# How IP Telephony Breaks the Internet

Assumptions

Henning Schulzrinne

Dept. of Computer Science

Columbia University

New York, New York

(sip:)schulzrinne@cs.columbia.edu

Sprint IP Retreat

May 20–21, 2001

# Overview

- VoIP = $\boxed{\text{traffic/QoS}}$, signaling, services

- reliability issues

- breaking the Internet architecture

# VoIP

- carrying voice (and multimedia) over IP

- strict separation signaling – media traffic ($\leftrightarrow$ PSTN)

- future: high-rate codecs, video

- (typically) *not* PC-based voice

- starting to displace traditional PBX in greenfield installations

- likely to see widespread use in 3G (UMTS R5) wireless

# Example: Pingtel SIP phone
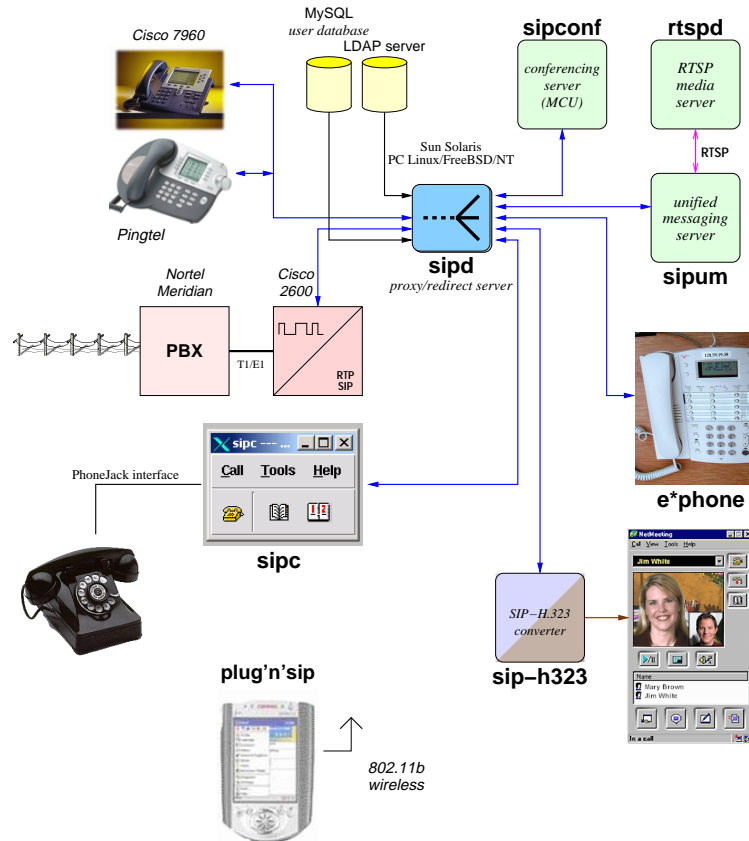
# Example: Cisco and 3Com SIP phones



Cisco

3Com

# Example: Columbia CS Phone System

Expand existing PBX via IP phones, with transparent connectivity

# The phone works — why bother with VoIP?

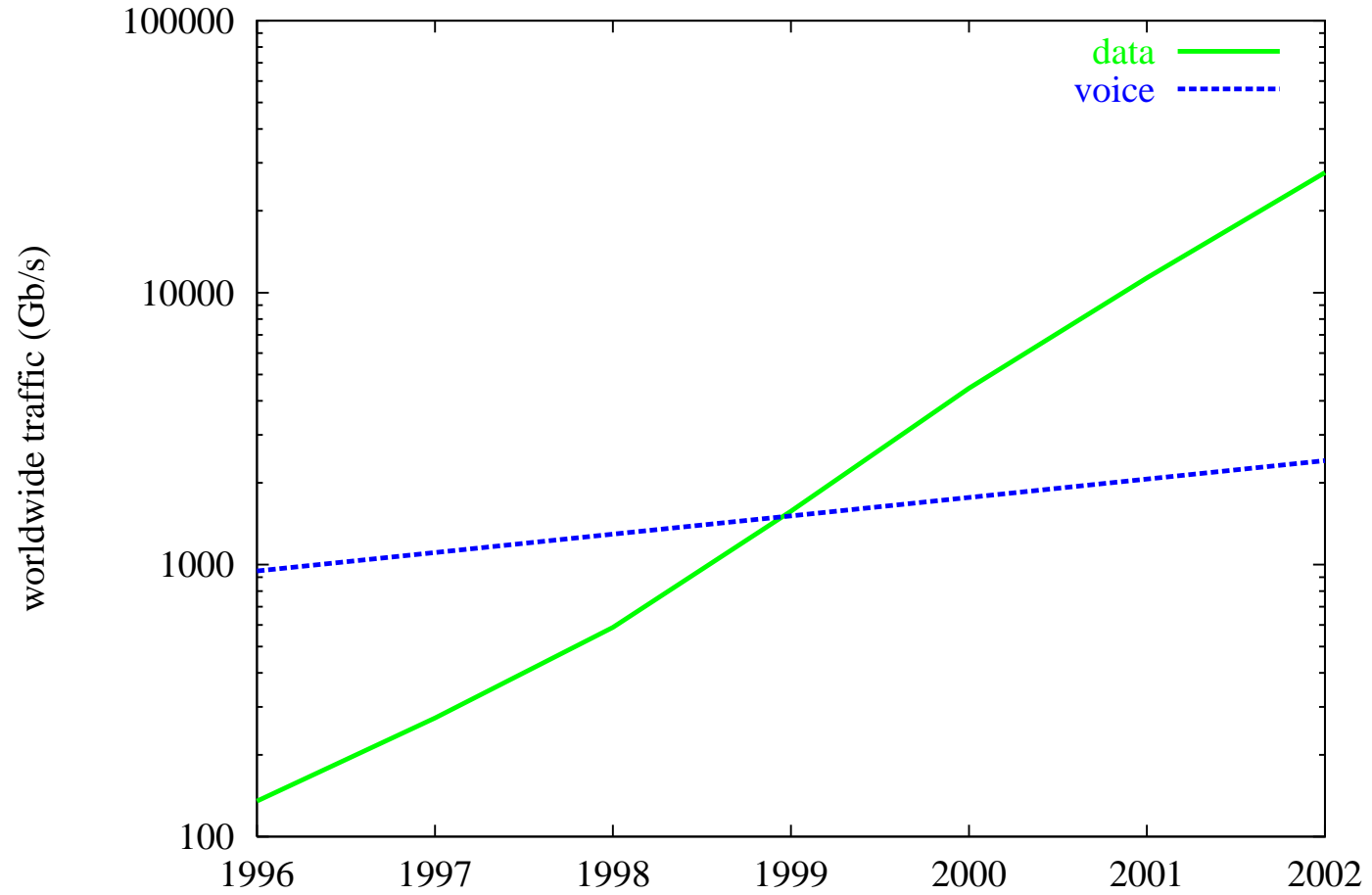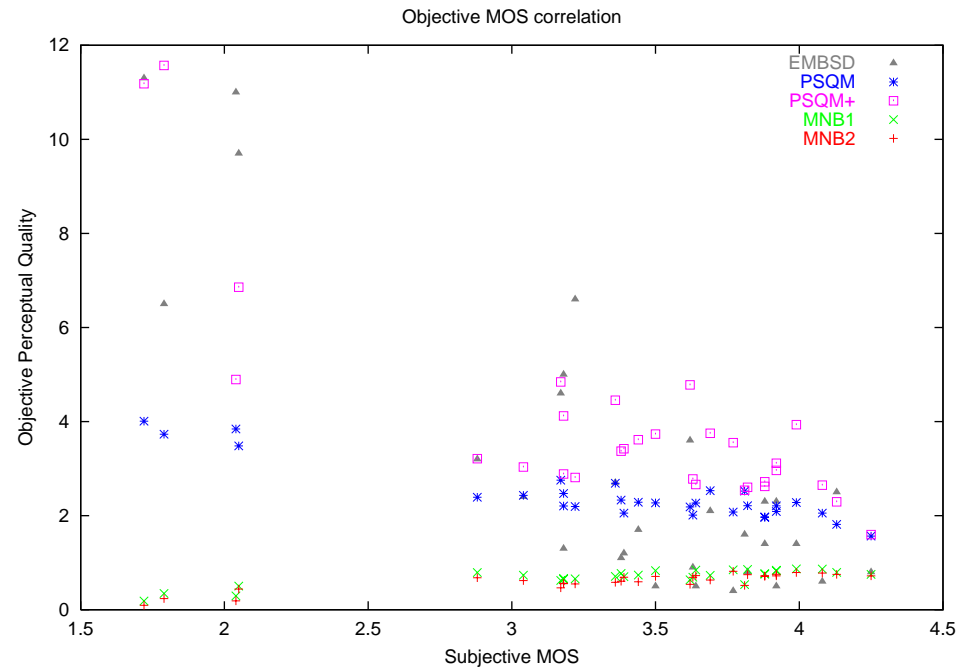| user perspective | carrier perspective |
|---|---|
| <ul><li>variable compression: tin can to broadcast quality</li><li>security through encryption</li><li>caller, talker identification</li><li>better user interface</li><li>internat. calls: TAT transatlantic cable = $0.03/hr</li><li>local calls: possibly cheaper (local access fees)</li><li>easy: video, whiteboard, …</li></ul> | <ul><li>silence suppression ⇒ traffic ↓</li><li>shared facilities ⇒ management, redundancy</li><li>advanced services (simpler than AIN and CTI)</li><li>separate fax, data, voice</li><li>cheaper switching</li><li>better management platforms</li></ul> |

# Audio Codecs

| Codec | rate | quality (MOS) | min. delay |
|---|---:|---|---|
| G.723 | 5.3 | 3.7 | 37.5 ms |
| | 6.3 | 3.98 | 37.5 ms |
| G.729 | 8.0 | 4 | 15 ms |
| AMR | 4.75-12.2 | | 20 ms |
| AMR-WB | 6.6-23.85 | 7 kHz | |
| G.728 | 16.0v | 4 | 5.625 ms |
| G.722 | 32.0 | 7 kHz | 40 ms |
| G.711 | 64.0 | $\mu$-law, MOS 4.3 | var. |

# Voice and Data Traffic

# Objective vs. Subjective MOS

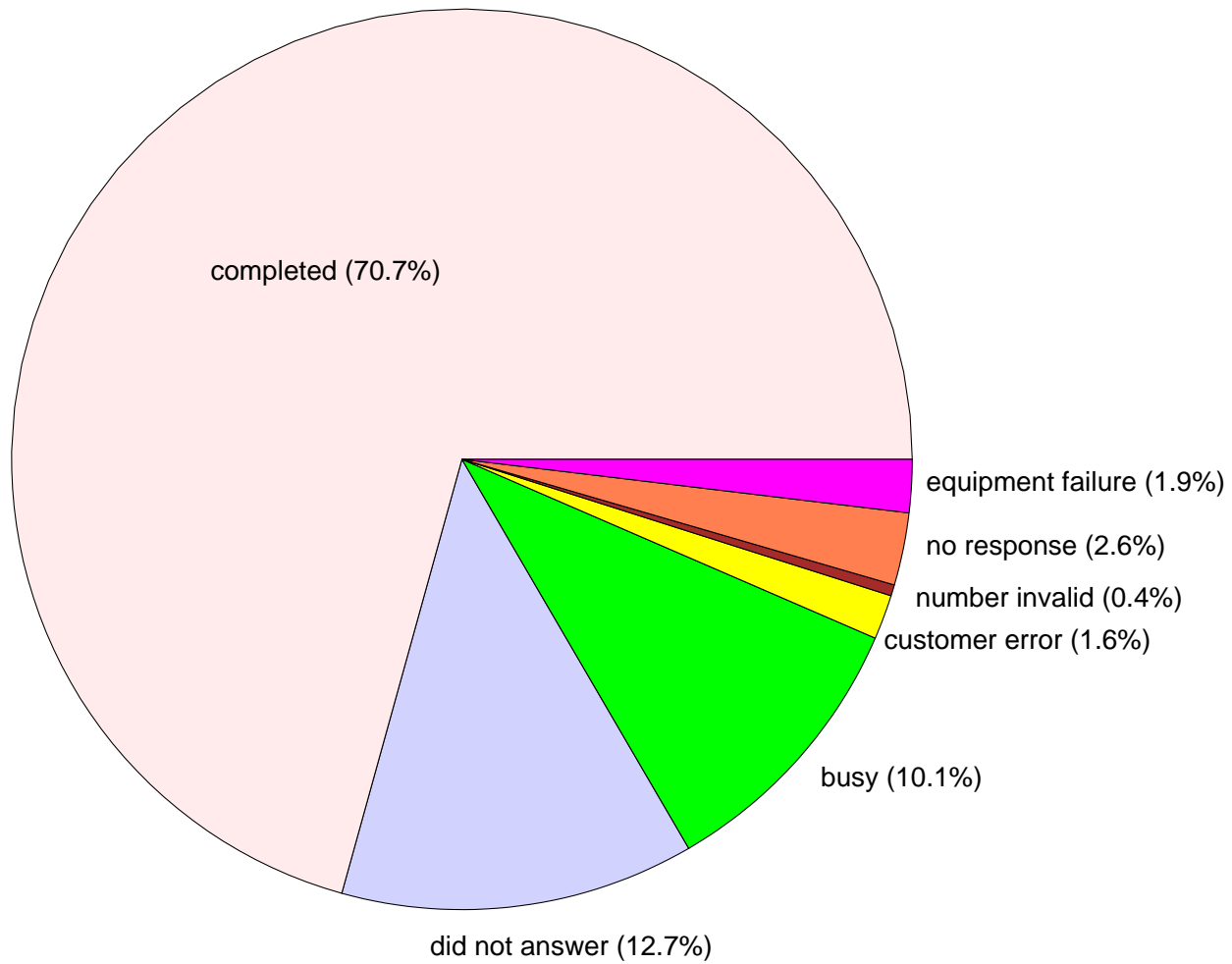Objective MOS tools don't always handle loss impairments correctly:



Objective MOS correlation

# Traffic (1998)

Measured in Dial Equipment Minutes (DEM) or bandwidth:

|                | GDEM | bandwidth (Gb/s) |
|----------------|------|------------------|
| Local          | 2986 | 364              |
| Intrastate toll| 422  | 51               |
| Interstate toll| 555  | 68               |

PBX: typically, about 10% utilization per phone ⇢ 6.4 kb/s per employee (128 Mb/s for 20,000 person campus)

# Call Attempts

completed (70.7%)

equipment failure (1.9%)

no response (2.6%)

number invalid (0.4%)

customer error (1.6%)

busy (10.1%)

did not answer (12.7%)

# Call Setup Delay



8.5

ringback    pick up

15.5    10.9

disconnect

38.1

no answer

1.7

off-hook

8.5

ringback    pick up

8.6    10.9

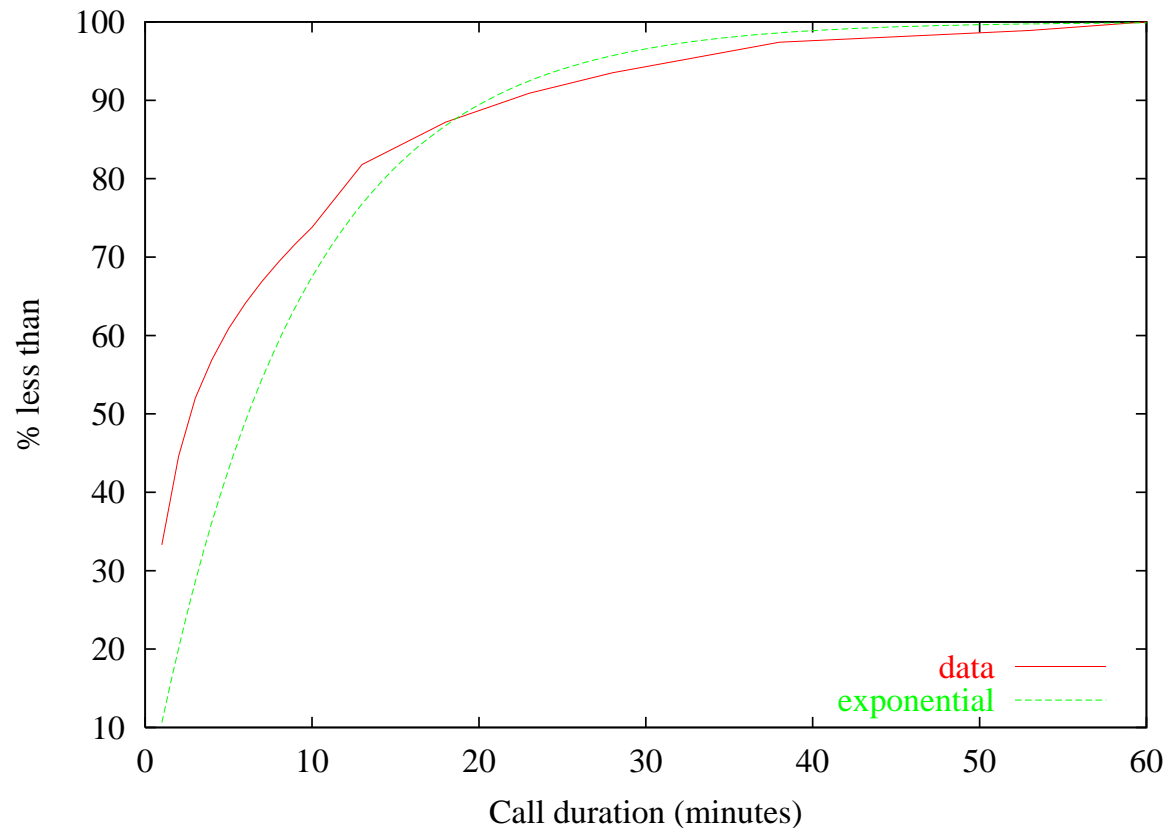disconnect

38.1

no answer

busy    disc.

start dialing    10.5    4.6

# The Three-Minute Myth

Local calls are about 2.4 minutes on average, but long distance calls are much *longer*, about 8.9 minutes:

# Calls Get Longer with Distance

| distance (mi) | % calls | duration (min.) |
|---|---|---|
| 1 – 10 | 5.1 | 4.6 |
| 11 – 22 | 20.2 | 5.1 |
| 23 – 55 | 23.2 | 5.9 |
| 56 – 124 | 13.3 | 7.7 |
| 125 – 292 | 12.1 | 9.4 |
| 293 – 430 | 4.6 | 10.4 |
| 431 – 925 | 9.7 | 11.9 |
| 926 – 1910 | 8.5 | 11.9 |
| > 1910 | 3.2 | 11.2 |
| average | 310 | 7.8 |
| median | 60 | 3.0 |

# Aside: Cost of Bandwidth

- T3 Internet access: $16,000/month

- or 0.05c/minute (for 64 kb/s) for full utilization (bogus)

- typically, assume peak-to-average ratio of 4 (17% during busy hour) ➠ 0.2c/minute

- may be better if data and voice load are offset

- lack of current traffic statistics

# Why Aren't We Junking Switches Right Now?

What made other services successful?

**email:** available within self-contained community (CS, EE)

**web:** initially used for local information

**IM:** instantly available for all of AOL

All of these . . .

- work with bare-bones connectivity ($\geq$ 14.4 kb/s)
- had few problems with firewalls and NATs
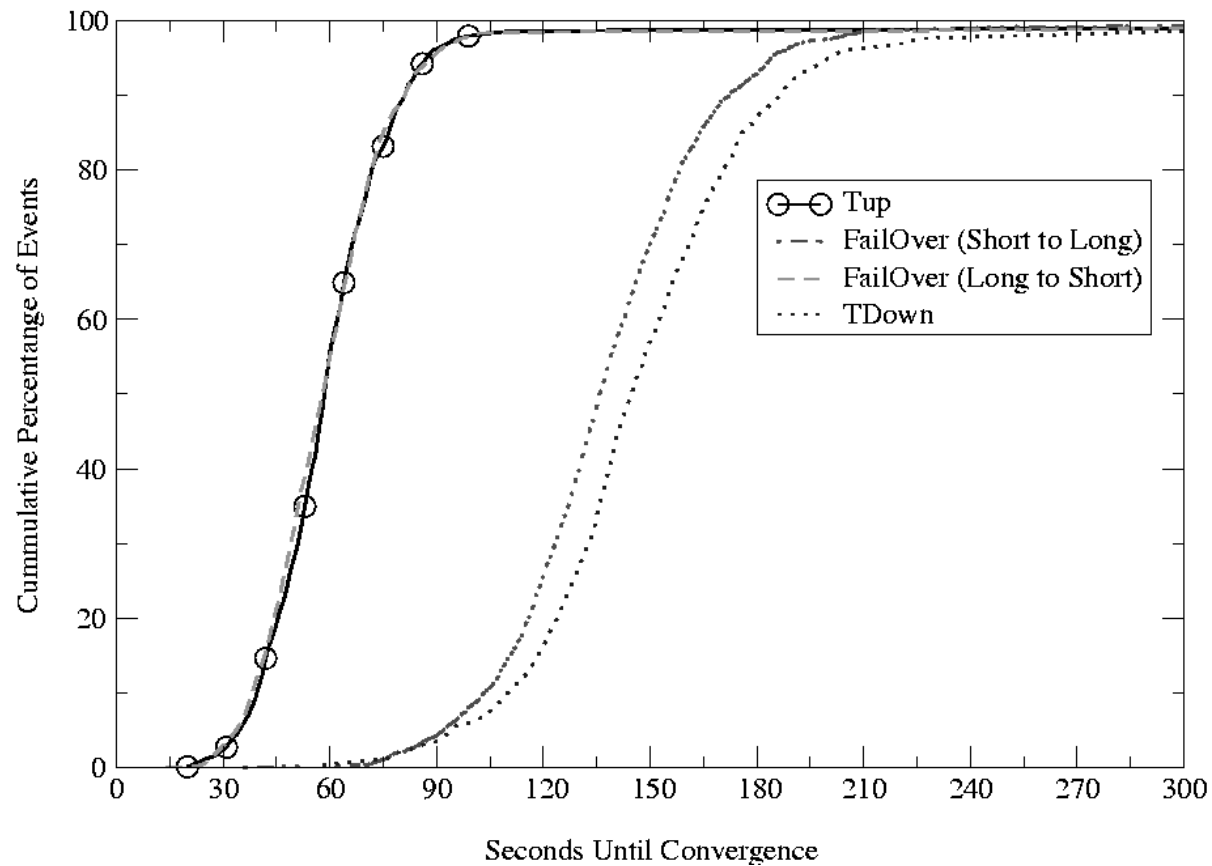- don't require a reliable network

# Reliability Issues

- $-$: software updates require "scheduled downtime"

- $+$: but signaling servers can be made redundant much easier than SS7 SCPs

- BGP convergence times of several *minutes*: 2 minutes to withdraw routes, 30 minutes to advertise routes

- "80% of withdraws take more than a minute"

- no clear IP reliability definition – reachability of any node? some large subset? "local calls"?

# BGP Convergence Times

(From Abha Ahuja's IETF50 plenary talk and Geoff Huston's talk)
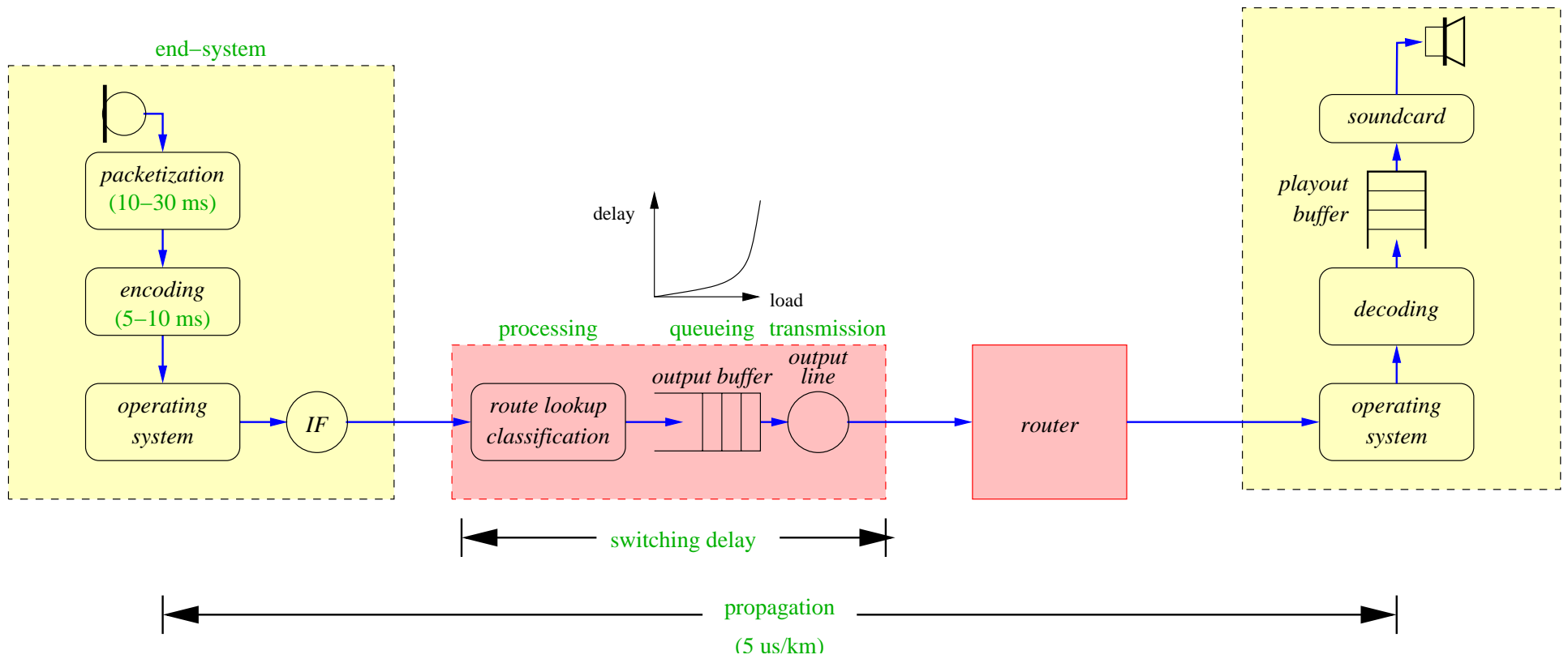
## Failures, Fail-overs and Repairs

# Reliability: Power

- more decentralized ⇒ harder to provide power coverage

- need power for Ethernet switches, phones – ≈ 7W/phone (48V)

- Ethernet powering (spare pairs), tandem or integrated into switch

- also useful for wireless base stations

- Columbia approach: separate power circuit for wiring closets

# Reliability: Denial-of-Service

- denial-of-service and attacks more likely than with traditional phones

- but traditional phones (including 800#) also subject to auto-dialers

- different scenarios:
  - external attack ⇒ can be filtered
  - internal compromise ⇒ spoof DiffServ, RSVP

- disadvantage of integration: no secondary channel

- thus, maybe keep authorized RSVP "circuits"

# Sources of Delay



end–system

packetization
(10–30 ms)

encoding
(5–10 ms)

operating
system

IF

delay

load

processing

queueing

transmission

route lookup
classification

output buffer

output
line

router

switching delay

propagation
(5 us/km)

soundcard
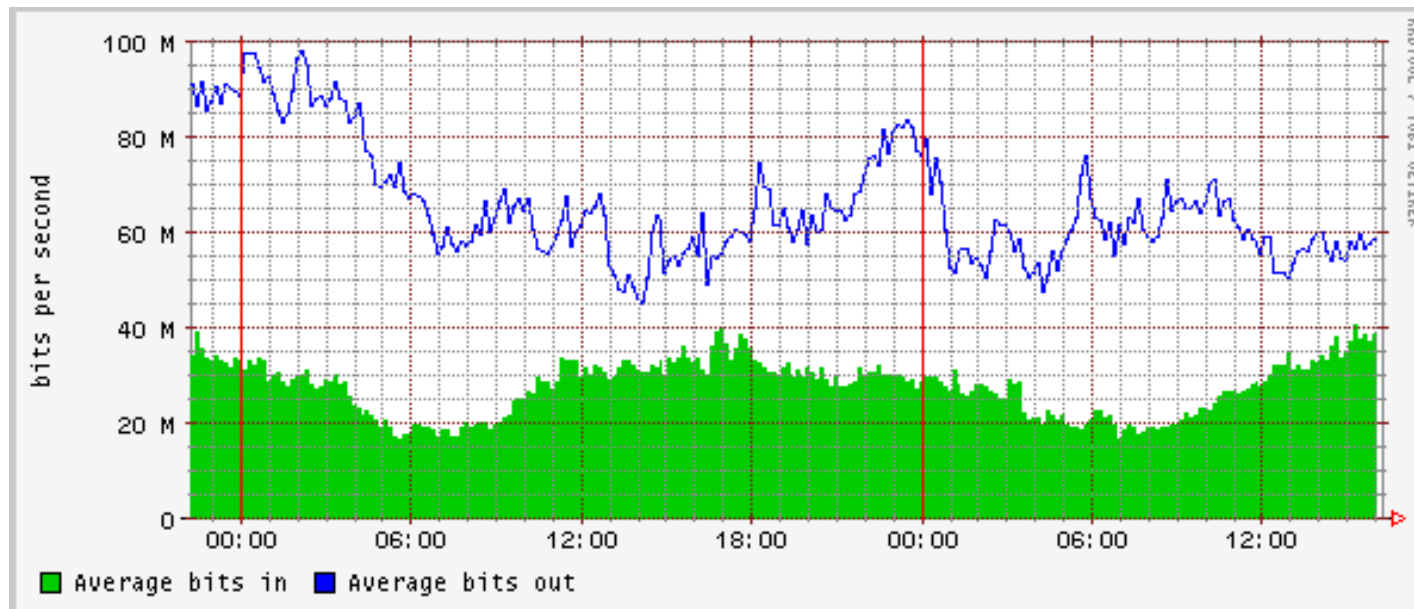
playout
buffer

decoding

operating
system

# QoS: Local Area Network

- typically, very low average utilization (few %)

- very little packet loss (a few packets a day)

- but long delay spikes (300 ms) due to Ethernet collisions if heavy file transfer

- ➠ avoid hubs, even for single office

- ➠ Ethernet prioritization

# QoS: Access Network

- usually, bottleneck (1:10 concentration)

- usually, asymmetrically loaded, depending on web traffic

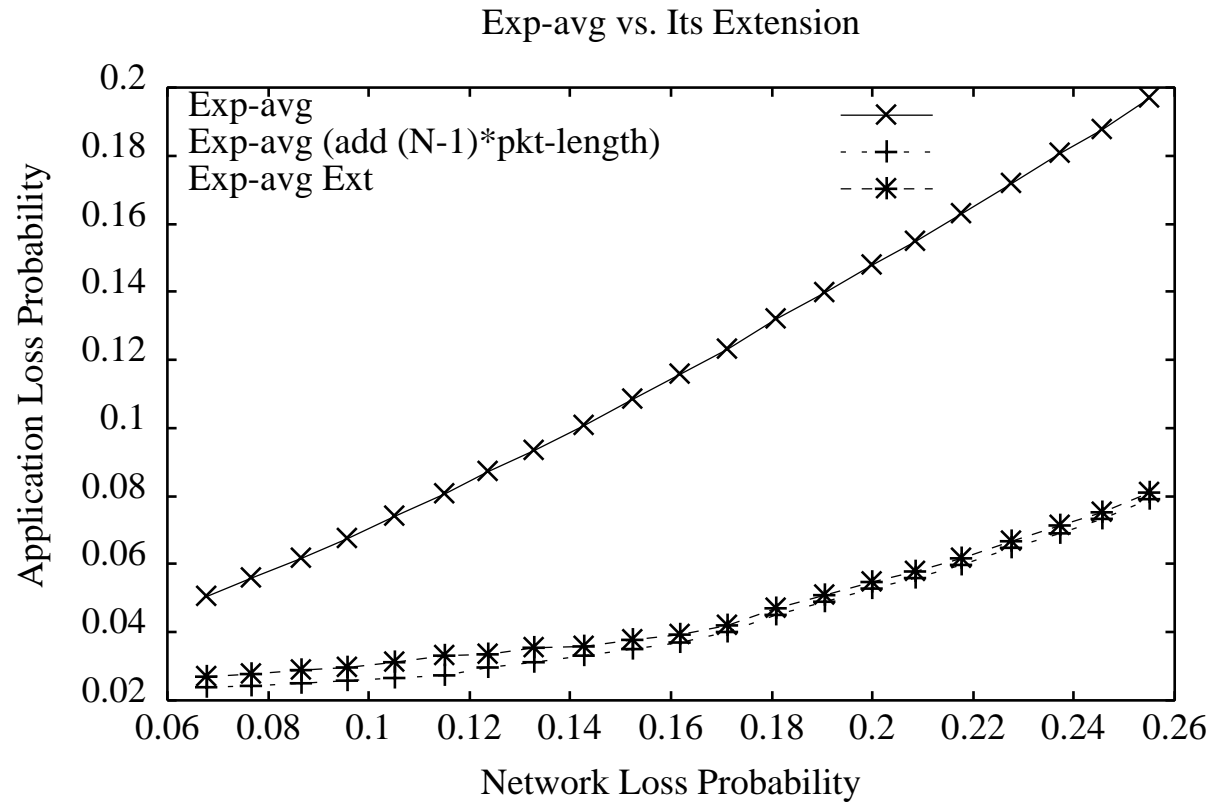- solution: TOS marking (supported by most phones)

# QoS: Wide-Area Network

- existing SLAs and measurements mostly useless: just averages

- e.g.: steady loss of 5% acceptable, one-second bursts of 20% not

- application loss = f(network loss, FEC, jitter, playout delay)

- need rough equivalent of "severely errored seconds"

- however, bursts of loss ⇛ interruptions

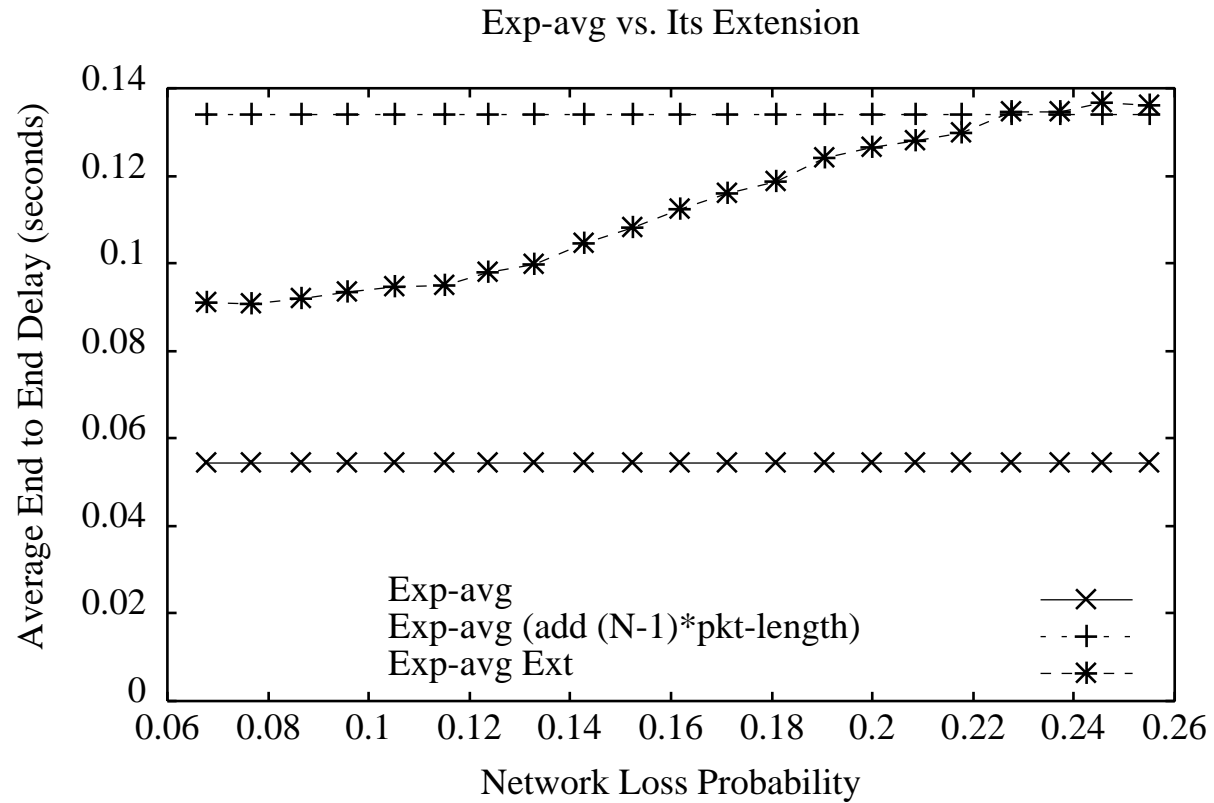- two types of carriers: "classical IP" vs. "voice heritage"?

# Resource reservation

- airline vs. subway: reserve if $> 1\%$ of bottleneck?

- resource reservation likely for upstream cable channel

- RSVP far too complex simple end systems and without multicast

- separate problem: need reserve/commit for VoIP? coupling with application-layer signaling?

- no harm in having several resource reservation protocols

- congestion pricing (RNAP, M2I), including holding costs

# Example: Adaptively Virtual Exponential Average



Exp-avg vs. Its Extension

# Example: Playout Delay

# Architectural Problems for VoIP

VoIP breaks architectural assumptions underlying recent additions:

- NATs: only work for client-server (and TCP)

- VPNs, mobile IP: encapsulation overhead

- firewalls: assume clients inside, servers outside

# Conclusion

- motivation for VoIP

- traffic characteristics

- QoS metrics

- new resource reservation mechanisms?