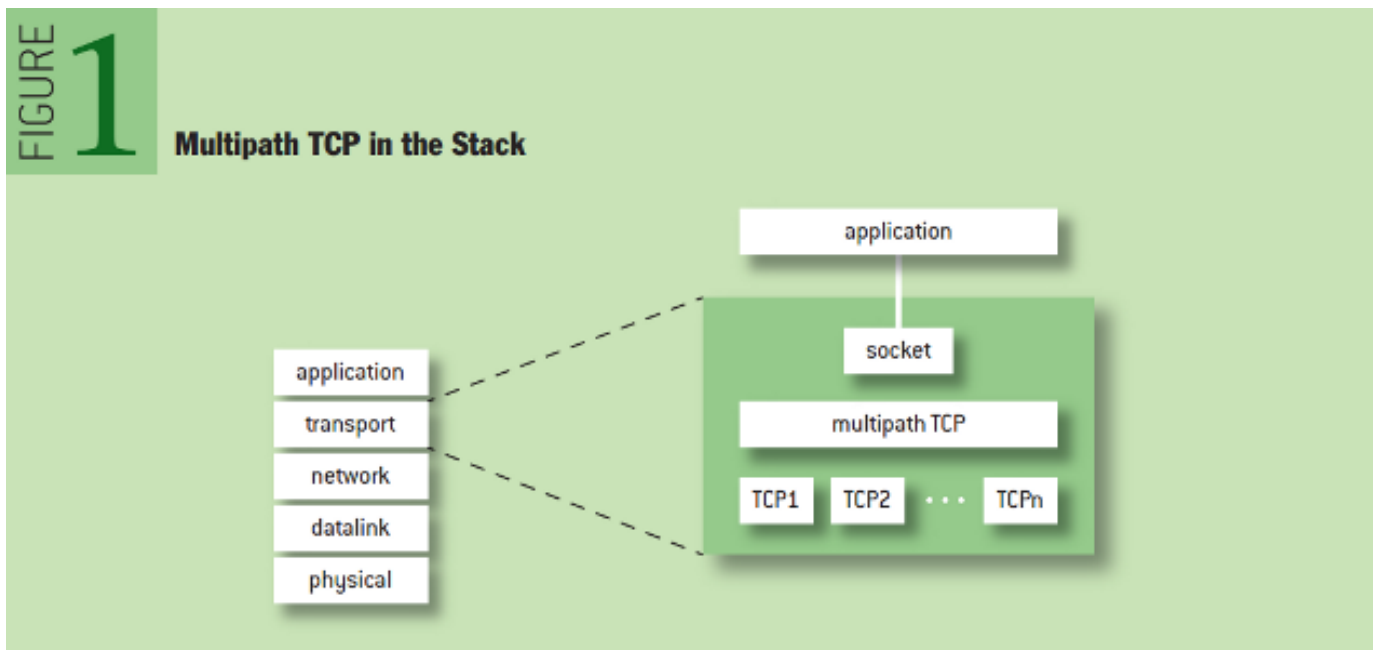# CS 204:
# Multipath TCP

Jiasi Chen

Lectures:  MWF 12:10-1pm in WCH 139

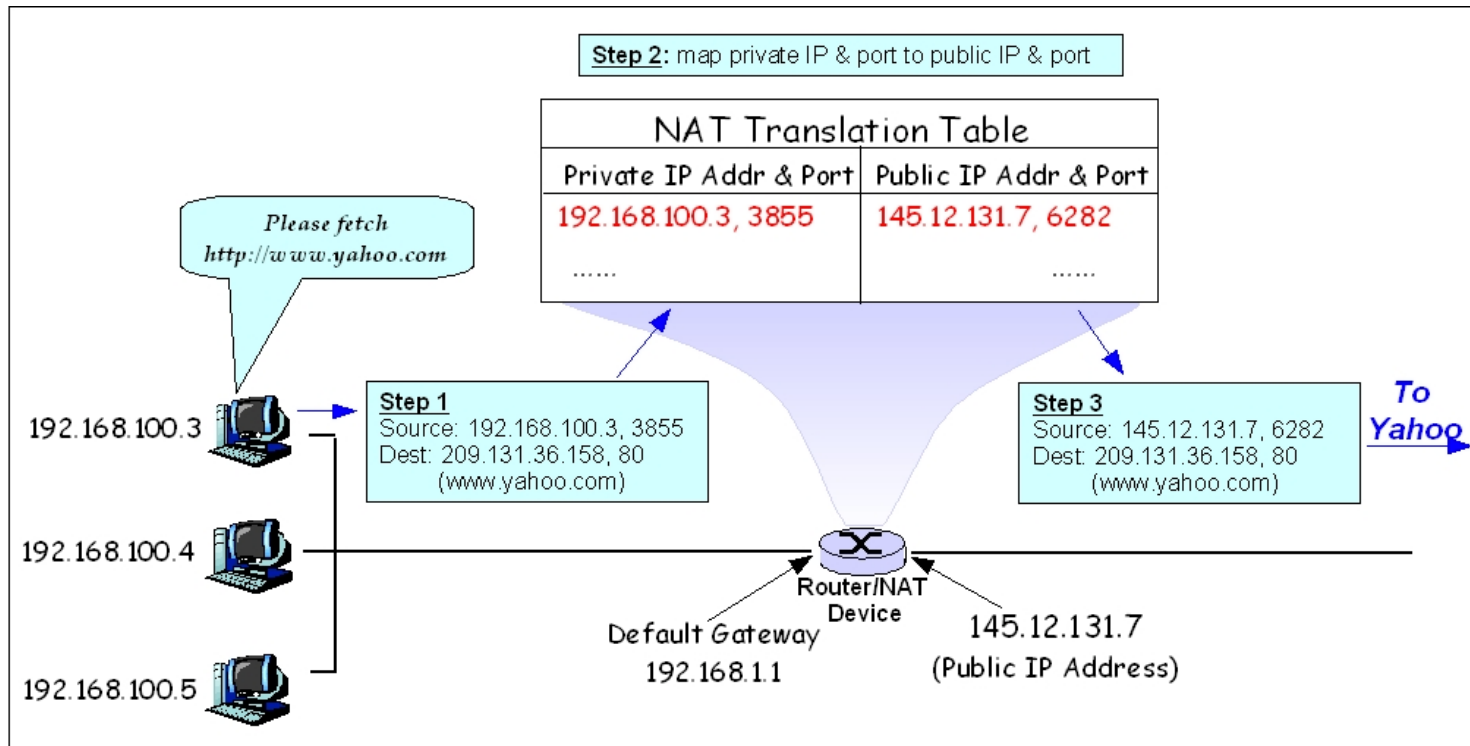http://www.cs.ucr.edu/~jiasi/teaching/cs204_spring16/

# Goals

- Use the available network paths at least as well as regular TCP, but without starving TCP.

- Usable as regular TCP for existing applications.

- Enabling MPTCP must not prevent connectivity on a path where regular TCP works.
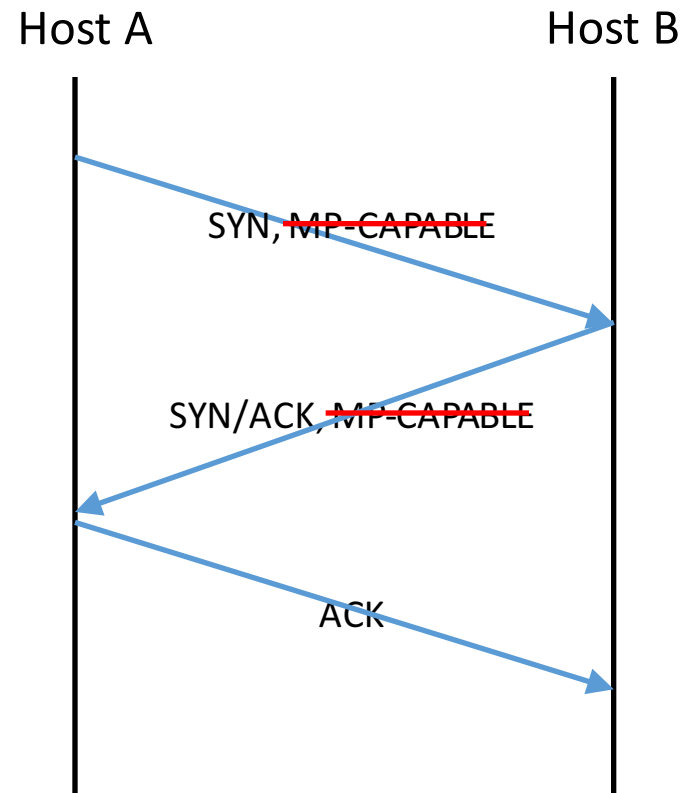
# Network Stack



FIGURE 1

**Multipath TCP in the Stack**

application
transport
network
datalink
physical

application
socket
multipath TCP
TCP1   TCP2   · · ·   TCPn

# Network Address Translators

# Connection Setup

- Use MP-CAPABLE flag to indicate sender has MPTCP capability

- **Problem:** Middleboxes remove TCP options

- Option removed on msg 1?

- Option removed on msg 2?

Host A                                    Host B
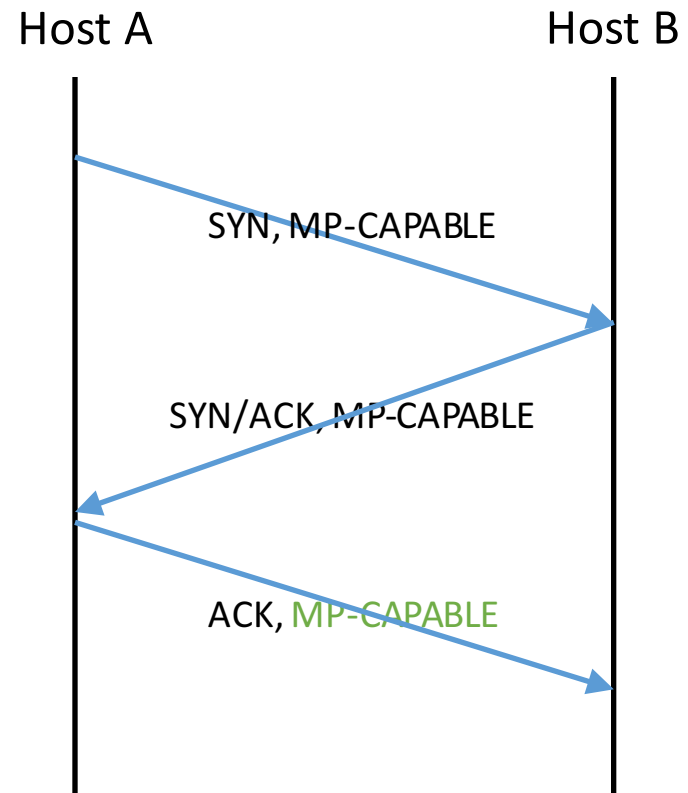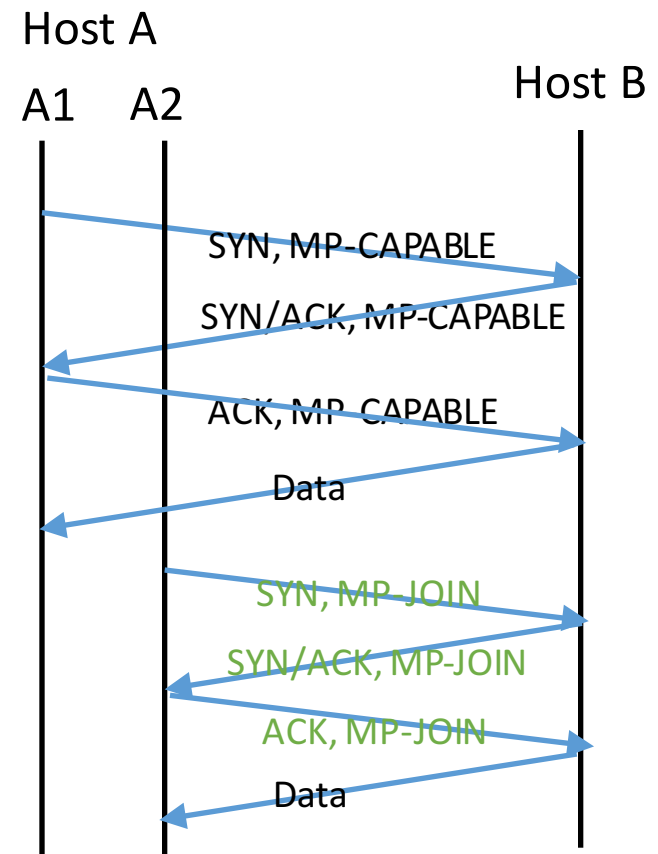
SYN, ~~MP-CAPABLE~~

SYN/ACK, ~~MP-CAPABLE~~

ACK

# Connection Setup

- Use MP-CAPABLE flag to indicate sender has MPTCP capability

- Problem: Middleboxes remove TCP options

- Option removed on msg 1?
  → fall back to TCP

- Option removed on msg 2?
  → host A and host B's views are inconsistent
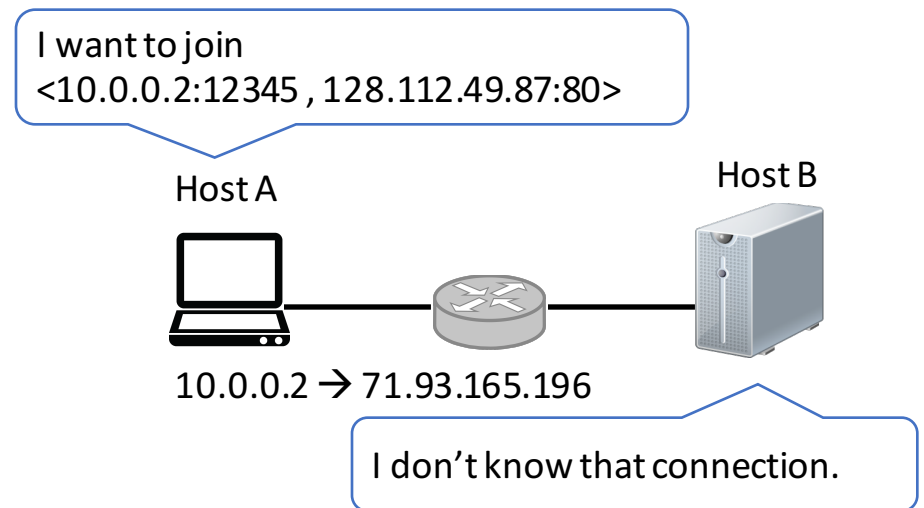  →add another MPT-CAPABLE to msg 3 if MP-CAPABLE recv'd in msg 2

Host A                                    Host B

SYN, MP-CAPABLE

SYN/ACK, MP-CAPABLE

ACK, MP-CAPABLE

6

# Adding New Flows: Naïve solution

- Host A has addresses A1 and A2
- Assume Host B knows these addresses and starts sending data to both

- Problem: Middleboxes will not allow data to be sent without SYN → need 3-way handshake for new subflows

Host A

A1     A2                          Host B

SYN, MP-CAPABLE

SYN/ACK, MP-CAPABLE

ACK, MP-CAPABLE

Data

SYN, MP-JOIN

SYN/ACK, MP-JOIN

ACK, MP-JOIN
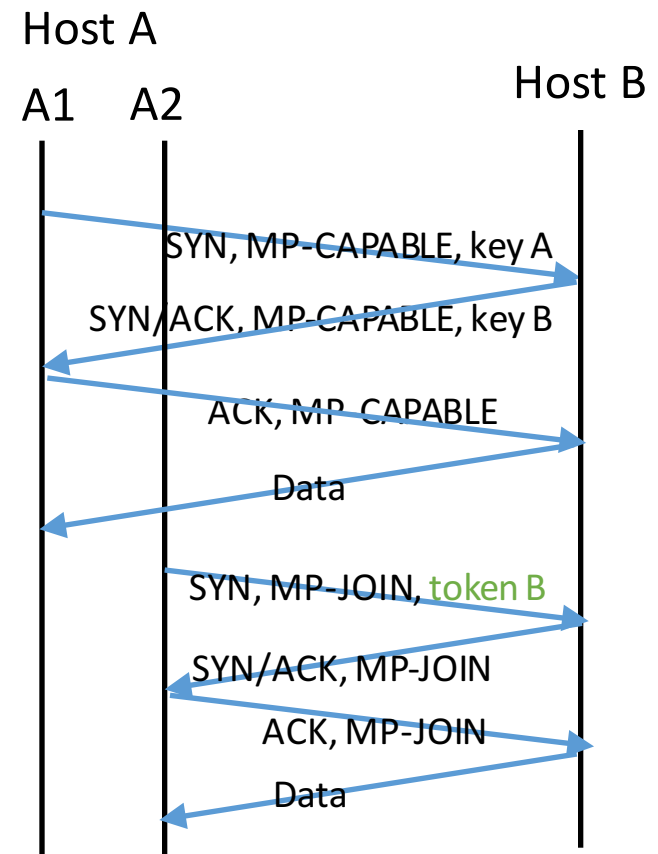
Data

# Adding New Flows: Identification

- TCP flows traditionally identified by <source IP, source port, dest IP, dest port>

- **Problem**: when adding new subflow to existing connection, don't know the source IP

I want to join
<10.0.0.2:12345 , 128.112.49.87:80>

Host A

Host B

10.0.0.2 → 71.93.165.196

I don't know that connection.
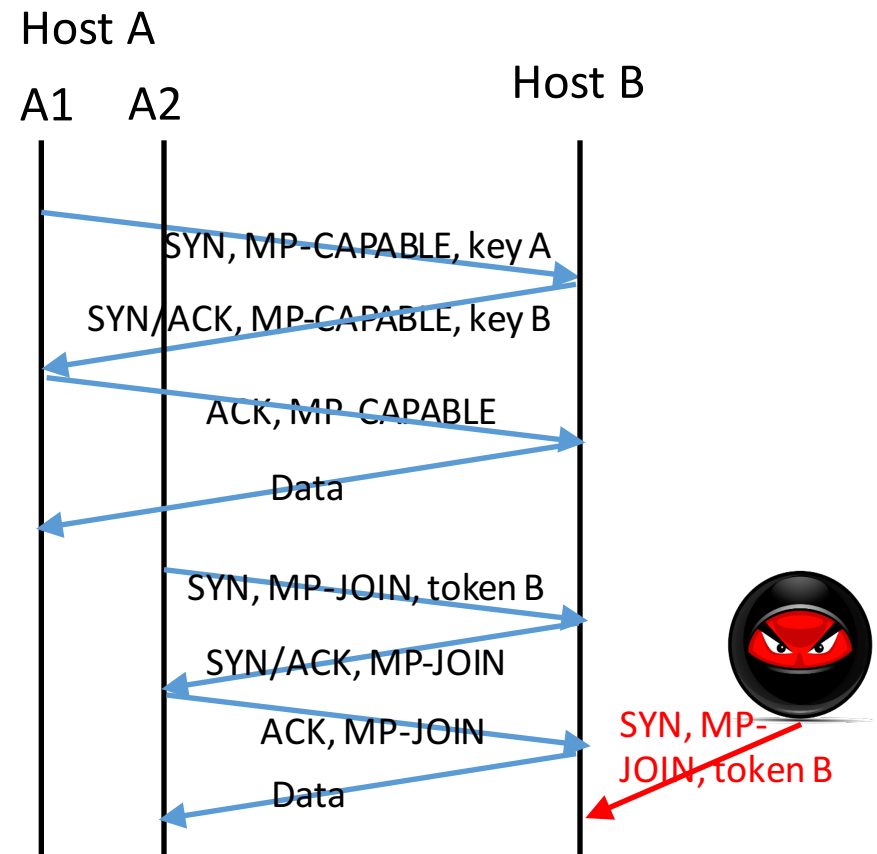
# Adding New Flows: Identification

- TCP flows traditionally identified by <source IP, source port, dest IP, dest port>

- Problem: when adding new subflow to existing connection, don't know the source IP
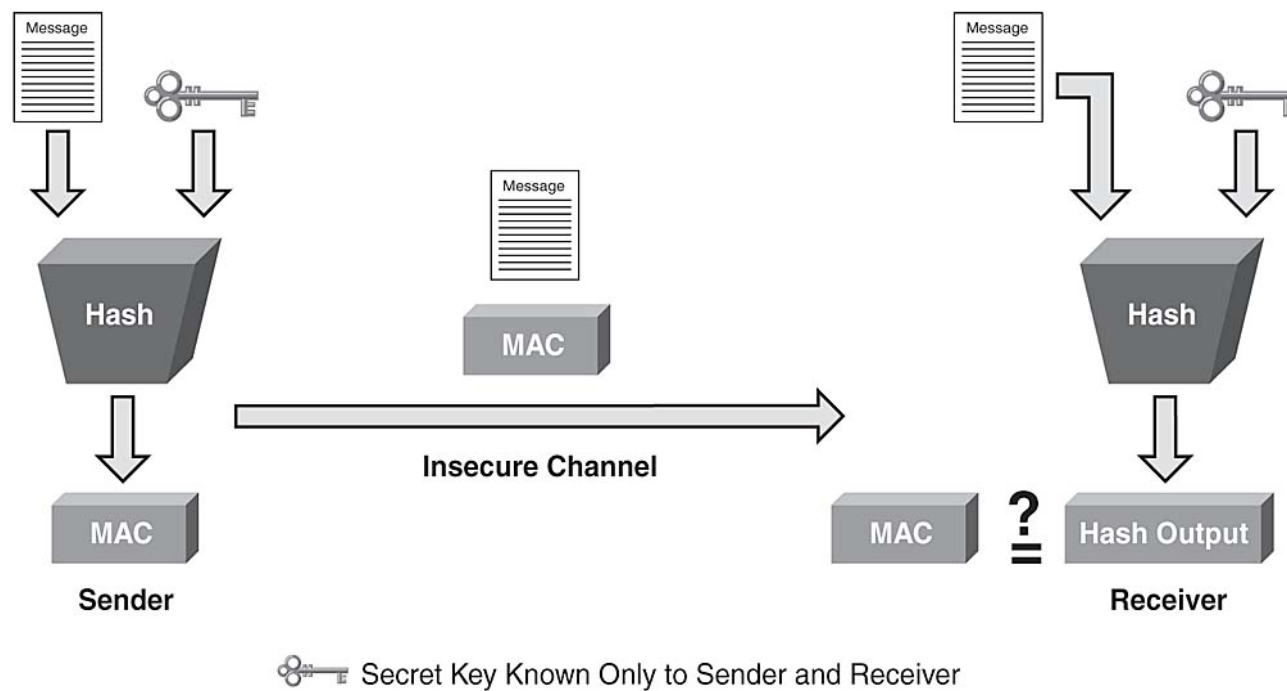  → add a token to identify the connection
  - token = hash(key)

Host A

A1    A2                    Host B

SYN, MP-CAPABLE, key A

SYN/ACK, MP-CAPABLE, key B

ACK, MP-CAPABLE

Data

SYN, MP-JOIN, token B

SYN/ACK, MP-JOIN

ACK, MP-JOIN

Data

# Adding New Flows: Authentication

- **Problem**: attacker could use the same token
  - → authentication using HMAC

Host A

A1    A2                        Host B

SYN, MP-CAPABLE, key A

SYN/ACK, MP-CAPABLE, key B

ACK, MP-CAPABLE

Data

SYN, MP-JOIN, token B

SYN/ACK, MP-JOIN

ACK, MP-JOIN
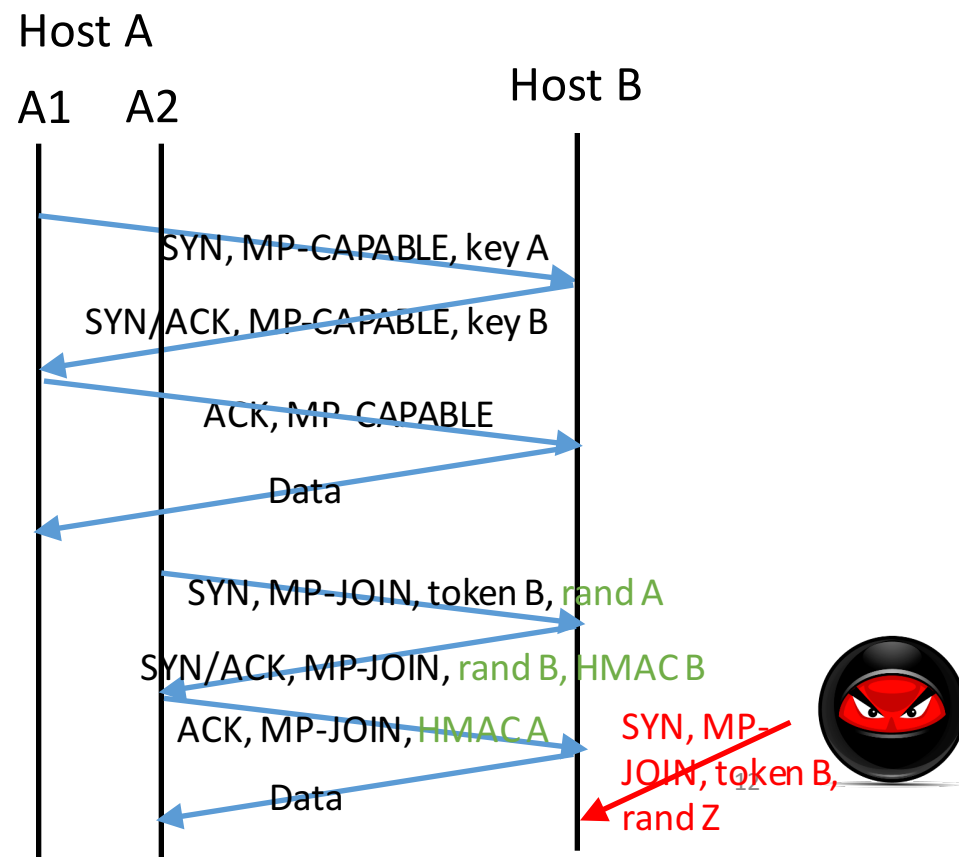
Data

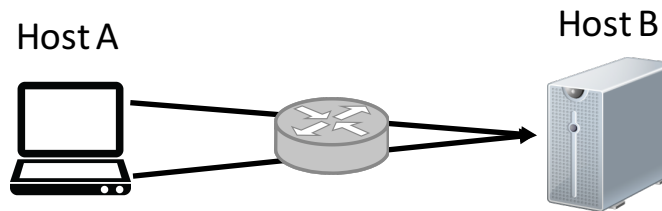SYN, MP-JOIN, token B

# Hash-based Message Authentication Code (HMAC)

# Adding New Flows: Authentication

- **Problem**: attacker could use the same token
  - → authentication using HMAC
  - HMAC = f(key, rand)
  - Attacker gets one change to guess the HMAC, otherwise rand changes

Host A

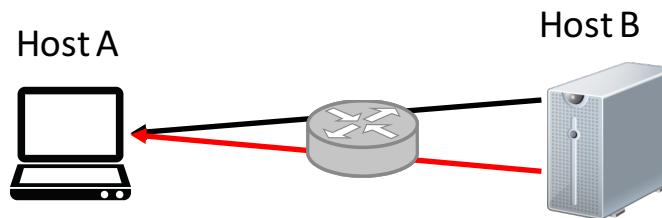A1    A2                              Host B

SYN, MP-CAPABLE, key A

SYN/ACK, MP-CAPABLE, key B

ACK, MP-CAPABLE

Data

SYN, MP-JOIN, token B, rand A

SYN/ACK, MP-JOIN, rand B, HMAC B

ACK, MP-JOIN, HMAC A

Data

SYN, MP-JOIN, token B, rand Z

# Adding New Flows: Addresses

- Implicit

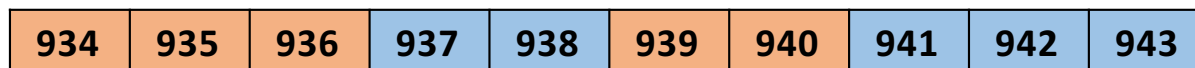Host A                                        Host B



- Explicit
  - Problem: second subflow can't reach client because of NAT
  - Server sends ADD_ADDR option

Host A                          Host B

# Sequence Numbers

- Naïve: Use one sequence of numbers, send subset those numbers on each subflow

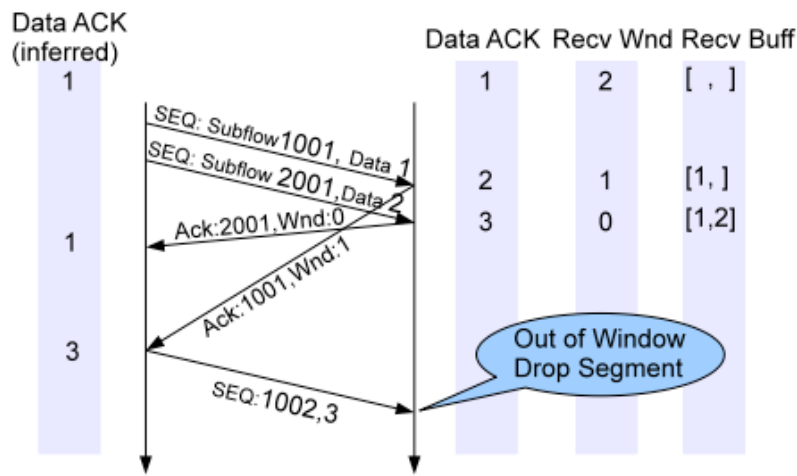| 934 | 935 | 936 | 937 | 938 | 939 | 940 | 941 | 942 | 943 |

Host A1
Host A2

- **Problem**: middleboxes re-initialize sequence numbers
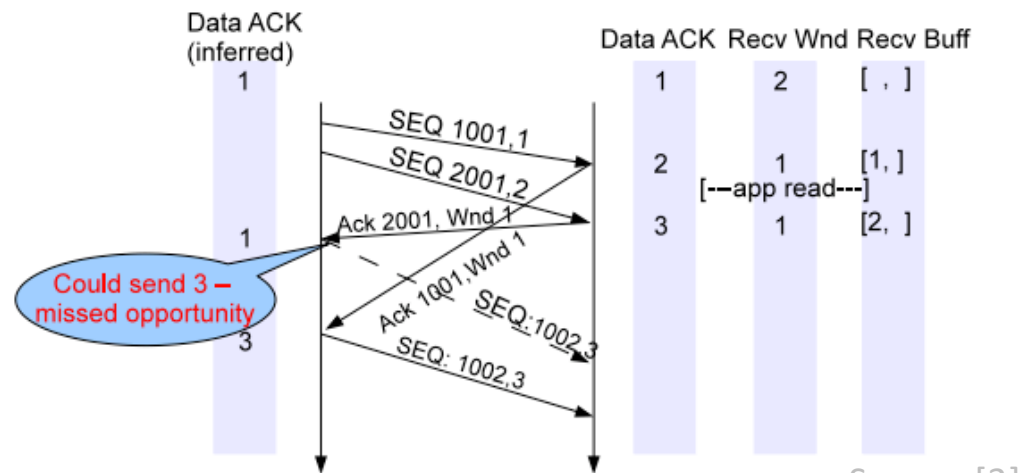- **Problem**: middleboxes don't like gaps in sequence numbers

→ use flow-level sequence numbers along with per-subflow sequence numbers

# Sequence Numbers: ACKs

- Flow-level sequence numbers needed
- Are flow-level ACKs needed? Can we infer them from subflow ACKs?
- Example: receive buffer size 2
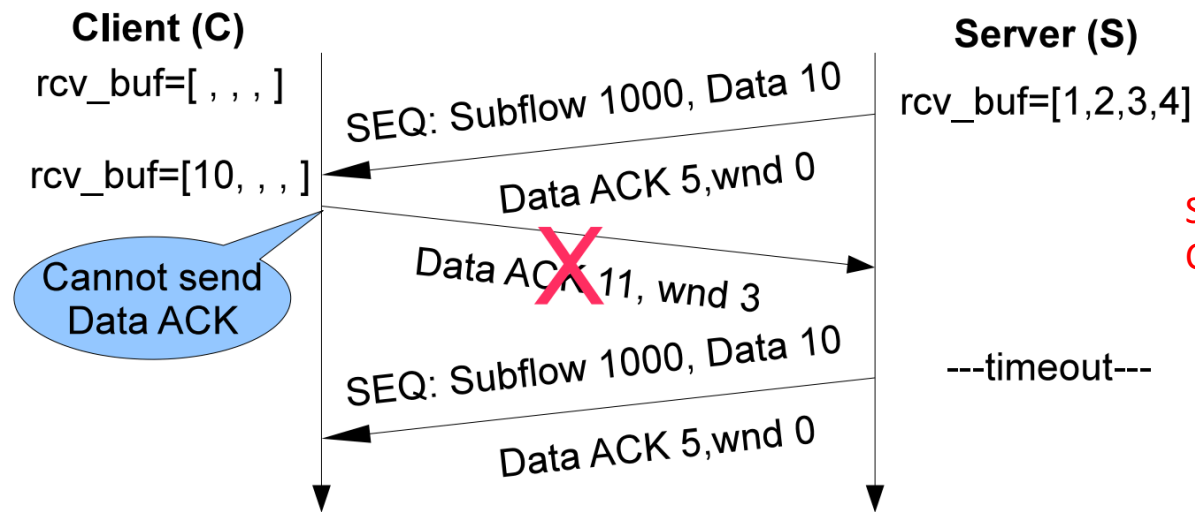


(a) Drops due to incorrect inference

(b) Stalls due to incorrect inference

Source: [3]

# Sequence Numbers: Mapping

- Mapping from subflow sequence number to data sequence number

- Naïve: On each packet, record absolute value of data sequence number

- TCP segmentation offload (TSO)
  - Divide large segments into smaller chunks
  - Performed by NICs to save CPU

- Problem: TSO copies same data sequence number onto multiple packets

  → record exact mapping between subflow and data sequence numbers

# Sequence Numbers: Encoding

- Option 1: Encode in data payload
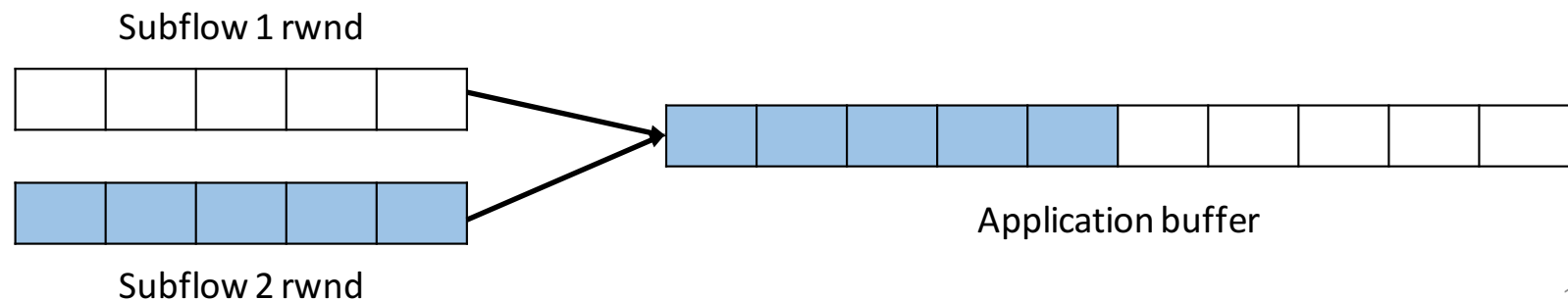- Problem: Data ACKs can get stuck from flow control



→ Encode data sequence numbers and ACKs in TCP options

# Flow Control

- Naïve: Use one receive window for each flow
  → one receive window for each subflow

- Problem: Subflow failure can lead to deadlock

  1. Application waiting for subflow 1's data
  2. Subflow 1 fails, doesn't send data
  3. No space left in subflow 2's rwnd to transmit new data

→ One receive window for the overall flow

Subflow 1 rwnd
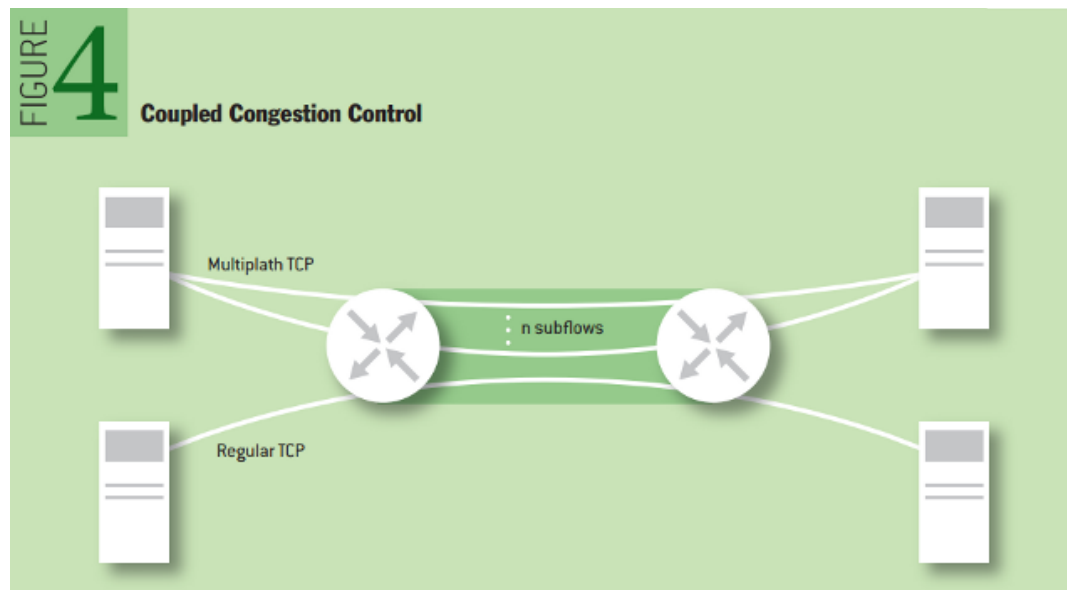
Subflow 2 rwnd

Application buffer

# Retransmissions

- What if data on a subflow times out?
  - Can resend on a different subflow

- Still need to retransmit on the original subflow
  - No holes in subflow sequence numbers for middlebox compatibility
  - Wastes bandwidth

- Protocol not defined by RFC
  - Aggressive: Re-transmit every packet not received on a different subflow
  - Conservative: Re-transmit after fixed number of retries on the original subflow

# Congestion Control

- Naïve: use TCP congestion control separately on each path
- Problem: Not TCP-friendly



FIGURE 4

**Coupled Congestion Control**

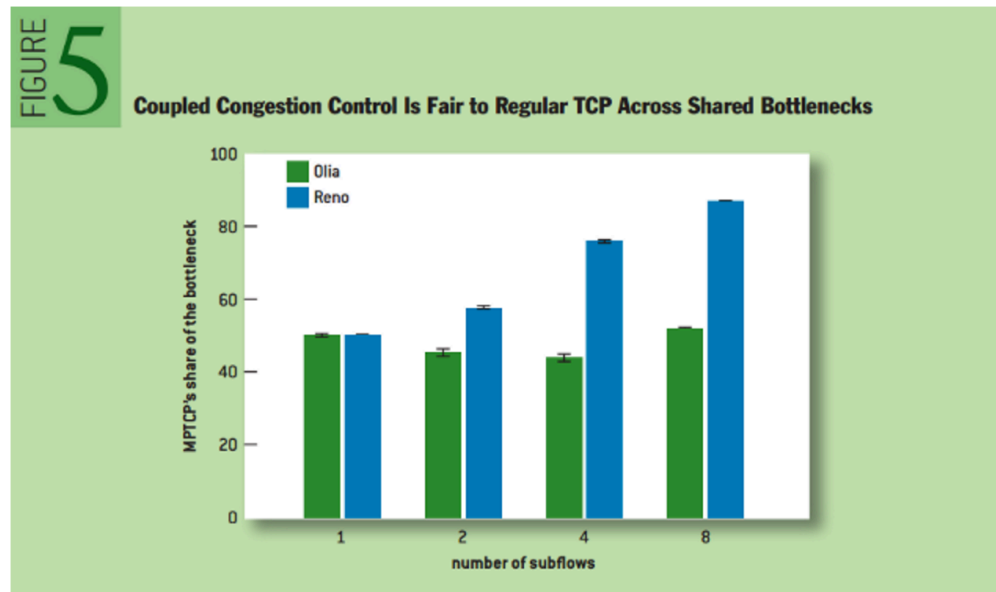Multiplath TCP

n subflows

Regular TCP

For example:
2 clients
Client A has 2 MPTCP subflows
Client B is regular TCP

Client A will receive 2/3 of capacity

# Congestion Control

- Solution: Congestion control coupled across subflows
  - Many algorithms developed



FIGURE 5

**Coupled Congestion Control Is Fair to Regular TCP Across Shared Bottlenecks**
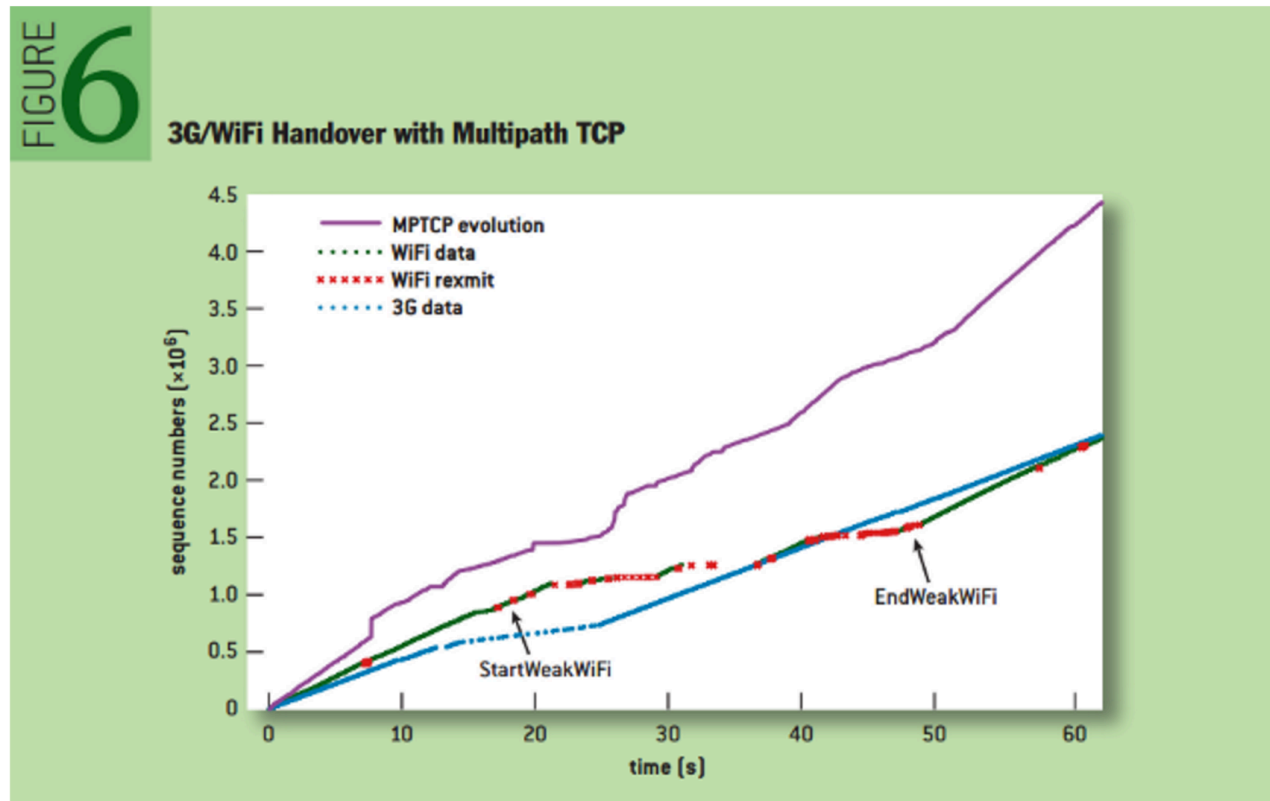
Source: [2]

# Scheduling

- When there is space in both congestion windows, which subflow to transmit on?
  - Round-robin
  - Lowest-RTT first

- ACK-clocked
  - Round-robin: if cwnd has space, send even if out of RR order?
  - Lowest-RTT first: if cwnd has space, send on higher-RTT subflow?

# Practical Example



FIGURE 6

**3G/WiFi Handover with Multipath TCP**

Source: [2]

# Who Uses MPTCP?

- iOS 7 for Siri
  - Primary TCP connection over WiFi
  - Backup TCP connection over cellular data

- Use cases
  - Smartphones with 4G and WiFi for connectivity
  - Data center servers with multiple high-speed links for load balancing

- Linux kernel available

# Paper Discussion

- How computationally expensive is it?

- Is TCP-friendliness too restrictive?

# Sources

1. "Multipath TCP," Christoph Pasch and Olivier Bonaventure, *ACM Queue*, 2014.

2. TCP Extensions for Multipath Operation with Multiple Addresses, RFC 2684.

3. "How Hard Can It Be? Designing and Implementing a Deployable Multipath TCP," Raiciu et al., *NSDI* 2012.