# IP Telephony: Opportunities

◆ Packet telephony & intelligent end terminals coupled with adequate bandwidth (especially access) provide tremendous opportunity
  – allows us to innovate in the communication services supported

◆ Take advantage of data and voice for an enhanced user experience
  – Browser-enabled telephones; click-to-dial

◆ Enhance our service provider role beyond basic telephony
  – Maintain and administer profiles for call handling, offer profile customization
    » e.g., handling calls from selected callers, call forwarding
  – Maintain personal directories, customized directories for small businesses
  – allow for customized handling of group calls, conferencing
    » manage group communication (e.g., chat) in customized manner

◆ IP dis-intermediates the service provider
  – how does the service provider play a role, *and derive revenue*?

# Distributed Open Signaling Architecture Framework



DCS-Proxy+GC

DCS-Proxy+GC

Announcement Server

MTA | CM — Access — CMTS/ER

Signaling Transport (IP)

Media transport (IP)

CMTS/ER — CM | MTA

PSTN G/W

PSTN — Local LD

**MTA = Media Terminal Adapter**

**CM = Cable Modem**

**CMTS/ER = Cable Modem Termination System/ Edge Router**

◆ Designed as a complete end-to-end signaling architecture for PacketCable

- Philosophy: encourage features and services in intelligent end-points, wherever technically and economically feasible

- "DCS Proxies" designed to be scalable transaction servers

- Resource management protocol provides necessary semantics for telephony

- "Gates" at network edge allow us to avoid theft of service

4

# Protocol Specification Efforts in PacketCable

◆ CableLabs: A group funded by multiple Cable Operators

   – Supports an activity to rapidly develop standards for Services over Cable

   – Major push now to standardize IP Telephony

◆ Initial effort was to provide support for

   – simple phones for consumer telephony: multiple lines per subscriber

     » based on SGCP/MGCP

       ⇨ limited to using the constrained user interface, provide traditional telephony

◆ Related efforts:

   – Separate Distributed Call Signaling Protocol

     » exploiting intelligence at end-points, address needs of provider

       ⇨ developing a SIP profile, with minimal extensions

   – Separate Dynamic Quality of Service Protocol

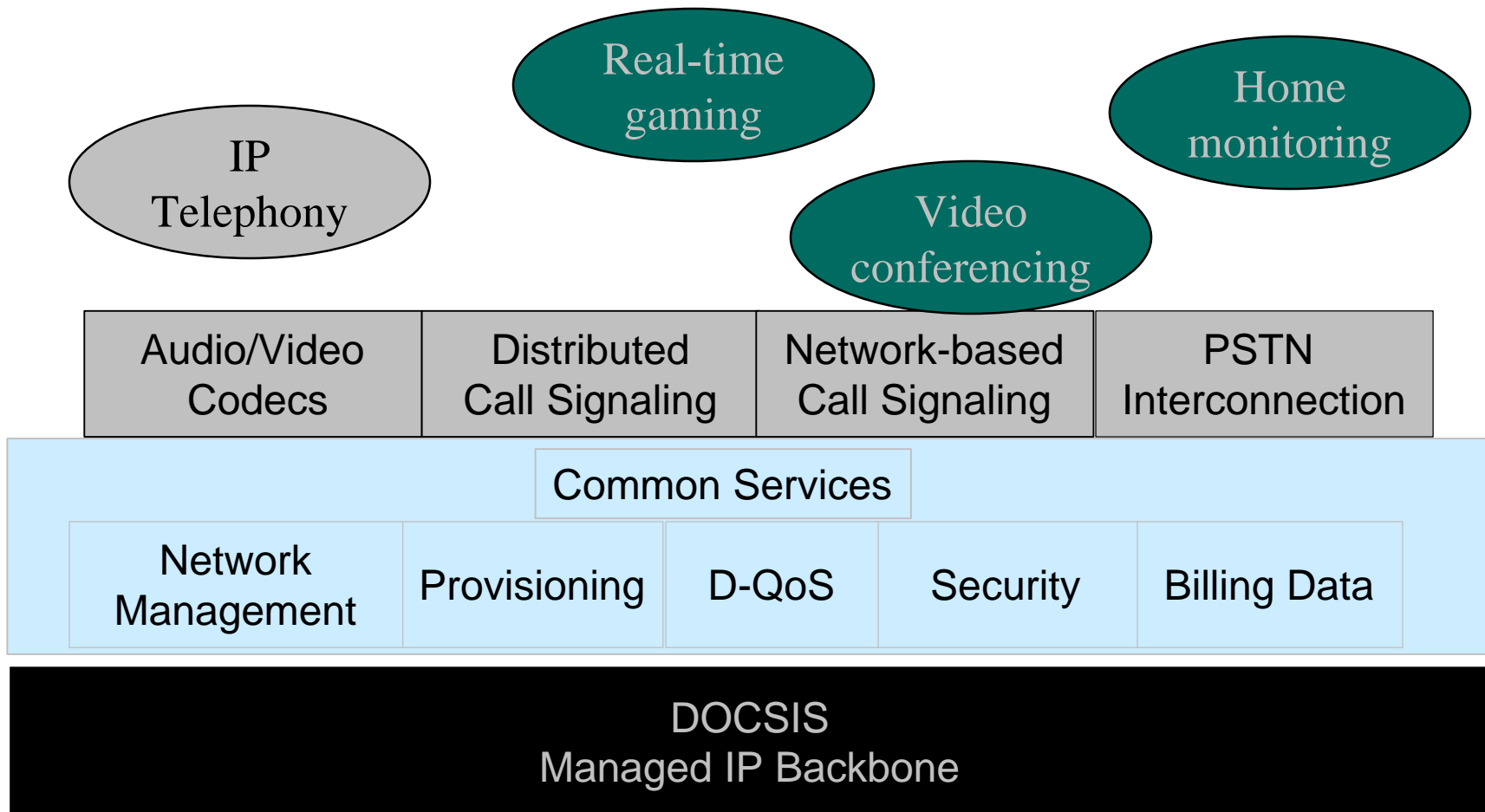     » based on supporting per-flow resource management on access network

# Extensions to SIP: Suggestions covered in the IDs

◆ We've found SIP to be a useful and flexible call signaling framework to incorporate the needs of a service provider

◆ Goal: Minimize local, cable environment specific solutions

- – use existing protocols, as far as possible

- – general applicability, beyond telephony, support interactive real-time streams

◆ We found the need to incorporate a small number of extensions to SIP and profile usage

- – hooks for resource management

- – support for privacy and anonymity

- – support for Local Number Portability

- – support for Billing, Operator services, law enforcement

- – support for communication of call state to end-points

◆ Hopefully, our solutions are applicable to other domains

  » e.g., wireless access, ...

6

# PacketCable ™ Layered Architecture

Real-time gaming

Home monitoring

IP Telephony

Video conferencing

| Audio/Video Codecs | Distributed Call Signaling | Network-based Call Signaling | PSTN Interconnection |
|---|---|---|---|

Common Services

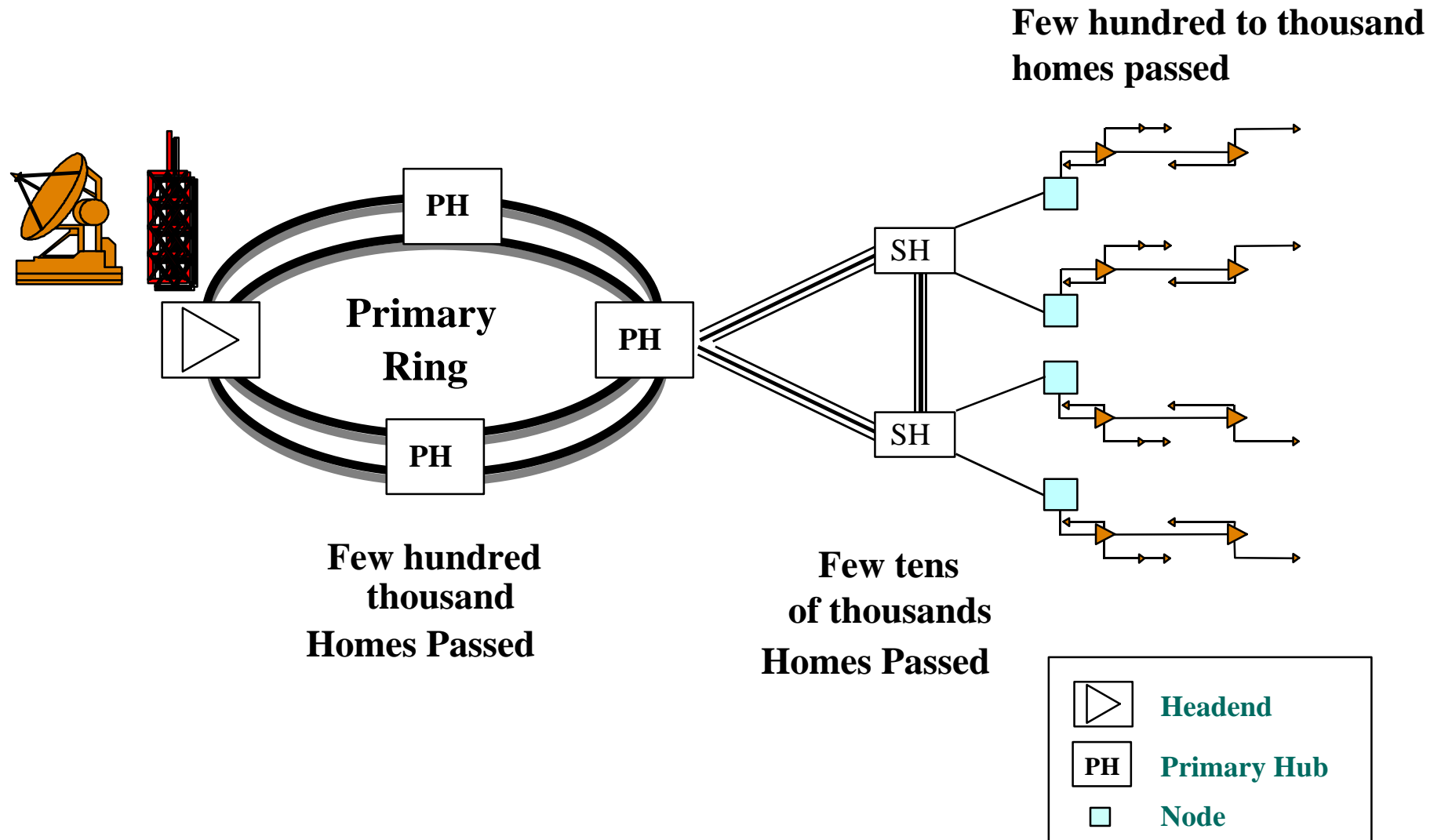| Network Management | Provisioning | D-QoS | Security | Billing Data |
|---|---|---|---|---|

DOCSIS
Managed IP Backbone

# HFC Media Capabilities

◆ Asymmetric Bandwidth Capabilities

◆ Downstream: multiple 27 Mbits/sec channels

– Operating in the range of 50-750 Mhz range

◆ Upstream: limited number of 2.0 - 2.5 Mbits/sec upstream channels

– Operating in the 5- 40 Mhz range

» number of channels typically limited by ingress noise

– typically 8-10 upstream channels supported on a given shared HFC cable

◆ Initial development of Medium Access Protocol: DOCSIS 1.0

– primarily supporting best-effort data service

» motivated by current Internet access

◆ Subsequent development of DOCSIS 1.1 enhancements

– enables support of a limited range of service classes

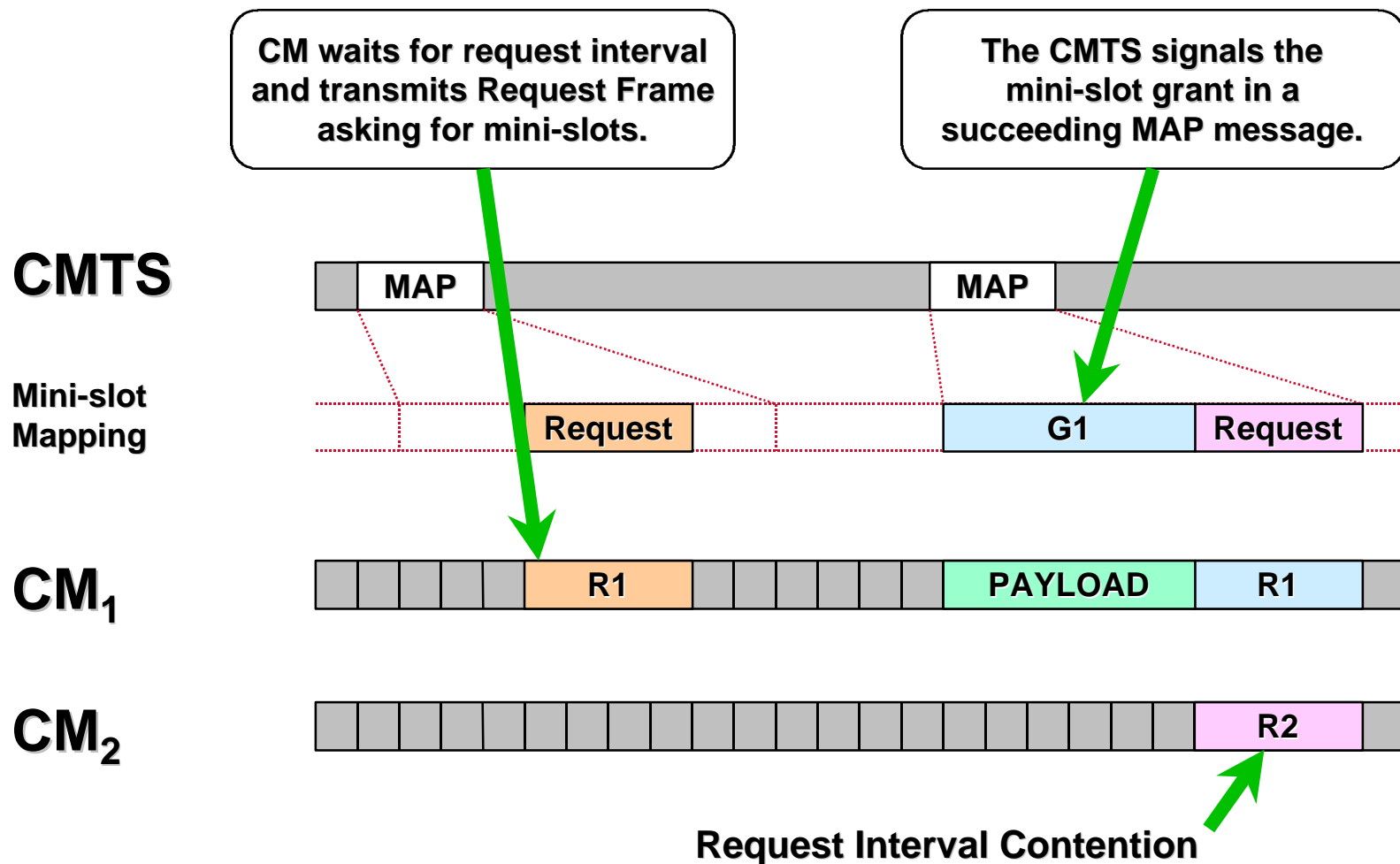» motivated, initially, to derive additional revenue from packet telephony

# Typical HFC Plant

**Few hundred to thousand homes passed**

PH

SH

**Primary Ring**

PH

PH

SH

**Few hundred thousand Homes Passed**

**Few tens of thousands Homes Passed**

| ▷ | **Headend** |
| --- | --- |
| PH | **Primary Hub** |
| ◻ | **Node** |

# DOCSIS 1.0 Media Access Protocol

◆ DOCSIS is a MAC protocol that is based on the CMTS scheduling

◆ DOCSIS 1.0 supports only best-effort service

- Cable Modem requests access to transmit a certain amount of data

  » requests using contention slots

- CMTS schedules modem transmission, and sends down a map on downstream channel

- modem may piggyback request with transmission of frame

- frames not fragmented

◆ Delay to transmit packet can be large and highly variable

# Request / Grant Mechanism

CMTS periodically sends a MAP downstream to CMs providing authorization to transmit on a given mini-slot
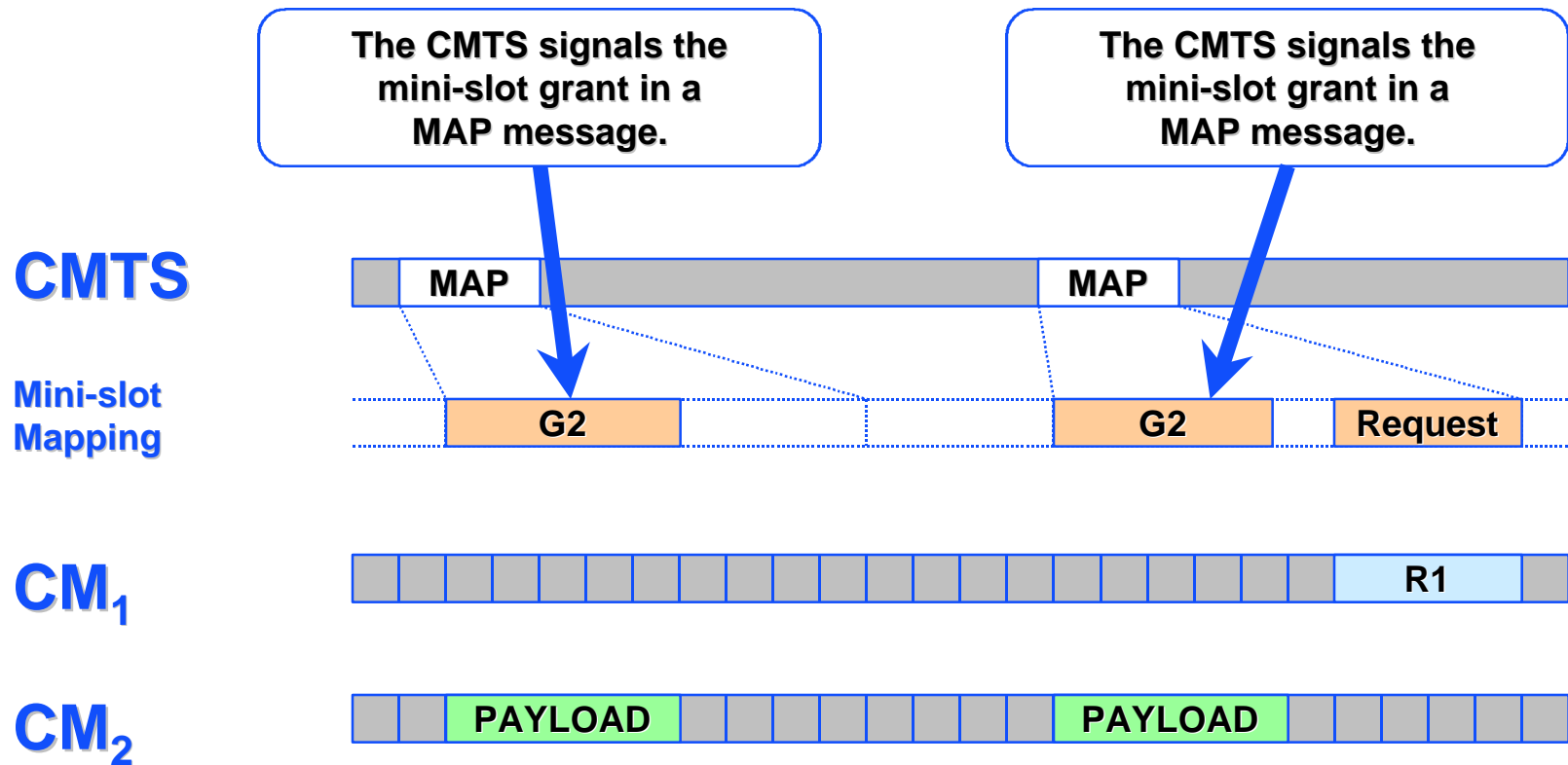
**CM waits for request interval and transmits Request Frame asking for mini-slots.**

**The CMTS signals the mini-slot grant in a succeeding MAP message.**

**CMTS**

| | MAP | | | MAP | |

**Mini-slot Mapping**

| Request | | G1 | Request |

**CM₁**

| | | | | | R1 | | | | | | PAYLOAD | R1 | |

**CM₂**

| | | | | | | | | | | | | | | | | R2 | |

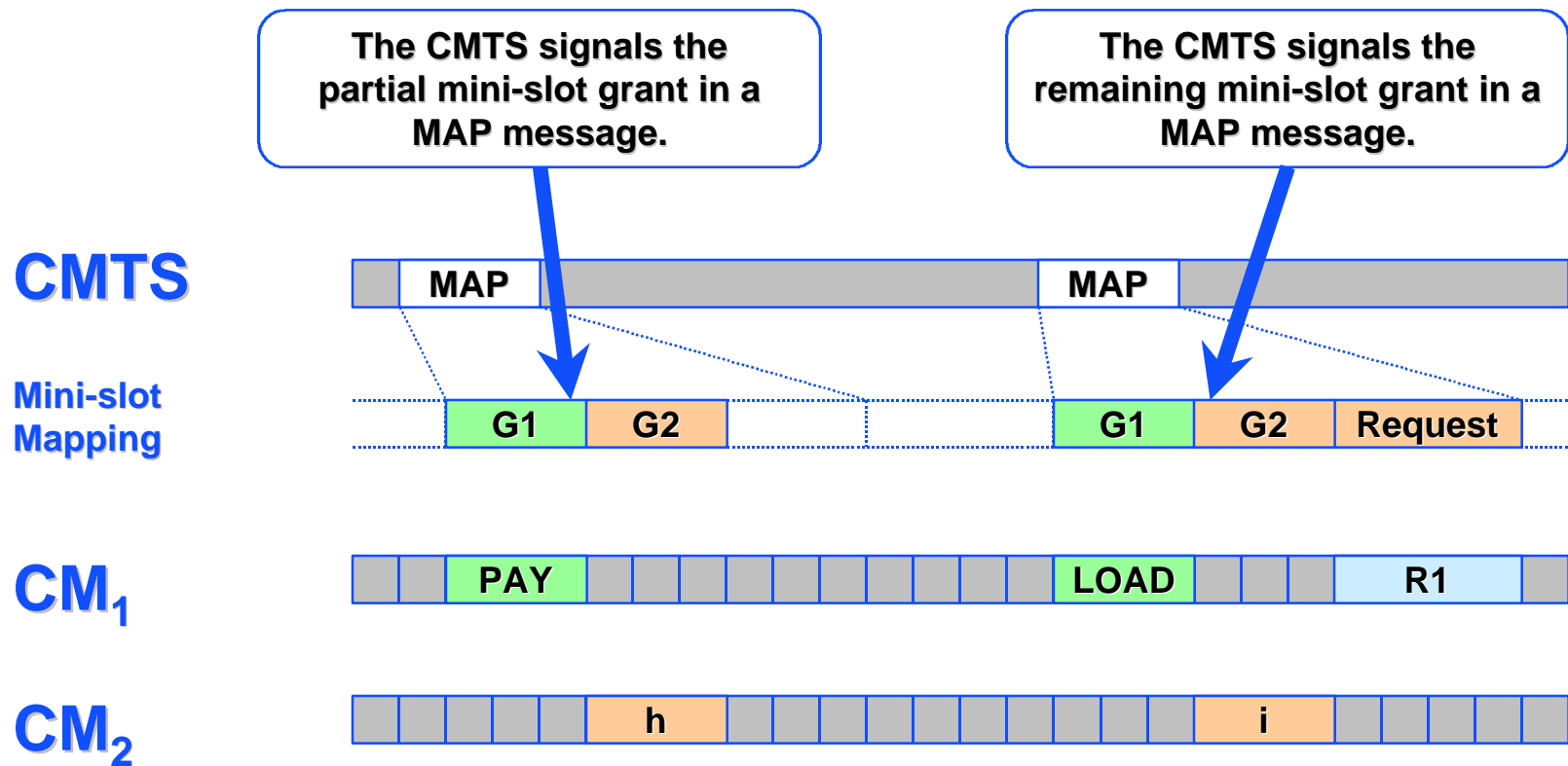**Request Interval Contention**

11

# DOCSIS 1.1 Enhancements

◆ Some Quality of Service support introduced in DOCSIS 1.1

- Support for Isochronous service and real-time polling

- Support for fragmentation of packets

◆ DOCSIS 1.1 introduces "unsolicited" grants

- CM/MTA negotiates with the CMTS to have a periodic grant of a certain size

  » Admission control performed on the request to provide isochronous service

- CMTS allocates grants to the CM/MTA periodically, with as little jitter as possible

  » vendor creativity and expertise helps in limiting jitter

◆ DOCSIS 1.1 introduces fragmentation

- CM negotiates with CMTS for maximum MTU size sustainable
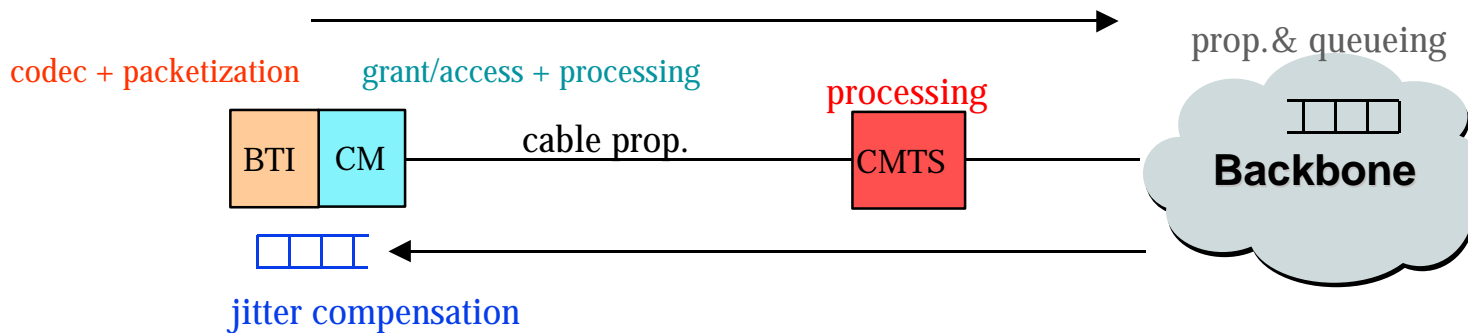
- CM fragments packets to be limited to that size

# Unsolicited Grant Mechanism

The CMTS signals the mini-slot grant in a MAP message.

The CMTS signals the mini-slot grant in a MAP message.

**CMTS**

| MAP | | MAP | |

**Mini-slot Mapping**

| G2 | | G2 | Request |

**CM$_1$**

| | | | | | | | | | | | | | | | | | | R1 | |

**CM$_2$**

| | PAYLOAD | | | | | | PAYLOAD | | | | | |

# Unsolicited Grants combined with Fragmentation



The CMTS signals the partial mini-slot grant in a MAP message.

The CMTS signals the remaining mini-slot grant in a MAP message.

**CMTS**

MAP      MAP

**Mini-slot Mapping**

G1   G2      G1   G2   Request

**CM$_1$**

PAY      LOAD      R1

**CM$_2$**

h      i

# Sources of Delay

codec + packetization    grant/access + processing         processing

prop.& queueing

cable prop.

| BTI | CM |    CMTS    **Backbone**

jitter compensation

◆ G.711: 80 Bytes at 64 Kbps => 10 ms packetization delay at BTI and gateway

◆ Sources of delay:

  – packetization, echo-canceller look-ahead & transmit processing delay in coder

  – processing delay in CM; Wait for Grant from CMTS

  – cable plant propagation delays; processing and queueing in CMTS

  – backbone propagation and queueing delay

◆ Jitter compensation

  – Worst case delay analysis needs to add the maximum jitter to delay PLUS the delay in the jitter compensation buffer at receiver
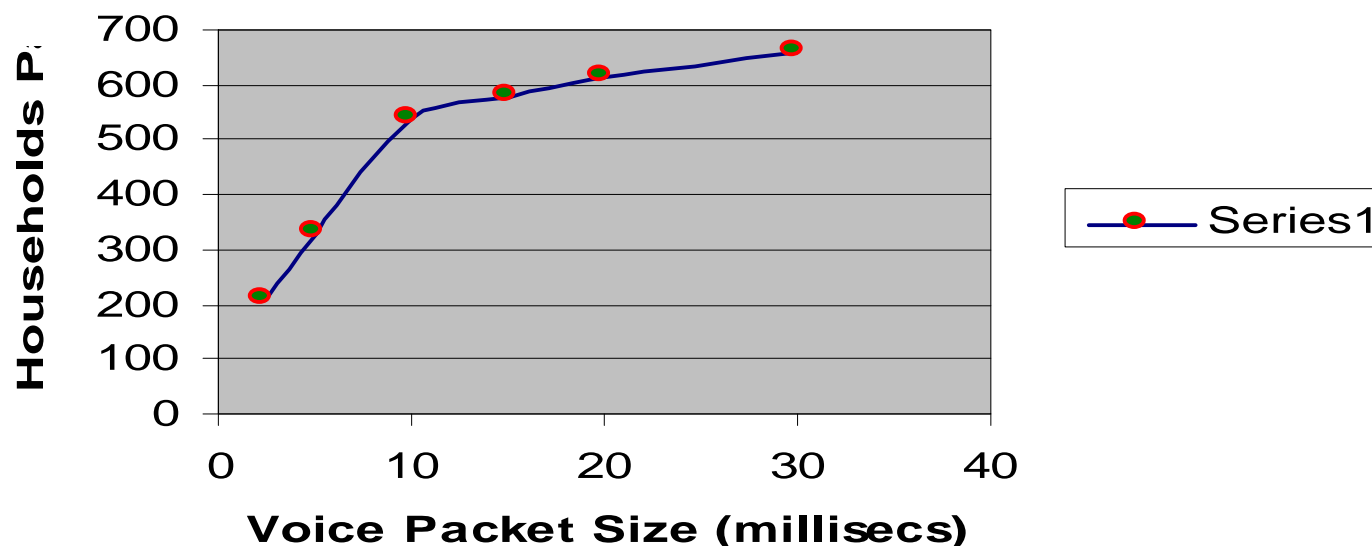
15

# Variation of Round-Trip Delay vs. Voice Packet Size

**Voice Pkts. Size vs. Delay**



◆ Typical (not restoration) backbone delays of 96 msec (4760 miles):

◆ 10 millsec. packetization results in: 199 msecs round trip

◆ 5 millisec. packetization results in: 172 msecs round trip

16

# Voice Packet Size vs. Capacity

**Voice Packet Size vs. Capacity (HHP)**



◆ Taking best case, 12 byte RTP, CRC, DOCSIS ovhd=11, PHS=2, but no UDP,IP,Enet headers #, homes passed with 25% take rate:

   – rapid reduction in capacity going below 10 millisecond voice packet size.

     » 208 HHP w/2.5 ms pkts; 328 HHP w/5 ms pkts; 536 HHP w/10 ms. pkts

   – Point of diminishing return beyond 10 msecs (header amortization insignificant)

# Summary of Delay and Capacity

◆ Round-Trip delay can be maintained below 300 milliseconds with

– keeping voice packet size at 20 milliseconds or less

– managing the backbone queueing delays

◆ Capacity is impacted by voice packet size

– payload header suppression helps

– choosing larger voice packet sizes doesn't result in substantial additional capacity

◆ We need to be concerned both with delay and capacity, and need to manage resources carefully.

# Requirements from a Service Provider's Perspective

◆ SIP will enable lots of these new services; but we also have to meet needs of current users.

◆ Need for differentiated quality-of-service is fundamental
  – must support resource reservation and admission control, where needed

◆ Allow for authentication and authorization on a call-by-call basis

◆ **Can't trust** CPE to transmit accurate information or keep it private

◆ Need to guarantee privacy and accuracy of feature information

  – e.g., Caller ID, Caller ID-block, Calling Name, Called Party

    » privacy may also imply keeping IP addresses private

◆ Protect the network from fraud and theft of service

  – critical, given the incentive to bypass network controls and billing

◆ We must be able to operate in large scale, cost-effectively

# Distributed Call Signaling Architecture

◆ Enhances SIP With Carrier Class Features

– Explicit recognition of need for Resource Management

– Privacy of "name", "number" and "address" of subscriber

– Don't retain Call state in network proxies

◆ Tight Coupling Between Call Signaling And QoS Control

– Authorize a call and allocate resources precisely when needed

» prevent Call Defects: don't ring the phone if resources are unavailable

– provide the ability to bill for usage, without trusting end-points

» prevent Theft Of Service: associate usage recording and resource allocation, ensuring non-repudiation

– ensure quality requirements for service are met (e.g., don't clip "Hello")

◆ Care taken to ensure untrusted end-points behave as desired
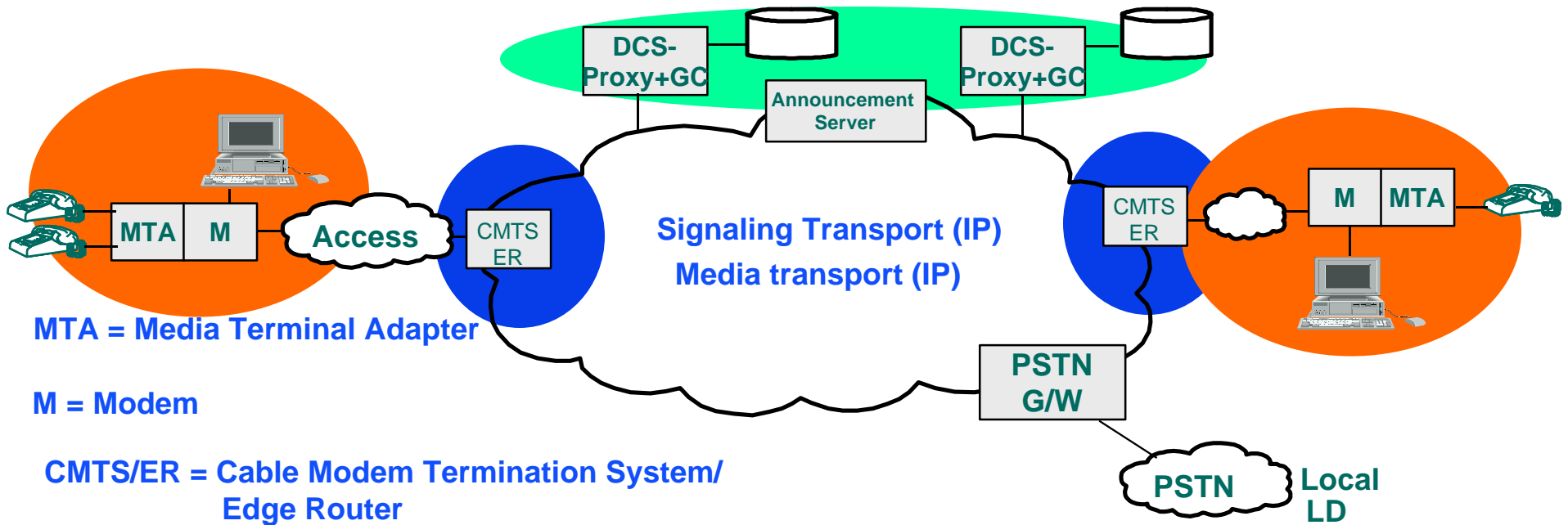
– Privacy mechanisms built into architecture

# Signaling Performance Requirements

◆ Short post-dial delay

- no perceptible difference in post-dial delay compared to circuit-switched network

◆ Short post-pickup delay

- delay from when the user picks up a ringing phone and the voice path being cut-through should be small

◆ Probability of Blocking: a metric to which provider may engineer net

◆ Probability of Call Defect (i.e., call that has both parties invited to and then fails due to lack of resources) needs to be much smaller

- target rates not necessarily under the control of the provider

◆ Flexibility in deployment of DCS-Proxies

- start small: possibly have a smaller number of proxies

- flexibility for the provider in placement of proxies

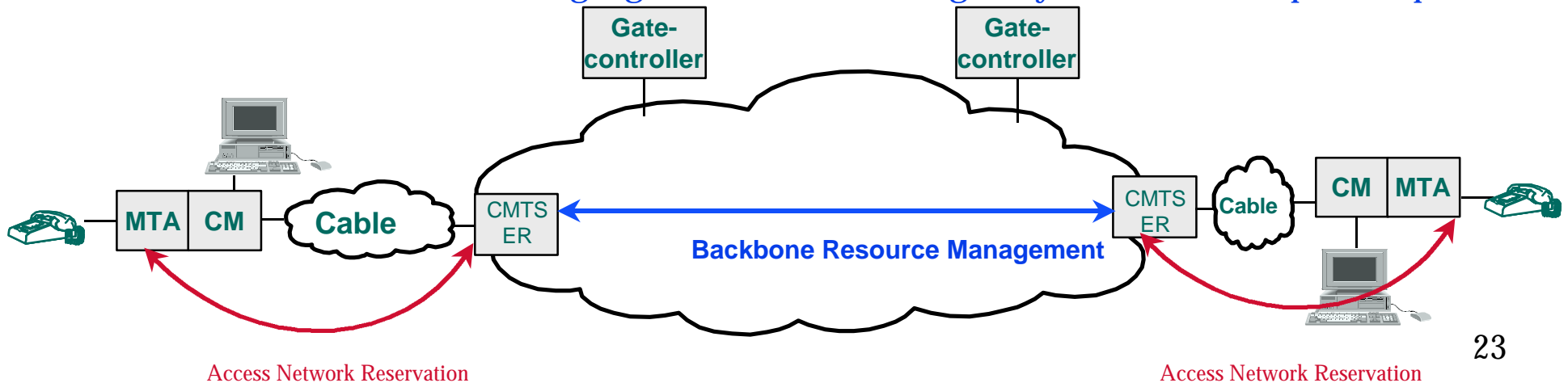⇨ reflected in the need for end-end signaling after "pick-up"

# DCS Architecture



**Signaling Transport (IP)**
**Media transport (IP)**

DCS-Proxy+GC

DCS-Proxy+GC

Announcement Server

MTA | M

Access

CMTS ER

CMTS ER

M | MTA

PSTN G/W

PSTN | Local LD

MTA = Media Terminal Adapter

M = Modem

CMTS/ER = Cable Modem Termination System/ Edge Router

**Call State**

**Connection State**

**Transaction State**

22

# Dynamic Quality of Service Framework

◆ Network considered to have multiple segments

 – each segment performs its own resource management, using protocols most appropriate for segment (e.g., RSVP/DOCSIS 1.1 on access; Diffserv on Backbone)

◆ Access network likely to be resource constrained

 – use per-flow signaling and allocation, on a call-by-call basis, in concert with call-leg manipulation

 – backbone network resources may be managed differently

◆ Two-phase resource management

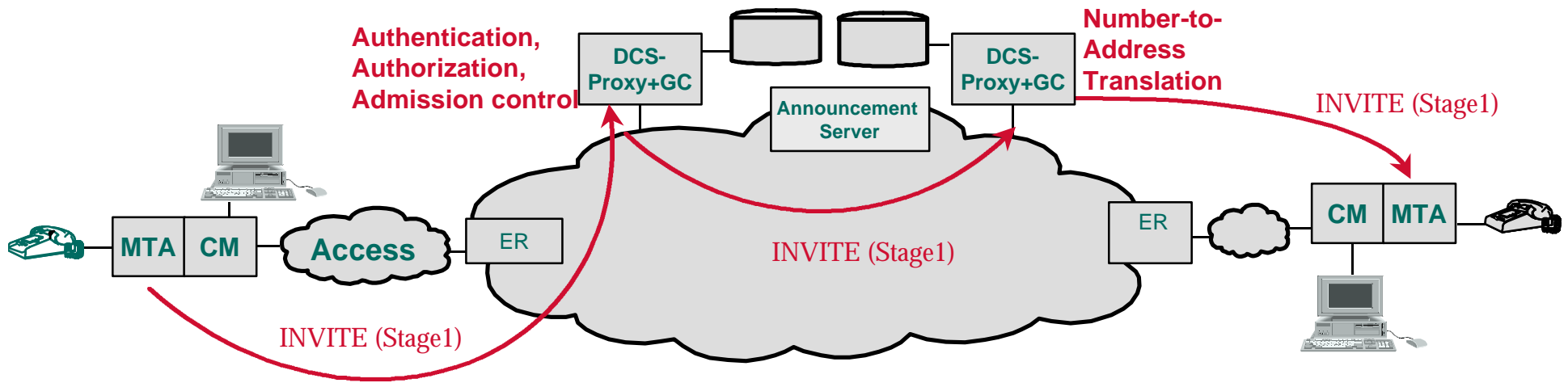 – have resource before ringing, but enable billing only after far end picks up.



Gate-controller

Gate-controller

MTA  CM

Cable

CMTS ER

Backbone Resource Management

CMTS ER

Cable

CM  MTA

Access Network Reservation

Access Network Reservation

23

# "Gates" and Edge Routers

◆ "Gates" in edge routers opened for individual calls

  – call admission control and policing implemented in edge routers

    » gate are packet filters that already exist in edge routers: "allow a call from this source to this destination"

      ⇨ for a particular range of traffic parameters, and a particular duration, etc.

  – however, *policy* is controlled by the gate controller

◆ Gate controller manipulates a gate after Call Setup is authorized

  – setting up gate *in advance* of reservation request allows a GC to be stateless

◆ MTA makes a resource reservation request by signaling to edge router

  – edge router admits the reservation if consistent with gate parameters

  – edge router generates usage recording events based on reservation state

◆ Accounting info stored at the edge router to generate usage events

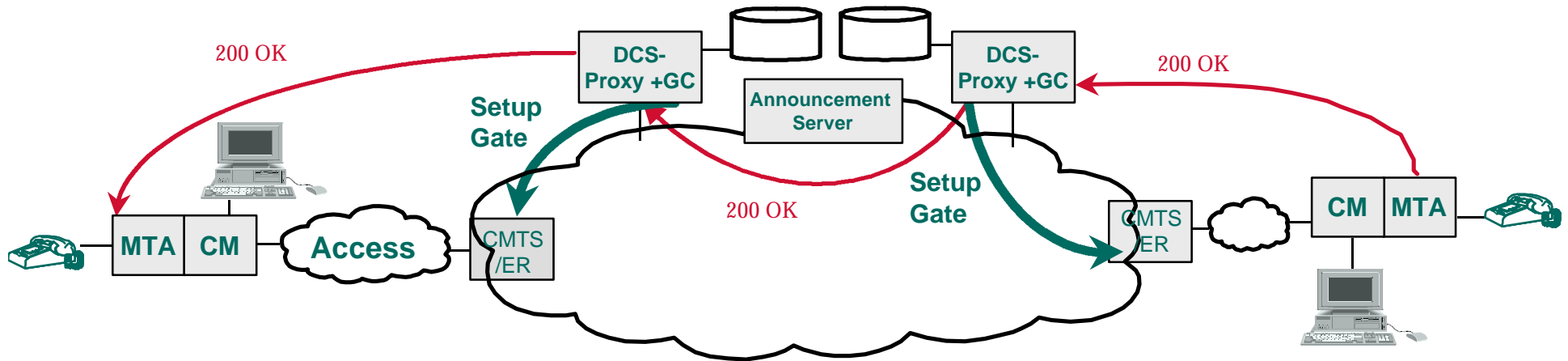    » opaque info' sent to record keeping servers for tracking usage and billing
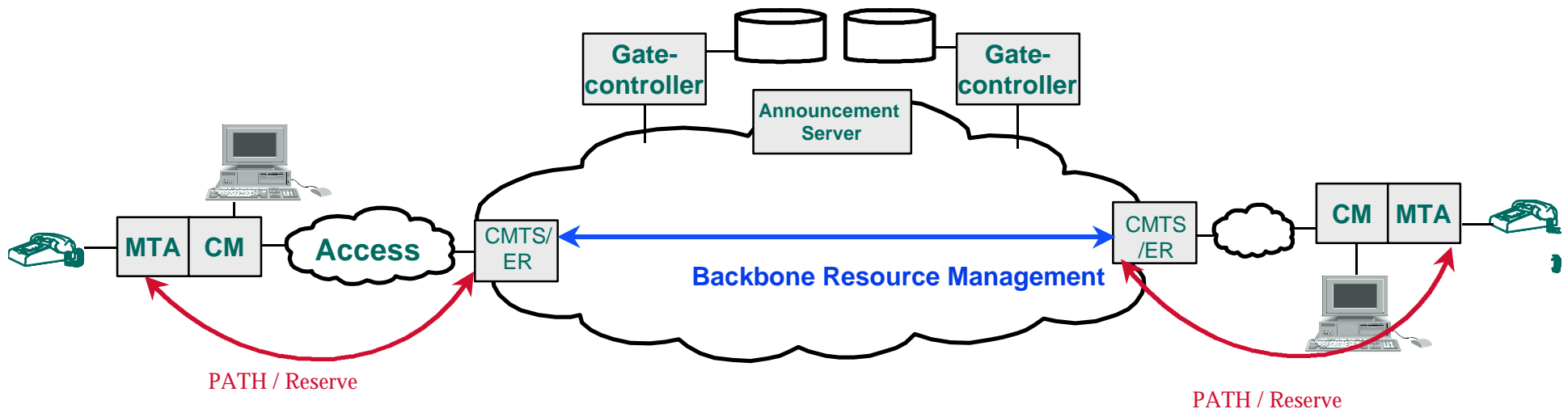
# Example Call Flow



Labels in diagram:
- **Authentication, Authorization, Admission control**
- **DCS-Proxy+GC**
- **Announcement Server**
- **DCS-Proxy+GC**
- **Number-to-Address Translation**
- INVITE (Stage1)
- MTA | CM
- **Access**
- ER
- INVITE (Stage1)
- INVITE (Stage1)
- ER
- CM | MTA
- INVITE (Stage1)

◆ MTA issues an INVITE to destination E.164 (or other) address

◆ Originating DCS-proxy performs authentication and authorization

◆ Terminating DCS-proxy translates dest. number to local IP address

   – no resources allocated: don't know yet "what" resources needed to "where"

   – provider may choose to block a call if resources are unavailable

      » but $P$(blocking) may be $\geq P$(call defect)

        ⇨ call defect: when the call fails after the parties are notified

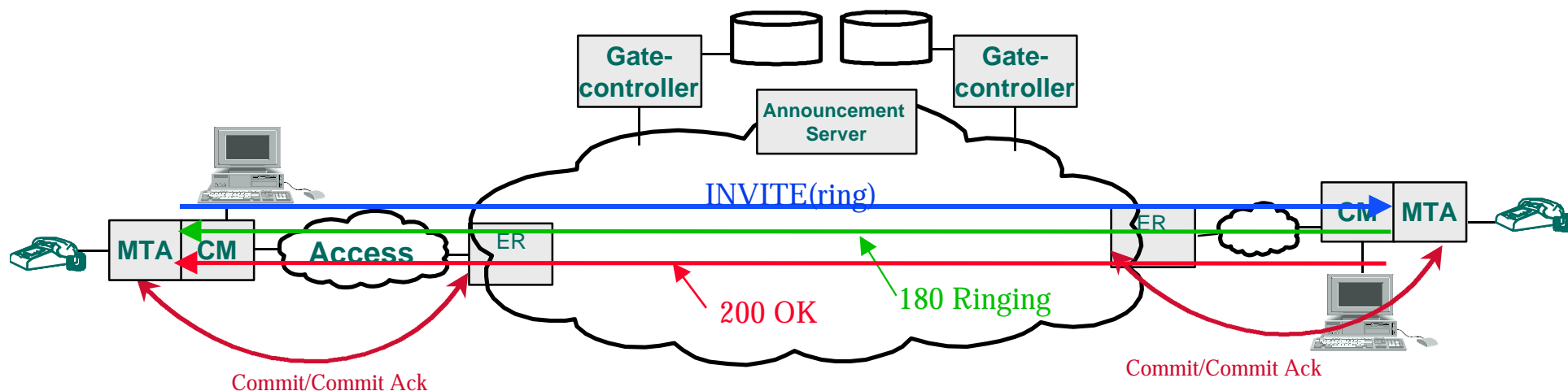# Example Call Flow (contd…)



- ◆ 200 OK conveys call parameters and (gate id) to originating MTA
- ◆ Gate controllers setup "gates" at edge routers as part of call setup
  - gate is described as an "envelope" of possible reservations issued by MTA
  - gate permits reservation for this call to be admitted
- ◆ Gate Controller acts as policy server in COPS framework
  - policy decisions provided to CMTS based on call signaling
  - CMTS acts as policy enforcement point

26

# Resource Management: 1<sup>st</sup> Phase



Gate-controller

Gate-controller

Announcement Server

MTA | CM

Access

CMTS/ER

CMTS/ER

CM | MTA

**Backbone Resource Management**

PATH / Reserve

PATH / Reserve

◆ MTA initiates resource reservation

– access resources are "reserved" after an admission control check

» this insures that resources are available when terminating MTA rings

– backbone resources are "reserved" (e.g., explicit reservation or "packet marking")

◆ Originating MTA starts end-to-end handshake with terminating MTA

– originating MTA sends INVITE(ring), terminating MTA sends 180 RINGING, 200 OK
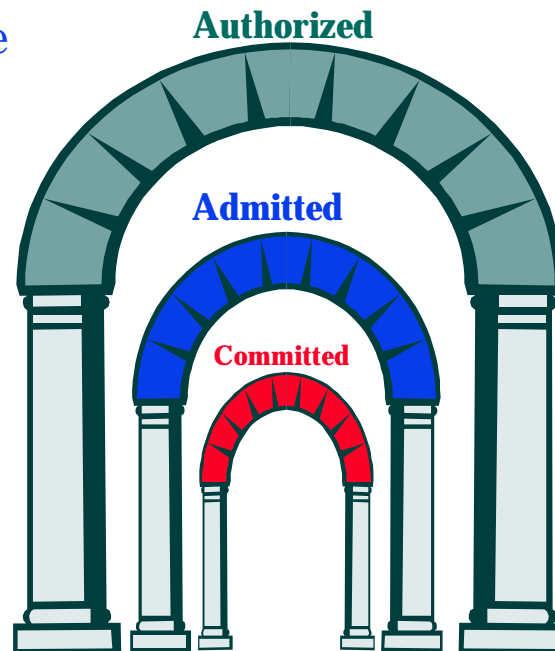
27

# Resource Management: 2nd Phase



- ◆ MTA knows voice path is established when it receives a 200 OK

- ◆ MTAs initiate resource "commitment"

  - – resources "committed" over access channel

    - » CMTS starts sending unsolicited grants; usage recording is started

  - – commitment deferred until far end pick up, to prevent theft of service; allow efficient use of constrained resources in access network

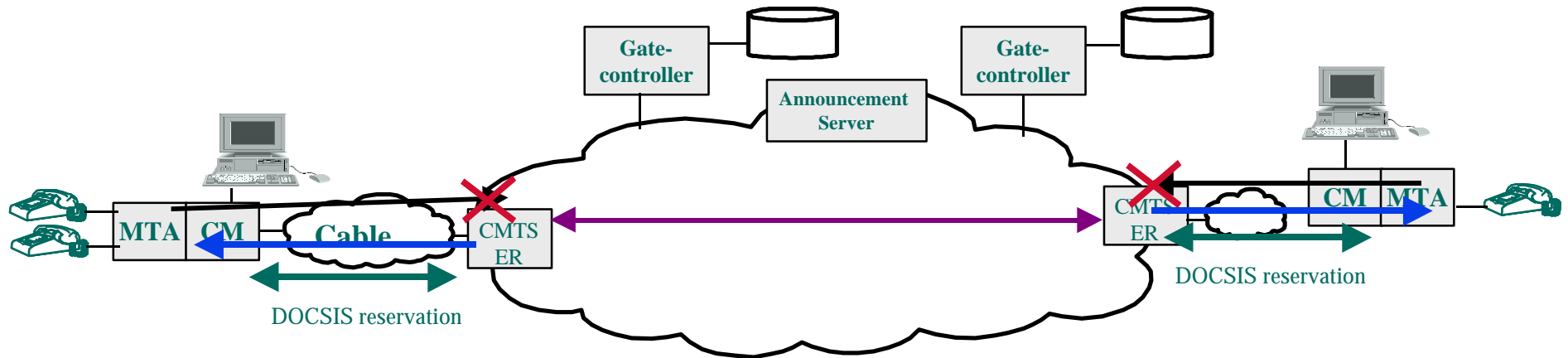- ◆ Commit opens the "gate" for this flow

# Resource Envelopes

◆ Authorization of resource usage done at call setup

   – Exercise Policy at GateController

   – "Authorized" Envelope.

◆ Later, capability negotiation enables end-points to reserve resources

   – "Admitted" Envelope

◆ End-points Commit to use resources when far end picks-up

   – "Committed" Envelope



Authorized

Admitted

Committed

# Using RSVP for Segmented Resource Allocation



**Client sends PATH message directed towards far endpoint**

**CMTS intercepts PATH message**

**CMTS reserves bandwidth on DOCSIS link**

**Backbone Resource Management ensures capacity available for call**

**On successful reservation, CMTS sends RESV message back to client**

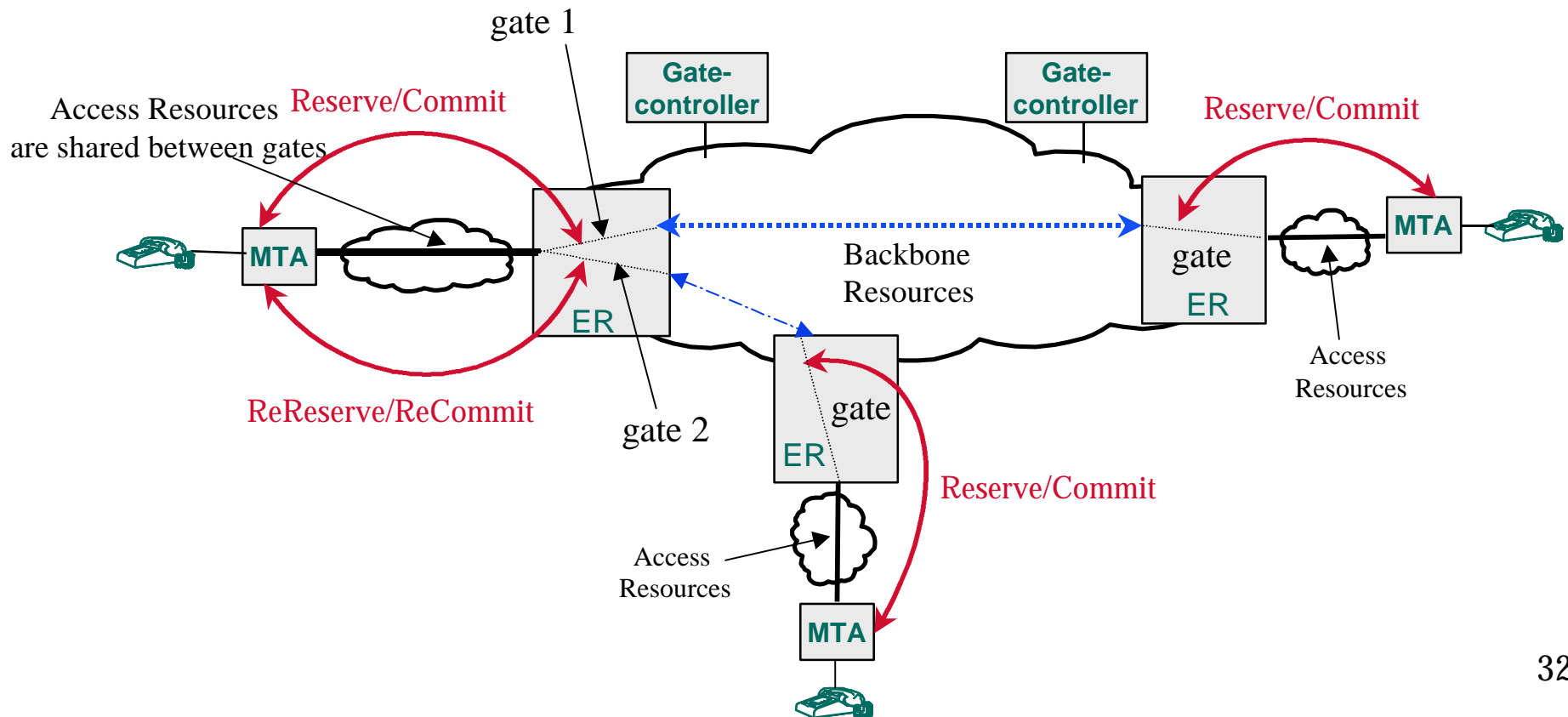minimal changes to RSVP: destination of PATH is destination of data

30

# Enhancements to the basic RSVP framework

◆ Bi-directional bandwidth is reserved with one PATH/RESV pair

  – Opaque objects added to RSVP messages to convey information for bi-directional bandwidth

◆ Resources identified explicitly

  – Resource_ID object

  – enables us to dynamically bind resources to a call

◆ Managing changes in resource usage during a call

  – e.g., Mid-call codec change; Call Waiting

    » Multiple flowspecs may be carried in one message

◆ Two-phase resource management support

  – Commit and Commit_Ack are unicast messages from client to CMTS

◆ Support included for resource reservation within the customer LAN

  – with traditional RSVP

31

# Call Waiting: Sharing Reservations Across Gates

◆ Access network resources are likely to be limited

- share resources between calls

- manage these resources carefully, in concert with call state

- de-commit or reuse resources when a call leg is on hold

# Summary

◆ DOSA introduced the concept of integrating QoS with call signaling

◆ DCS call signaling allows use of end-point intelligence to support new services and integration with other applications

◆ DCS proxies not required to be involved throughout call

  – simple transaction processor; less stringent reliability requirements;scalable

◆ Dynamic QoS provides the common underlying framework of QoS for call signaling protocols

◆ Two phase Reserve/Commit for managing resources

  – provides semantics that resources are available when phone rings, without billing for ringing

◆ Gates for each call: allows provider to manage access to resources

  – ensures that users who want toll quality go through network proxies

  – avoid theft of service with careful coordination between signaling and QoS