# Measuring the Evolution of Transport Protocols in the Internet
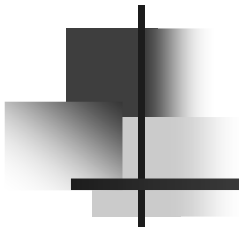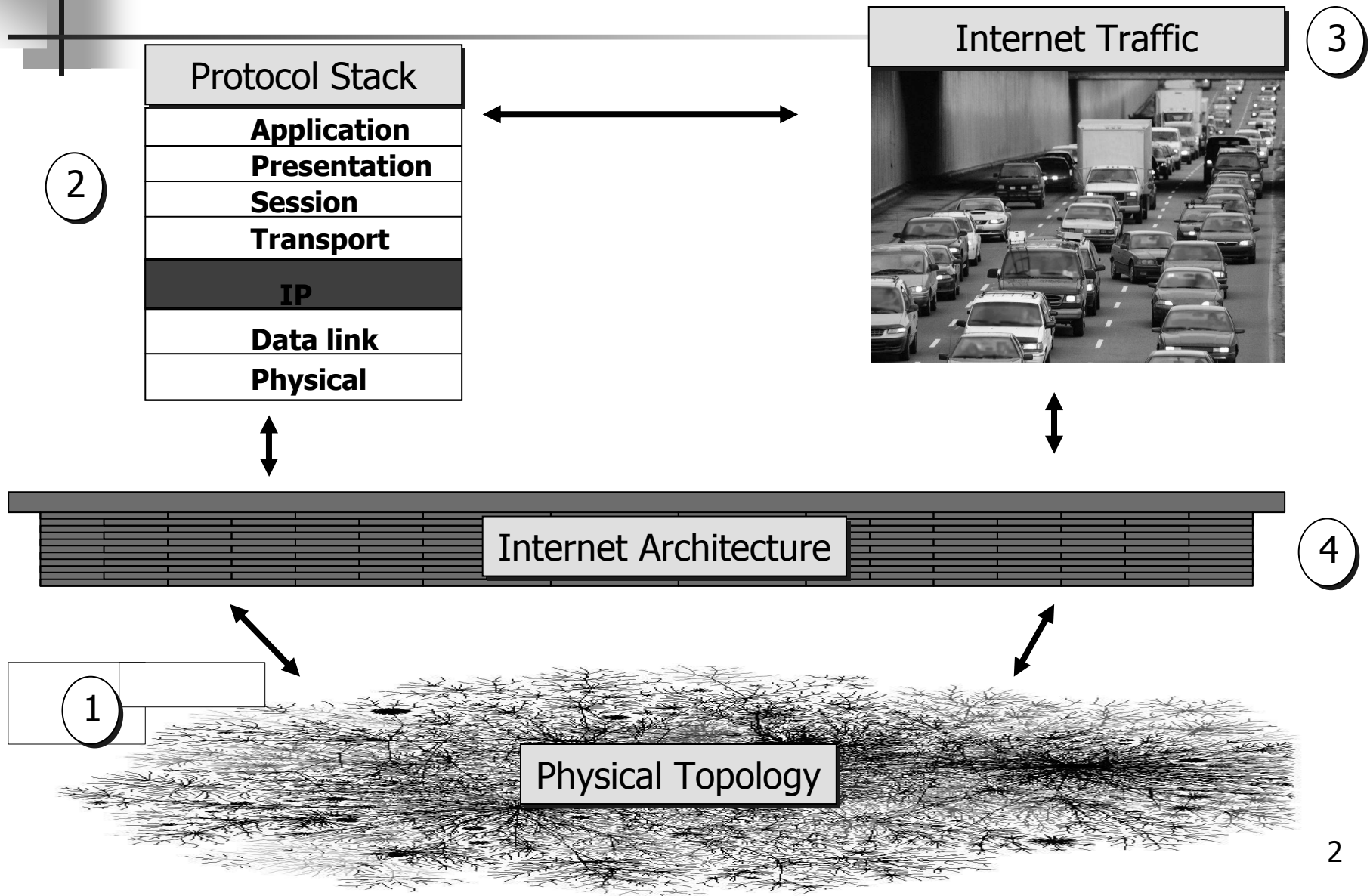
Alberto Medina
Mark Allman
Sally Floyd

# The Internet

**Protocol Stack**

| Application |
| --- |
| Presentation |
| Session |
| Transport |
| IP |
| Data link |
| Physical |

(2)

**Internet Traffic** (3)



**Internet Architecture** (4)
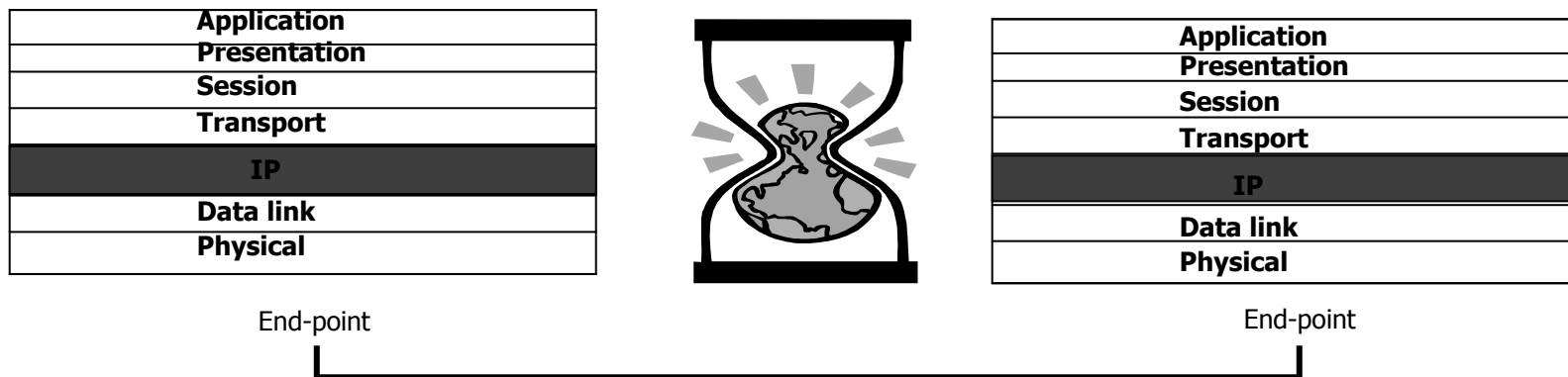
**Physical Topology**

(1)

2

# TCP Evolution

- Congestion Control behavior
    - No proper congestion control => congestion collapse
- Deployment, correctness of transport mechanisms
    - Assess correctness and behavior of newer additions
- Dynamics: Theory vs. Practice
    - Differences between <u>protocol specs</u> (theory) and their <u>implementation and its environment</u> (practice)
- Network Modeling
    - Aim at improving accuracy of network models

# Network Evolution

- ## Hourglass Model



| Application |
|---|
| Presentation |
| Session |
| Transport |
| **IP** |
| Data link |
| Physical |

End-point

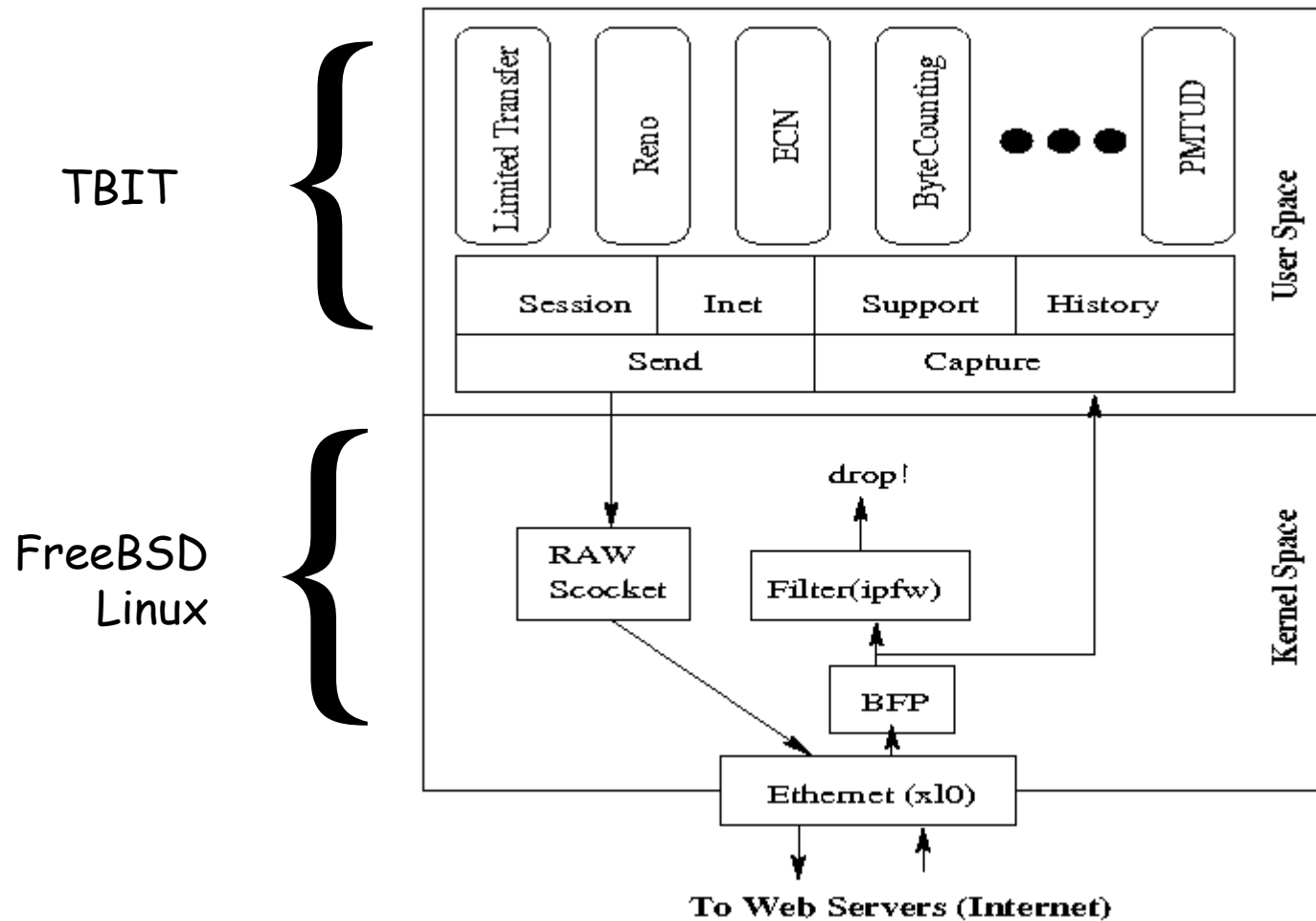| Application |
|---|
| Presentation |
| Session |
| Transport |
| **IP** |
| Data link |
| Physical |

End-point

- ## End-to-end principle
  - ### "Some functions can only be implemented completely and correctly end-to-end, with the help of the end points"
- ## Study effect of middleboxes on these principles
  - ### firewalls, load balancers, NATs, …
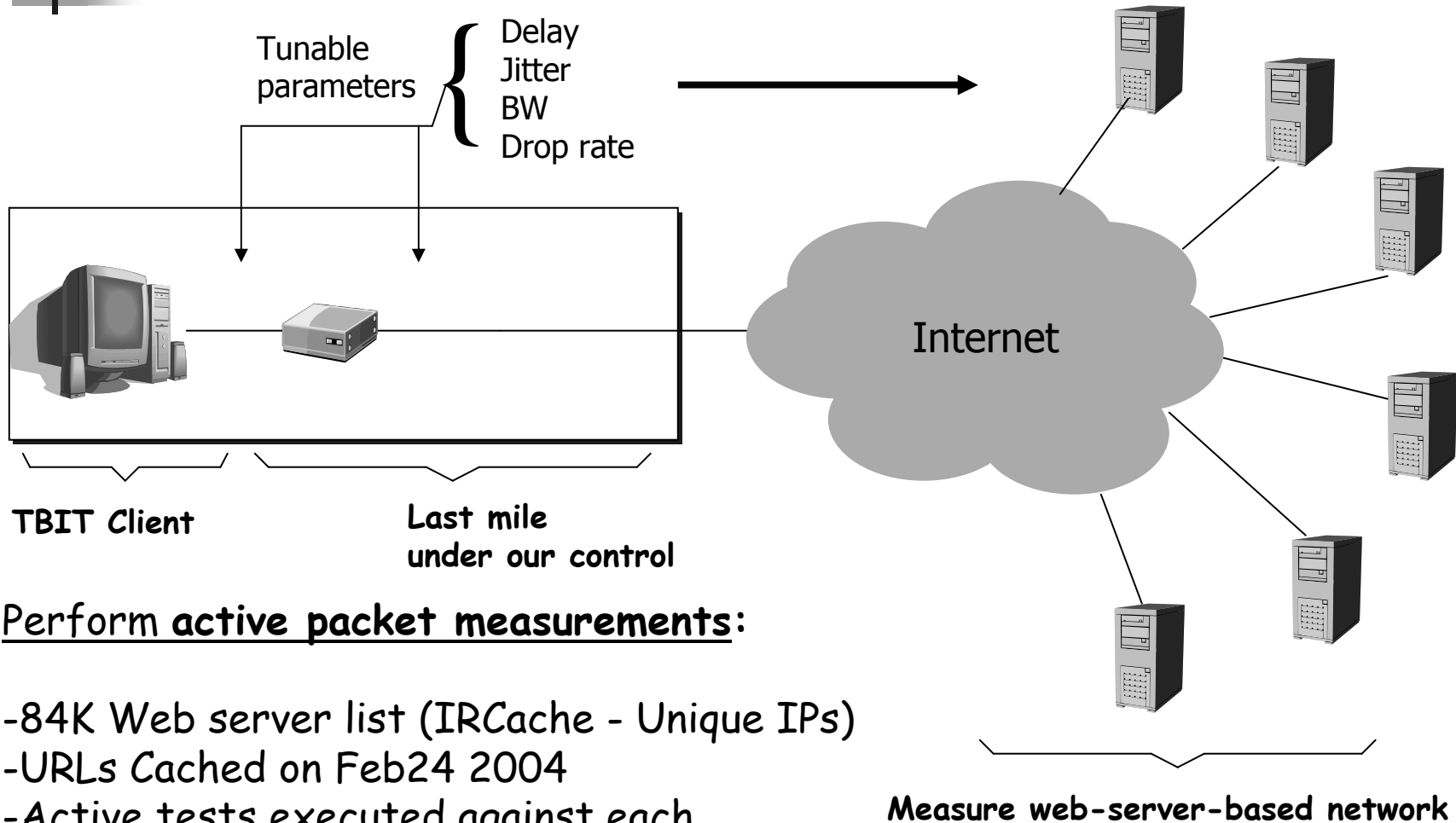
# Experimental Platform

- **Measuring TCP implementations**
  - <u>Passive measurements</u> of web clients
  - <u>Active measurements</u>
    - Web server mechanisms
    - Interactions with environment
- **Active measurements requirements**
  - Measure in-the-field Web servers
  - Employ only conformant TCP traffic
  - Unilateral control at measurement side
- **Employ "undercover" web clients…**

# Undercover Web Clients: TBIT

**TBIT** {

**FreeBSD Linux** {



Limited Transfer | Reno | ECN | ByteCounting | ● ● ● | PMTUD

User Space

Session | Inet | Support | History

Send | Capture

Kernel Space

drop!

RAW Scocket

Filter(ipfw)

BFP

Ethernet (xl0)

To Web Servers (Internet)

# Experimental Platform: Server Side

Tunable parameters

{ Delay
Jitter
BW
Drop rate

Internet

**TBIT Client**

**Last mile
under our control**

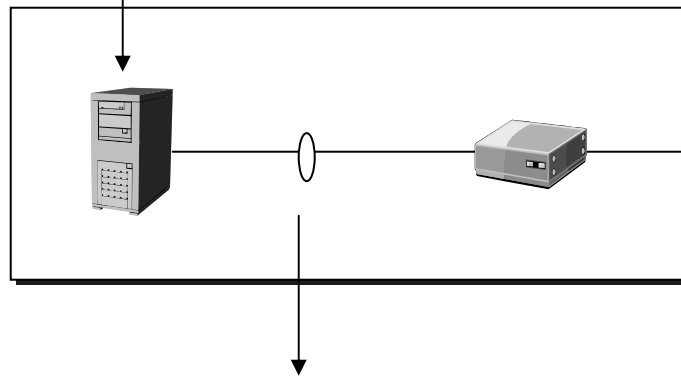**Measure web-server-based network**

Perform **active packet measurements**:

-84K Web server list (IRCache - Unique IPs)
-URLs Cached on Feb24 2004
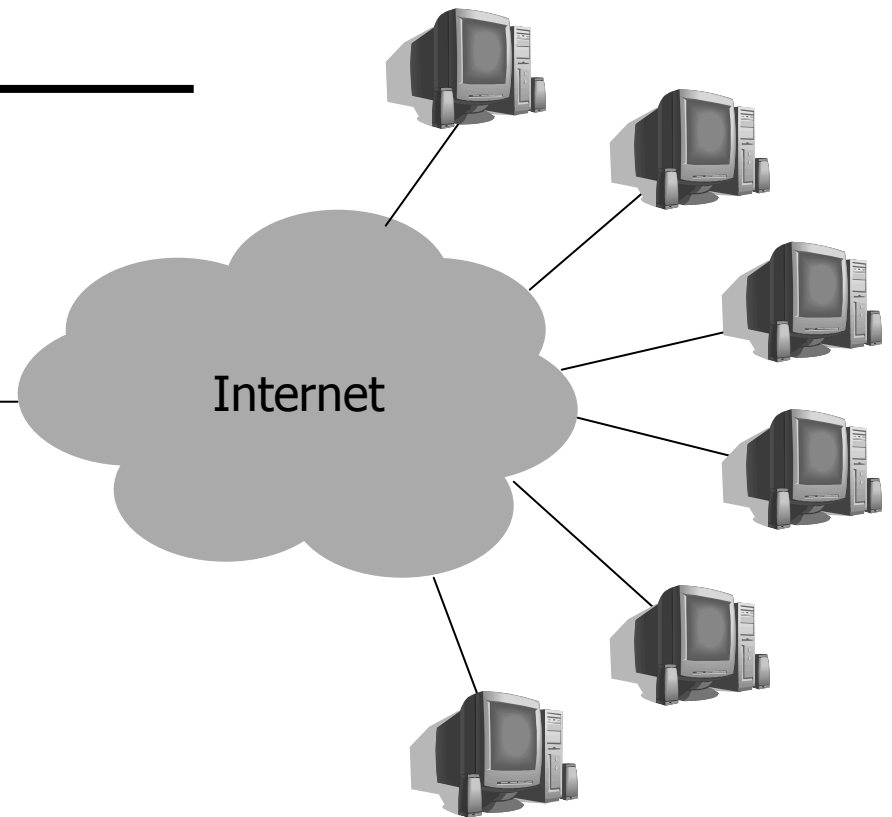-Active tests executed against each
  server in URL list

# Experimental Platform: Client Side

ICSI Web Server
(www.icsi.berkeley.edu)

Internet

Collect **passive packet traces**:

-To and from ICSI server's port 80
-Two-week collection time (Feb24-Mar10)
-206K Connections observed
-28K Clients (e.g. IP addresses)

8

```
                          ┌─────────────────────┐
                          │    MEASUREMENTS     │
                          └─────────────────────┘
```

| Deployment evolution | Internet Architecture Evolution | Tracking Changes |
|---|---|---|

| | | |
|---|---|---|
| TCP stacks Deployment | | |
| SACK Info Processing | ECN Deployment | Reordering, drops |
| SACK Info Generation | Path MTU Discovery | MSS Values |
| D-SACK | Use/Abuse of IP Opts | ICW Values/Perf. |
| Byte Counting | Use/Abuse of TCP Opts | RTO Values |
| Limited Transmits | Middlebox Discovery | Redirections |
| Window-Scale Option | | |
| Window Halving | | |
| Cong. Window Buildup | | |

# Talk Outline

- ✓ Motivation
- ✓ Measurement Platform
- ■ Active Measurements
  - ■ Deployment of Transport Mechanisms
  - ■ Middleboxes and Transport Protocols
- ■ Summary of Results
  - ■ Including client-side
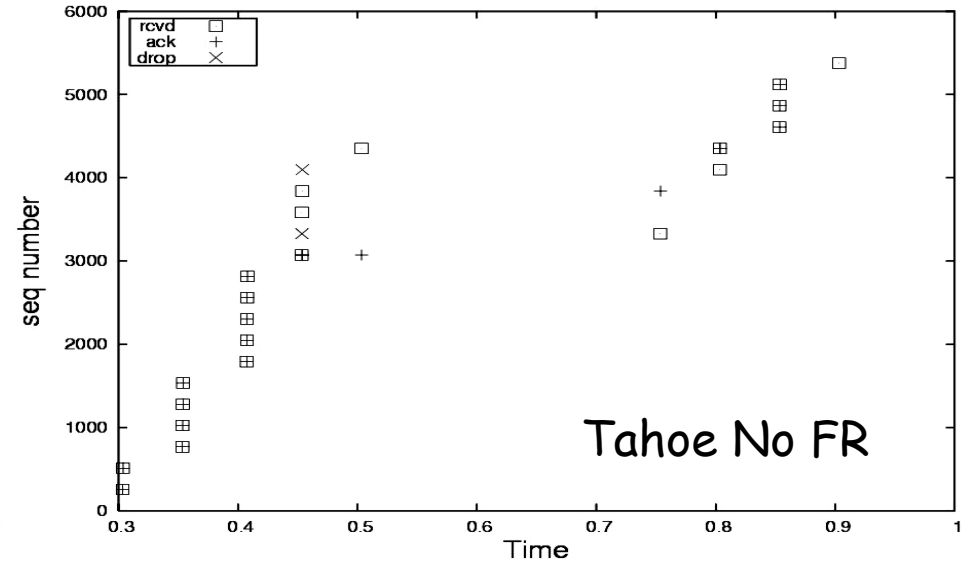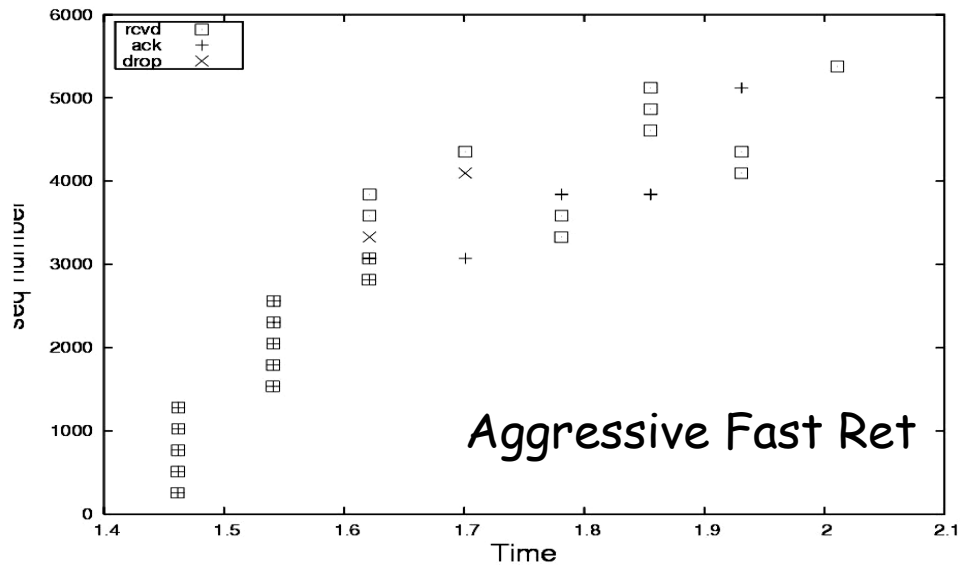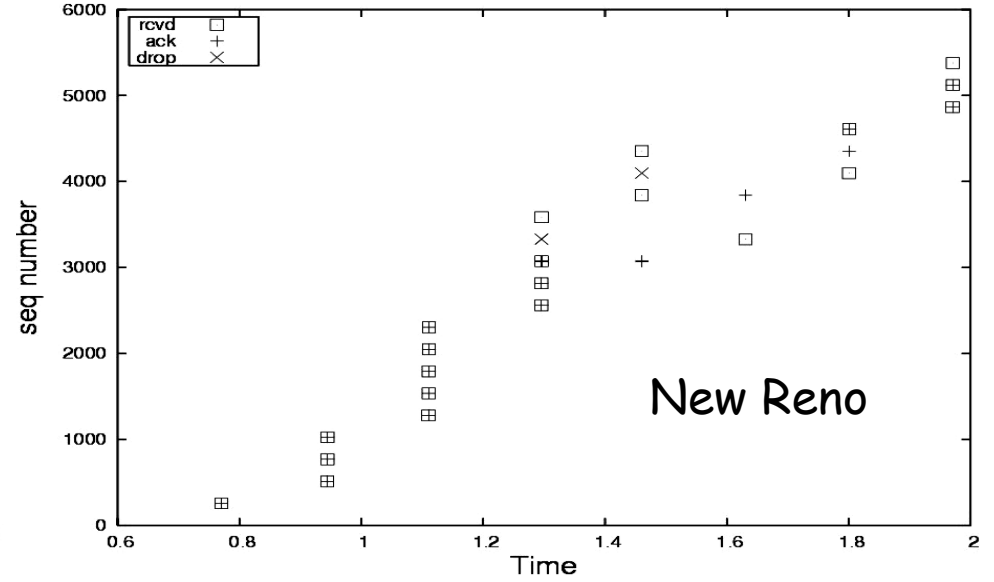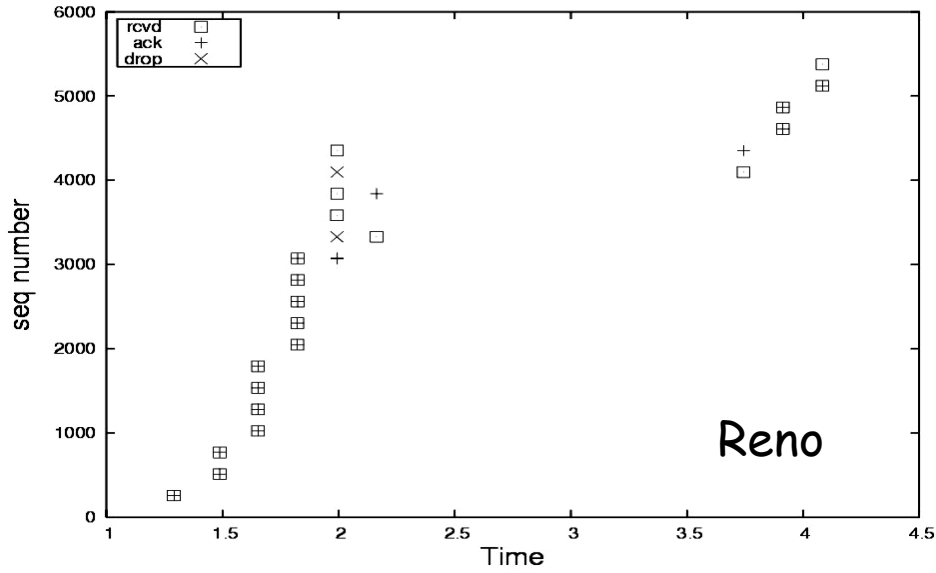- ■ Conclusions
- ■ Future Work

# Deployment of Transport Mechanisms

# Test: Assess Deployment of TCP stacks

- **Establish connection with Web server**
  - Use small MSS
  - Restrict Congestion window to 5 segments
- **Request web page**
- **Receive and ACK incoming packets, but…**
  - Drop packet $13^{th}$
  - Receive and ACK packets $14^{th}$ and $15^{th}$
  - Drop Packet $16^{th}$
  - Continue download until receiving packet $25^{th}$
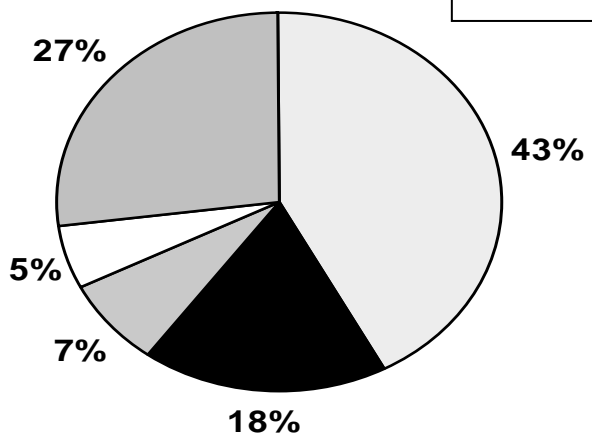
# TCP Behaviors observed
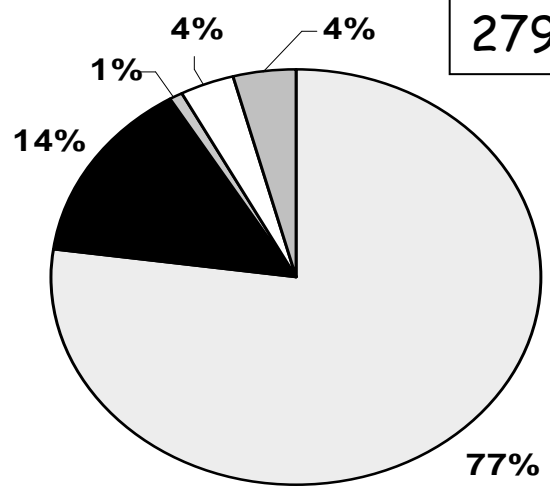
# TCP Deployment Results

**May 2001**

Total: 3728

- New Reno
- Reno
- Reno, Aggr FR
- Tahoe
- Tahoe, No FR

43%
27%
5%
7%
18%

**Feb 2004**

Total: 27914

- New Reno
- Reno
- Reno, Aggr FR
- Tahoe
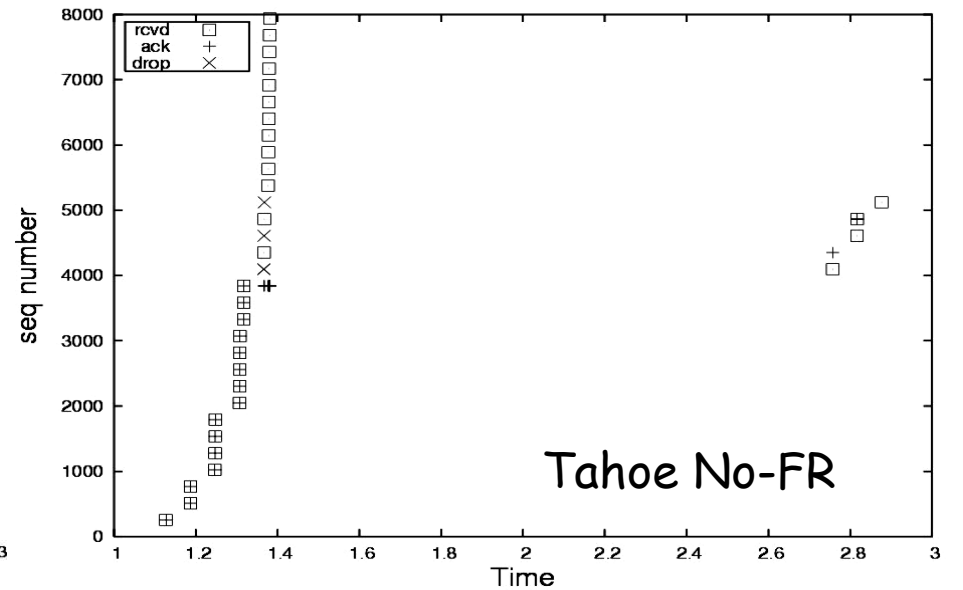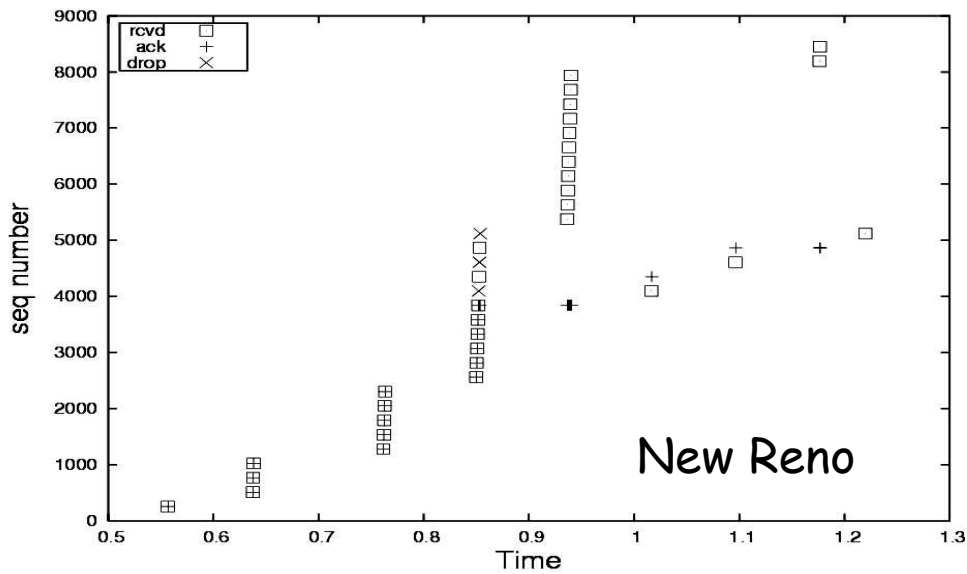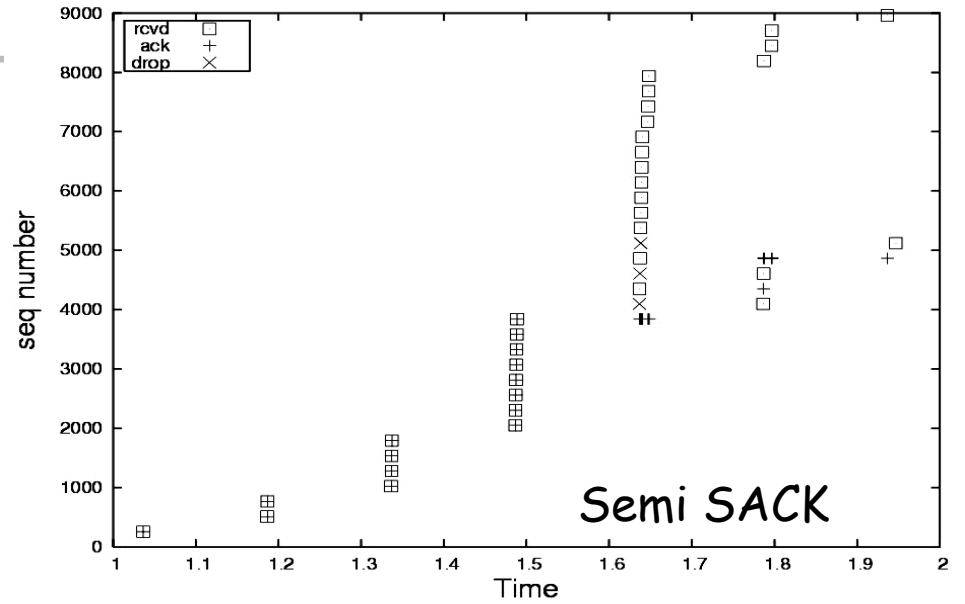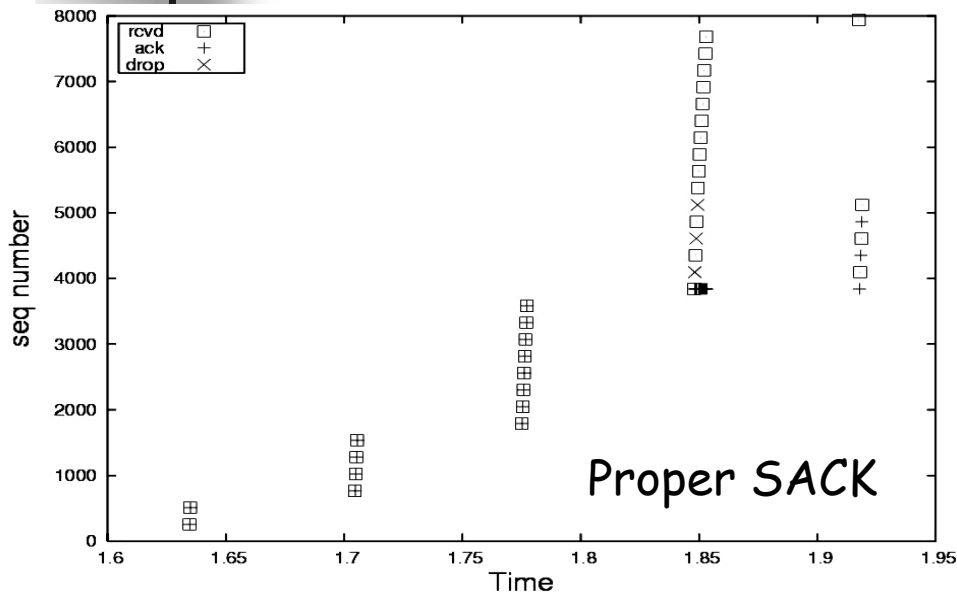- Tahoe, No FR

4%
4%
1%
14%
77%

- Deployment of New Reno increased significantly
- Buggy Tahoe without Fast Retransmit decreased
- Network simulations should use New Reno TCP

14

# Test: Assesing SACK Behaviors
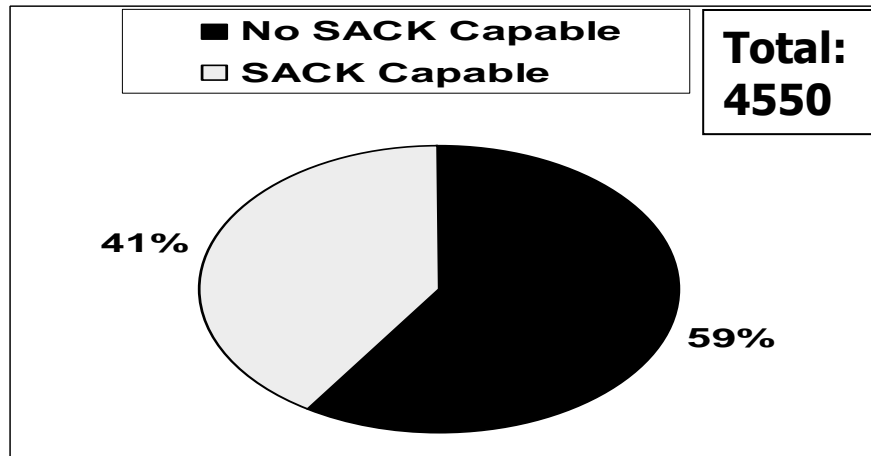
- Negotiate SACK-enabled connection
- If server not capable => NO SACK
- Request web page download
- Receive and ACK incoming packets but…
  - Drop packets $15^{th}$, $17^{th}$, and $19^{th}$
  - Continue receiving and ACKing packets normally (sending appropriate <u>SACK blocks</u> for "drops")
  - Observe retransmission behavior
  - Terminate test, close connection

# SACK Behaviors observed

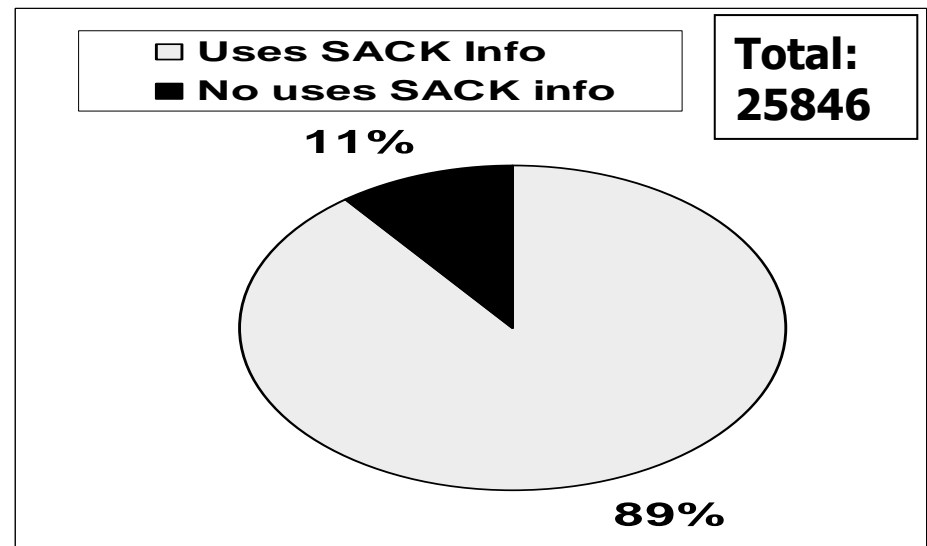# Sack Sender test: Results

## May 2001



**No SACK Capable**
**SACK Capable**

Total: 4550

41%
59%



**Uses SACK Info**
**No uses SACK info**

Total: 1309

42%
58%

## Feb 2004



**No SACK Capable**
**SACK Capable**

Total: 81283

30%
70%



**Uses SACK Info**
**No uses SACK info**

Total: 25846

11%
89%

# Generation of SACK Information

- Do servers generate accurate SACK information?
- Test:

Request: "GET / HTTP 1.1...

| G | X̶ | T | X̶ | / | X̶ |

| H | X̶ | T | X̶ | \b | X̶ | . | X̶ |

"Drop" X-marked packets and update sequence numbers appropriately

```
- - - - - - - - -    SYN (1)    - - - - - - →
←- - - - - - - -    SYN/ACK    - - - - - - - -
- - - - - - - - .      ACK        - - - - →
- - - - - - - -    REQ('G' (2))    - - - - - →
←- - - - - - -      ACK(3)      - - - - - - -
- - - - - - - -    REQ('T' (4))    - - - - →
←- - - - - -    ACK(3, SACK(3))    - - - - - -
- - - - - - -    REQ('/' (6))    - - - - - →
←- - - - -    ACK(3, SACK(3), SACK(5))    - - - - -
```

TBIT                    ...    ...    ...

# Sack Receiver Test: Results

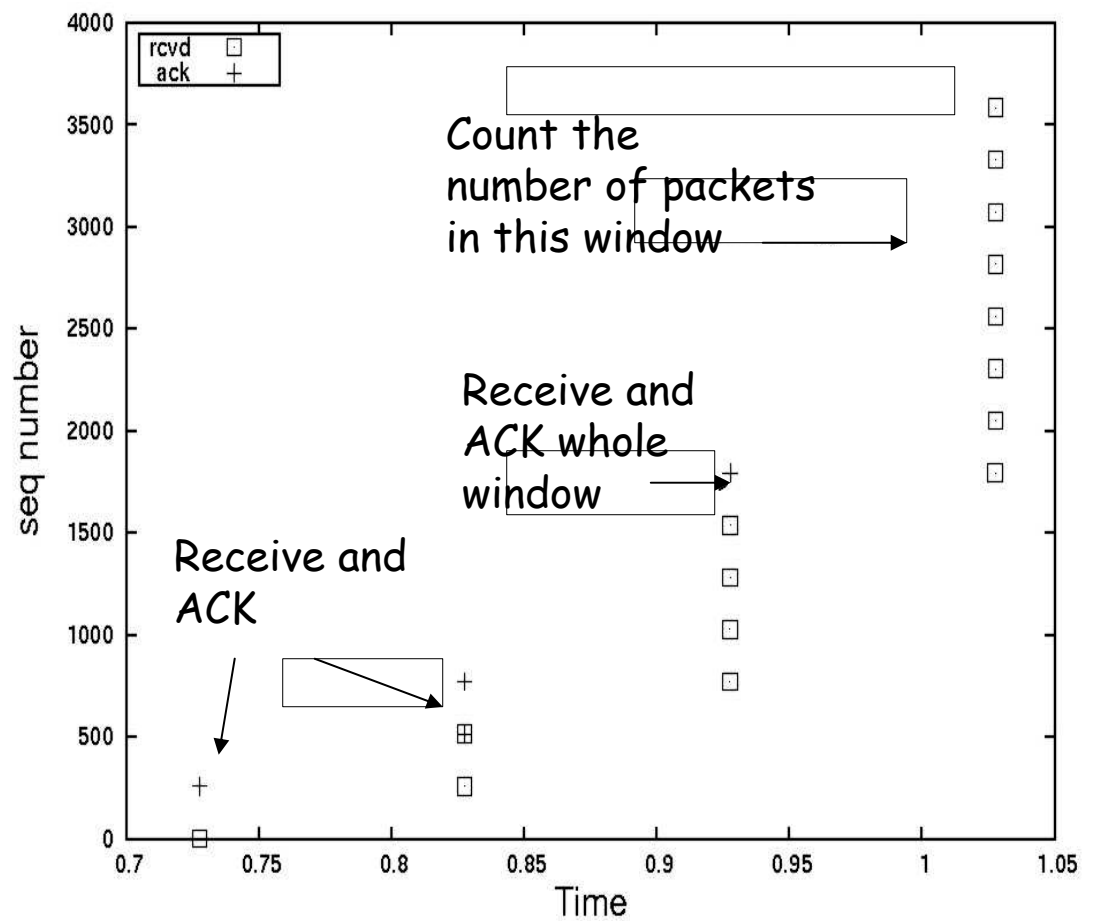| Type of Server | Servers (%) |
|---|---|
| Total servers | 84394 |
| I. Not SACK-Capable | 24361 (28.8%) |
| II. SACK blocks OK | 54650 (64.7%) |
| III. Shifted SACK blocks | 346 (0.5%) |
| I.V. Errors | 5037 (6%) |

- Shifting blocks could have been caused by:
  - NATs, Fingerprint scrubbers,...
- Such middlebox interactions affect any TCP-based communication

# Test: Appropriate Byte Counting (ABC)

- TCP Congestion Control
  - Slow start: increase CWND by one segment for each received ACK
  - Congestion avoidance: increase CWND by 1/MSS for each received ACK
- Drawbacks
  - Delayed ACKs reduce CWND opening rate
  - Mis-behaving receivers may induce servers to open CWND too fast
- ABC Proposal
  - Increase CWND based on bytes ACKed by incoming ACKs, instead of based on number of ACKs received
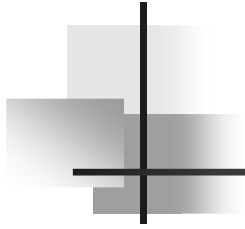
# ABC Test: Example for ICW = 1

- Receive and ACK packets 1, 2 and 3
- Wait for window of 4 packets to arrive
- ACK whole window
- Count number of packets received in next window

Count the number of packets in this window

Receive and ACK whole window

Receive and ACK

rcvd ⬚
ack +

seq number

Time

# ABC Test: Results

| Slow Start Behavior | Number (%) |
|---|---|
| Total number of servers | 44579 |
| I. Classified Servers | 23170 (52%) |
| I.A. Packet Counting | 15331 (51.9%) |
| I.B. ABC | 65 (0.1%) |
| II. Unknown behavior | 288 (0.6%) |
| III. Errors | 21121 (47%) |

- Notice a <u>5-year old</u> proposed mechanism addressing (1) <u>performance</u> concerns and (2) <u>security</u> issues and yet, not being deployed!

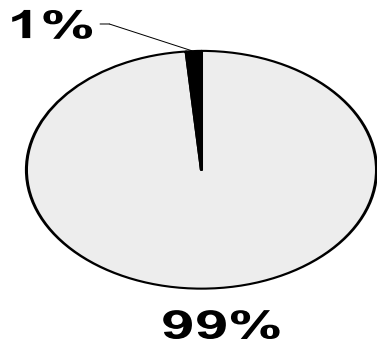# Middleboxes and Transport Protocols

# ECN Capabilities

- **ECN**: Explicit Congestion Notification
  - Allows routers to notify congestion to end nodes
- **TCP**: 2-way handshake
  - Active end: send ECN-Setup SYN (ECN_ECHO, CWR)
  - Passive end: send ECN-Setup SYN/ACK (ECN_ECHO)
- **IP**: 2-bit ECN field in IP header => 4 ECN CPs
  - 00: Not ECT
  - 01: ECT(1) – Sender is ECN capable
  - 10: ECT(0) – Sender is ECN capable
  - 11: Congestion Experienced (CE)

# ECN Test: Results

## May 2001

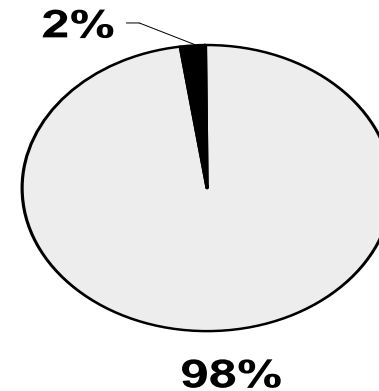□ Not ECN Capable  ■ ECN Capable
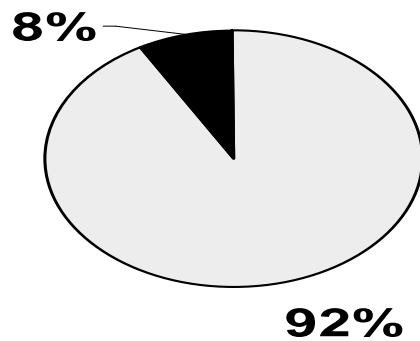
1%

99%

Total: 21879

## Feb 2004

□ Not ECN Capable  ■ ECN Capable

2%

98%

Total: 78733

□ No ECN Echo in ACK  ■ ECN Echo in ACK

8%

92%

Total: 277

□ No ECN Echo in ACK  ■ ECN Echo in ACK

26%

74%

Total: 1765

# ECN Test: Results (2)

## Blocking ECN (3194 Conns)

**■ With ECN □ Without ECN**

25%

75%

Total: 3194

## ECN Code Points (1765)

| ECN CPs data pkts | Number (%) |
|---|---|
| Received pkts w/ ECT 00 (Non-ECT) | 758 (42%) |
| Received pkts w/ ECT 01 (ECT(1)) | 0 (0%) |
| Received pkts w/ ECT 10 (ECT(0)) | 1167 (66%) |
| Received pkts w/ ECT 11 (CE) | 0 (0%) |
| Received both pkts with ECT00 & ECT 10 | 174 (10%) |

# Path MTU Discovery (PMTUD)



MTU=1500

MTU=256

MTU=1500

| SYN | GET "/" HTTP 1.1, MSS = 1500 | ① |

② | | SYN/ACK |

③ | 1500 Bytes of Data, DF=1 |

| ICMP Message | ④ |

⑤ | 256 Bytes of Data | 256 Bytes of Data |

# PMTUD Test: Results

**Legend:** □ No PMTUD □ Proper PMTUD ■ PMTUD Blackholes

20%

34%

46%

**Total: 71737**

- Observed a non trivial number of black holes
- No hope for new ICMP-based proposed mechanisms
  - Explicit corruption notification
  - Handoff notification

28

# Interference with TCP/IP Options

- TCP/IP options
  - Allow encoding additional information at end of packets
- Several concerns raised about using IP Options
  - Overhead, misalignment problems, DoS attacks
- Solutions to concerns
  - Range from OS patches to dropping "offending" packets
- Issue concerns protocol designers
  - Use of unused TCP/IP options in new proposals
  - Ex: QuickStart (QSR) IP Option
- TCP/IP Options tests
  - Evaluate connections with SYN-packet TCP/IP options
  - Evaluate connections with Mid-Stream TCP/IP options

# IP Options Test: Results

**IP Options – SYN Packets**



- Severe interference with known and unknown IP options
- Negative results for of new IP-option-based mechanisms

# Summary of Results: TCP Evolution

| TCP Mechanism | Conclusion |
|---|---|
| TCP Cong. Ctrl | ~ 2/3 use **NewReno** => Use it in ns |
| Loss Recovery | **SACK-cap Prevalence**: ~ 2/3 servers, 9/10 clients<br>Most claiming SACK, do SACK properly<br>**SACK info**:  (mostly) correct |
| DSACK | ~ ½ of SACK-capable servers, send D-SBs |
| ABC | Not deployed |
| LT | Not fully deployed (~1/4 of servers) |
| MSS | Most clients use ~ 1.4K bytes<br>Most servers accept << 1.4K bytes |
| RTO | Many servers use RTO < 1s |

# Summary of Results: TCP Evolution (2)

| TCP Mechanism | Conclusion |
|---|---|
| ICW | Many ICW = 1, Most used ICW = 2-4<br>Some gain from larger ICWs<br>No changes for reordering and losses |
| Window Scaling | Most servers support WS (shift count=0) |
| Window Halving | Most servers do proper window halving<br>Some servers use CWND without caring for RWND |
| Window Buildup | Most servers do no increase cwnd if not used |
| Advertised Window | Most clients surveyed advertise 64KB windows<br>Many clients advertise 8KB and 16KB |
| ECN | Very few servers using ECN (~2.3%)<br>1% Increase since 2001 |

# Summary of Results: Network Evolution

| Behavior | Conclusion |
| --- | --- |
| SACK | Small number of cases, web servers and clients receive SACK blocks with incorrect sequence numbers |
| ECN | Roughly 1% of refused connections |
| PMTUD | < ½ servers PMTUD-capable<br>Likely routers/middleboxes blocking ICMP messages for 1/6 of the servers |
| IP Options | Many failures (1/3) when IP RR or TS SYN options used<br>Majority of failures (70%) if unknown IP option used |
| TCP Options | More resilient and tolerant than unknown options |
| Reordering | Significant small-scale reordering |

# Conclusions

- Achieved set goals
  - Tracked deployment of transport mechanisms
  - Evaluated transport-network interactions
- Competition of interests complicates deployment priorities
  - ABC not implemented, LT implemented
  - PMTUD failing, no ECN deployment,...
- Pinpointed specific cases exemplifying how evolving network challenges end-to-end principle
  - Fundamental design principles of Internet have changed
  - Current network needs to evolve towards new reality

# Future Work

- Further TCP in-the-field behavior
  - Restart behavior after an idle period; Backoff behavior
  - Behavior in other environments (p2p, wireless,…)
- Study other protocols and mechanisms
  - <u>Existing</u>: UDP, FTP, HTTP, RTP, …
  - <u>New</u>: AQM, High-Speed TCP, SCTP, DCCP,…
  - TCP in other environments (P2P, web caching,…)
- Further exploration of Middlebox behavior/impact
  - Many open questions (e.g. How about PEPs?)
  - Detecting middleboxes
- Continuous Monitoring Platform
- Active measurements of client behaviors
- Unilaterally-controlled Active Measurements

# Contact & Information

- **People:**
  - Alberto Medina: medina@icir.org
  - Sally Floyd: floyd@icir.org
  - Mark Allman: mallman@icir.org
- **Software and data**
  - http://www.icir.org/tbit

*icir*

*i.c.s.i. center for internet research*