# Decoupling TCP from IP with Multipath TCP

Olivier Bonaventure

http://inl.info.ucl.ac.be

http://perso.uclouvain.be/olivier.bonaventure

July 2013

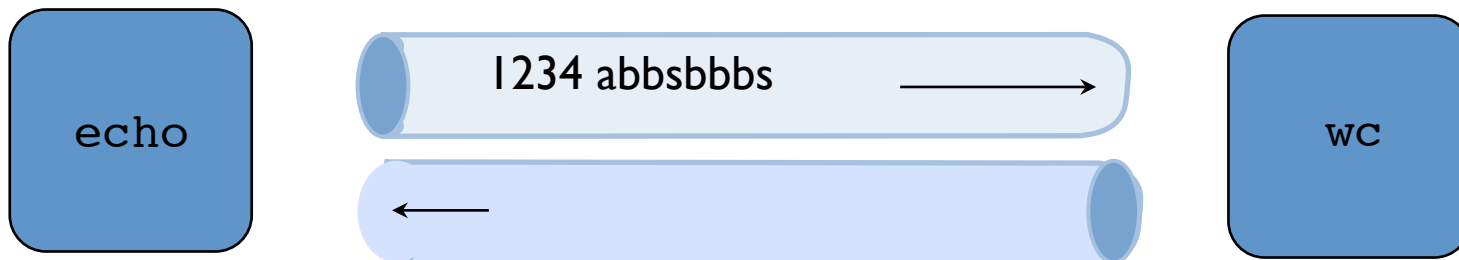# Agenda

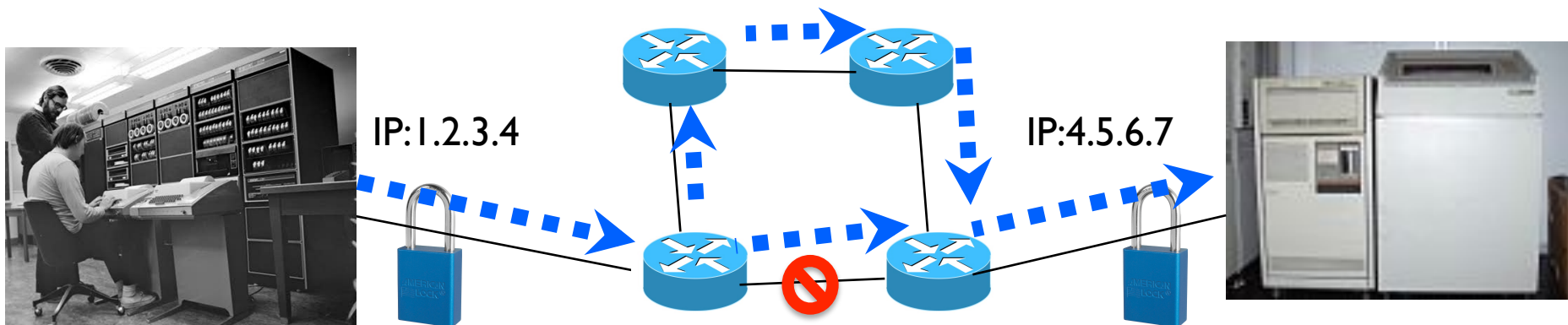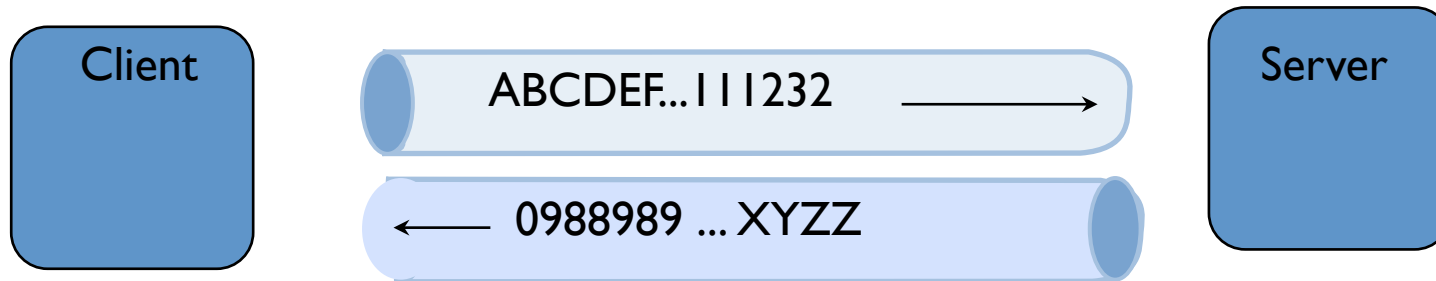→ <span style="color:red">The motivations for Multipath TCP</span>

- The changing Internet

- The Multipath TCP Protocol

- Multipath TCP use cases

# The Unix `pipe` model

```
Terminal — bash — 49×7
Last login: Tue Nov 13 10:07:47 on ttys006
You have new mail.
mbpobo:~ obo$ echo "1234 abbsbbbs" | wc -c
      14
mbpobo:~ obo$ 
```

echo

1234 abbsbbbs →

←

wc

# The TCP bytestream model

Client

ABCDEF...111232 →

← 0988989 ... XYZZ

Server

IP:1.2.3.4

IP:4.5.6.7
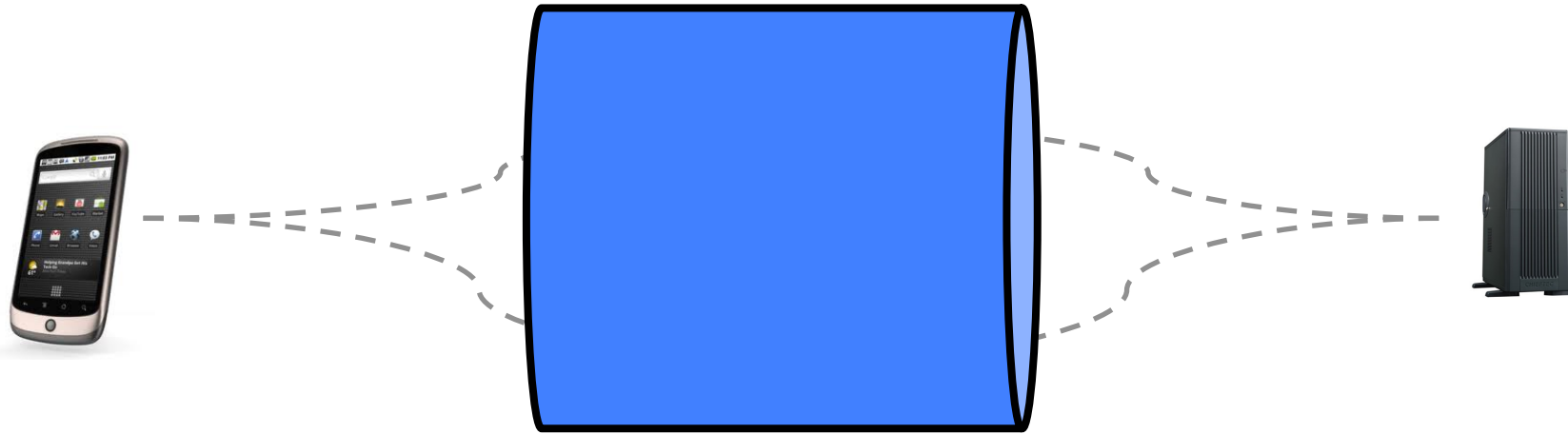
# Endhosts have evolved



**Mobile devices have multiple wireless interfaces**

# User expectations

# What technology provides

3G celltower

IP 1.2.3.4

# What technology provides



IP 1.2.3.4

IP 5.6.7.8

3G celltower

**When IP addresses change TCP connections have to be re-established !**

# Equal Cost Multipath



**ECMP implementation**
Packet arrival :
  $Hash(IP_{src}, IP_{dst}, Prot, Port_{src}, Port_{dst}) \bmod \#oif$

**Packets from one TCP connection follow same path**

**Different connections follow different paths**

G. Detal, Ch. Paasch, S. van der Linden, P. Mérindol, G. Avoine, O. Bonaventure, *Revisiting Flow-Based Load Balancing: Stateless Path Selection in Data Center Networks*, Computer Networks, April 2013
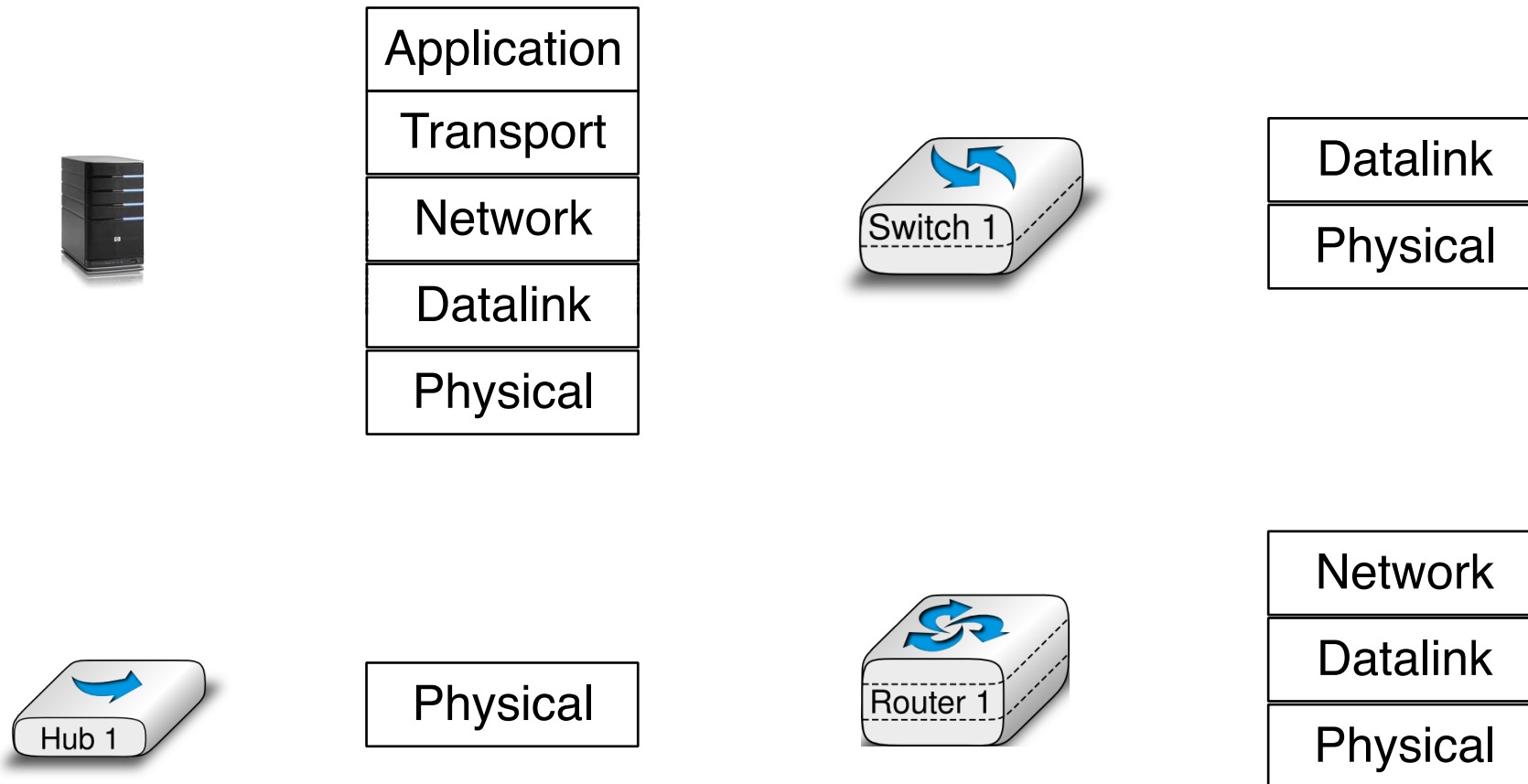
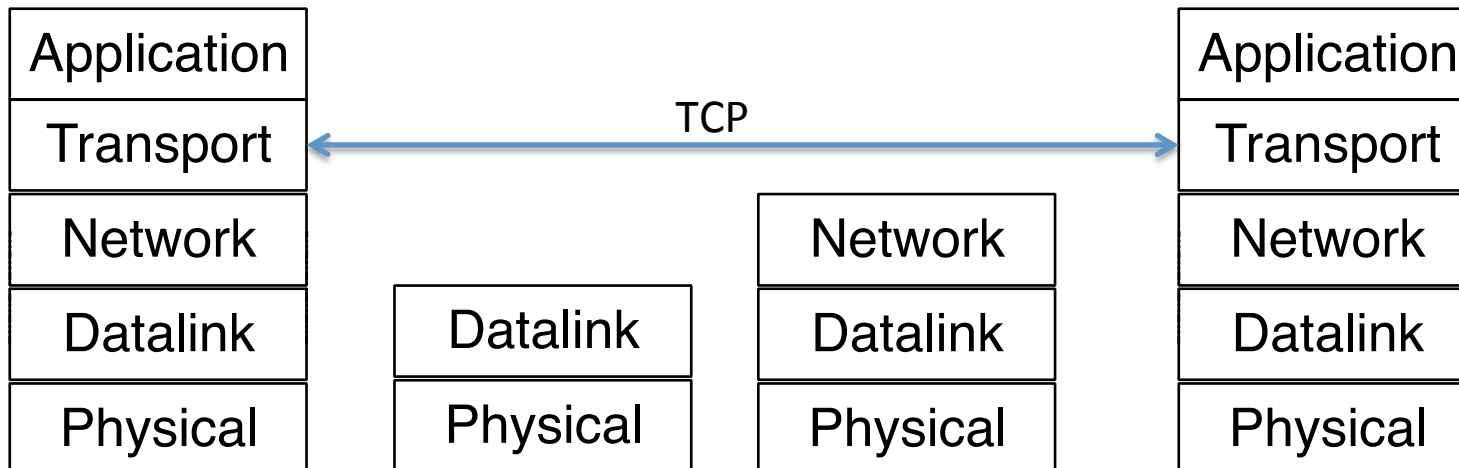# Datacenters

# Agenda

- The motivations for Multipath TCP

➡️ The changing Internet

- The Multipath TCP Protocol

- Multipath TCP use cases

# The Internet architecture
# that we explain to our students

| Application |
|---|
| Transport |
| Network |
| Datalink |
| Physical |

Switch 1

| Datalink |
|---|
| Physical |

Hub 1

| Physical |
|---|

Router 1

| Network |
|---|
| Datalink |
| Physical |

O. Bonaventure, Computer networking : Principles, Protocols and Practice, open ebook, http://inl.info.ucl.ac.be/cnp3

# The end-to-end principle

# In reality



Figure 1: Box plot of middlebox deployments for small (fewer than 1k hosts), medium (1k-10k hosts), large (10k-100k hosts), and very large (more than 100k hosts) enterprise networks. Y-axis is in log scale.
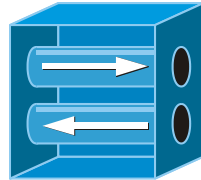
– almost as many middleboxes as routers
– various types of middleboxes are deployed

Sherry, Justine, et al. "*Making middleboxes someone else's problem: Network processing as a cloud service*."
Proceedings of the ACM SIGCOMM 2012 conference. ACM, 2012.
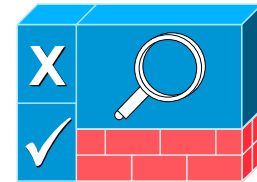
# A middlebox zoo
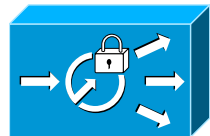
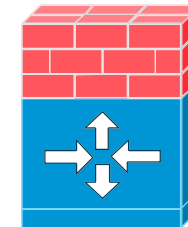Web Security Appliance

VPN Concentrator

SSL Terminator
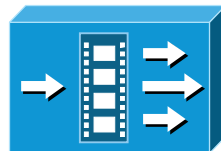
NAC Appliance

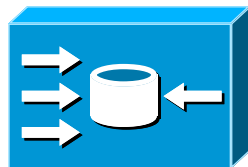ACE XML Gateway

PIX Firewall Right and Left

Cisco IOS Firewall

IP Telephony Router

Streamer

Voice Gateway

Content Engine

NAT

http://www.cisco.com/web/about/ac50/ac47/2.html

# How to model those middleboxes ?

- In the official architecture, they do not exist

- In reality...

# TCP segments processed by a router

**IP**

| Ver | IHL | ToS | Total length | |
|-----|-----|-----|--------------|---|
| Identification | | Flags | Frag. Offset | |
| TTL | Protocol | Checksum | | |
| Source IP address | | | | |
| Destination IP address | | | | |

**TCP**

| Source port | Destination port | |
|-------------|------------------|---|
| Sequence number | | |
| Acknowledgment number | | |
| THL | Reserved | Flags | Window |
| Checksum | Urgent pointer | | |
| Options | | | |
| Payload | | | |

Router 1 →

| Ver | IHL | *ToS* | *Total length* | |
|-----|-----|-----|--------------|---|
| Identification | | *Flags* | *Frag. Offset* | |
| **TTL** | Protocol | **Checksum** | | |
| Source IP address | | | | |
| Destination IP address | | | | |
| Source port | Destination port | |
| Sequence number | | |
| Acknowledgment number | | |
| THL | Reserved | Flags | Window |
| Checksum | Urgent pointer | | |
| Options | | | |
| Payload | | | |

# TCP segments processed by a NAT

# How transparent is the Internet ?

- 25th September 2010 to 30th April 2011

- 142 access networks

- 24 countries

- Sent specific TCP segments from client to a server in Japan

Table 2: Experiment Venues

| Country | Home | Hotspot | Cellular | Univ | Ent | Hosting | Total |
|---|---|---|---|---|---|---|---|
| Australia | 0 | 2 | 0 | 0 | 0 | 1 | 3 |
| Austria | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| Belgium | 4 | 0 | 0 | 1 | 0 | 0 | 5 |
| Canada | 1 | 0 | 1 | 0 | 1 | 0 | 3 |
| Chile | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| China | 0 | 7 | 0 | 0 | 0 | 0 | 7 |
| Czech | 0 | 2 | 0 | 0 | 0 | 0 | 2 |
| Denmark | 0 | 2 | 0 | 0 | 0 | 0 | 2 |
| Finland | 1 | 0 | 0 | 3 | 2 | 0 | 6 |
| Germany | 3 | 1 | 3 | 4 | 1 | 0 | 12 |
| Greece | 2 | 0 | 1 | 0 | 0 | 0 | 3 |
| Indonesia | 0 | 0 | 0 | 3 | 0 | 0 | 3 |
| Ireland | 0 | 0 | 0 | 0 | 0 | 1 | 1 |
| Italy | 1 | 0 | 0 | 0 | 1 | 0 | 2 |
| Japan | 19 | 10 | 7 | 3 | 2 | 0 | 41 |
| Romania | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| Russia | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| Spain | 0 | 1 | 0 | 1 | 0 | 0 | 2 |
| Sweden | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| Switzerland | 2 | 0 | 0 | 0 | 0 | 0 | 2 |
| Thailand | 0 | 0 | 0 | 0 | 2 | 0 | 2 |
| U.K. | 10 | 4 | 4 | 2 | 1 | 1 | 22 |
| U.S. | 3 | 4 | 4 | 0 | 4 | 2 | 17 |
| Vietnam | 1 | 0 | 0 | 0 | 1 | 0 | 2 |
| Total | 49 | 34 | 20 | 17 | 17 | 5 | 142 |

Honda, Michio, et al. "*Is it still possible to extend TCP?*" Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference. ACM, 2011.

# End-to-end transparency today

| Ver | IHL | ToS | Total length |
|---|---|---|---|

Identification ... Offset

Middleboxes don't change the Protocol field, but many discard packets with an unknown Protocol field

Acknowledgment number

| THL | Reserved | Flags | Window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|

Options

Payload

| Ver | IHL | *ToS* | *Total length* |
|---|---|---|---|
| **Identification** | | *Flags* | *Frag. Offset* |
| **TTL** | Protocol | | **Checksum** |

**Source IP address**

**Destination IP address**

| **Source port** | **Destination port** |
|---|---|

*Sequence number*

*Acknowledgment number*

| THL | **Reserved** | **Flags** | **Window** |
|---|---|---|---|
| **Checksum** | | | **Urgent pointer** |

*Options*

*Payload*

# Agenda

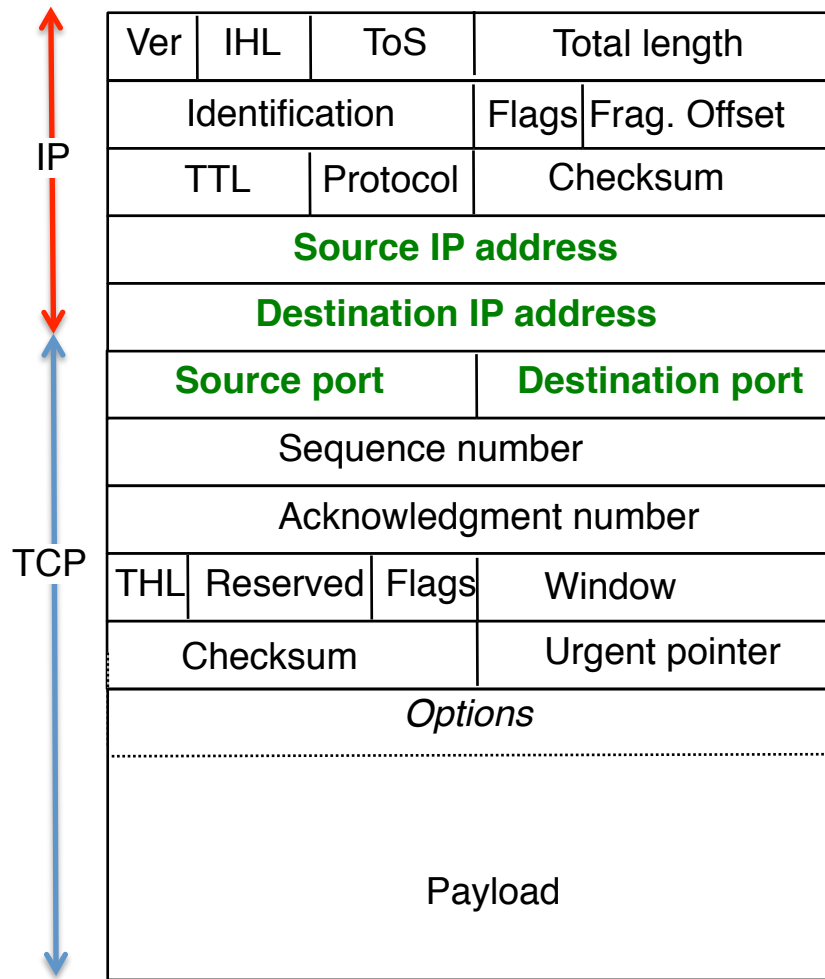- The motivations for Multipath TCP

- The changing Internet

→ The Multipath TCP Protocol

- Multipath TCP use cases

# Design objectives

- Multipath TCP is an *evolution* of TCP

- Design objectives
  - Support unmodified applications
  - Work over today's networks (IPv4 and IPv6)
  - Works in all networks where regular TCP works
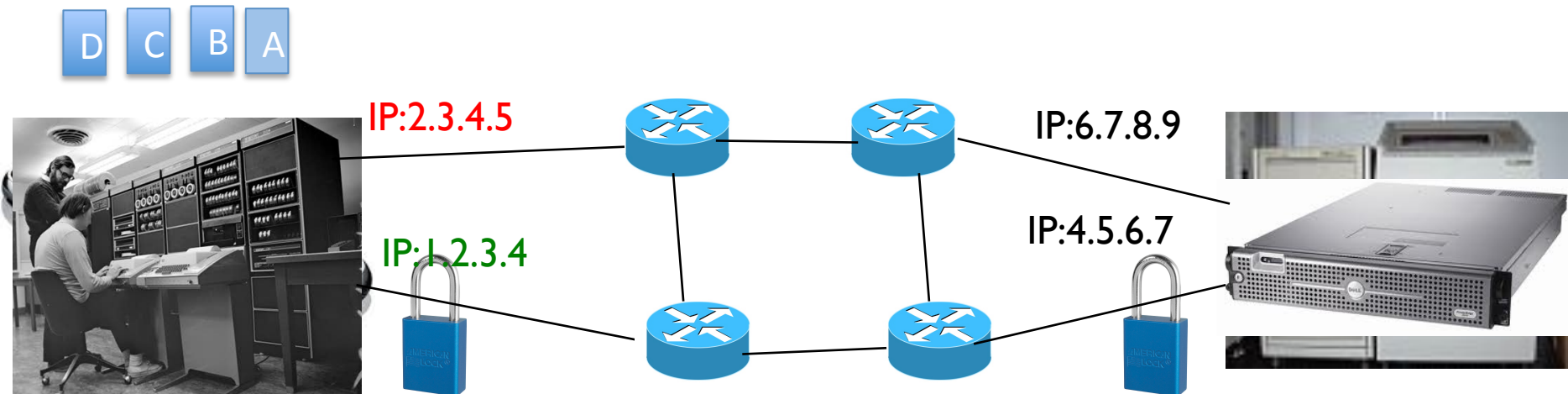
# Identification of a TCP connection

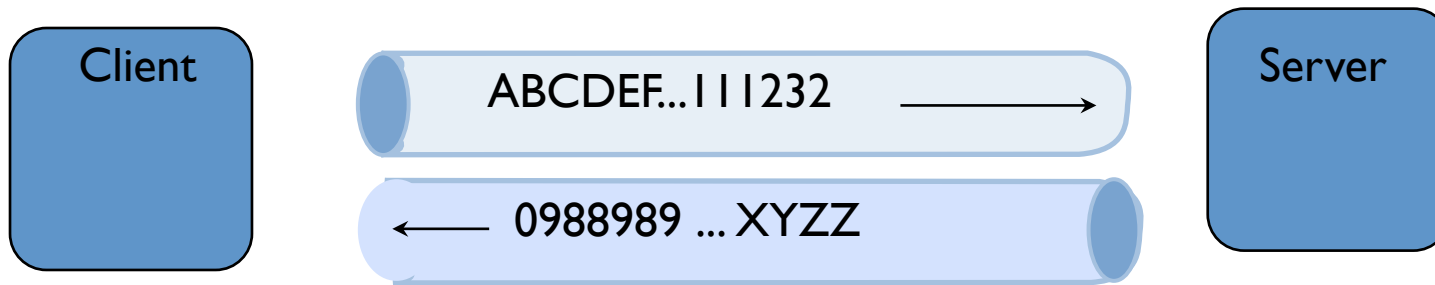| Ver | IHL | ToS | Total length | | |
|---|---|---|---|---|---|
| Identification | | | | Flags | Frag. Offset |
| TTL | | Protocol | Checksum | | |
| **Source IP address** | | | | | |
| **Destination IP address** | | | | | |
| **Source port** | | | **Destination port** | | |
| Sequence number | | | | | |
| Acknowledgment number | | | | | |
| THL | Reserved | Flags | Window | | |
| Checksum | | | Urgent pointer | | |
| *Options* | | | | | |
| Payload | | | | | |

IP

TCP

**Four tuple**

– $IP_{source}$

– $IP_{dest}$

– $Port_{source}$

– $Port_{dest}$

All TCP segments contain the four tuple

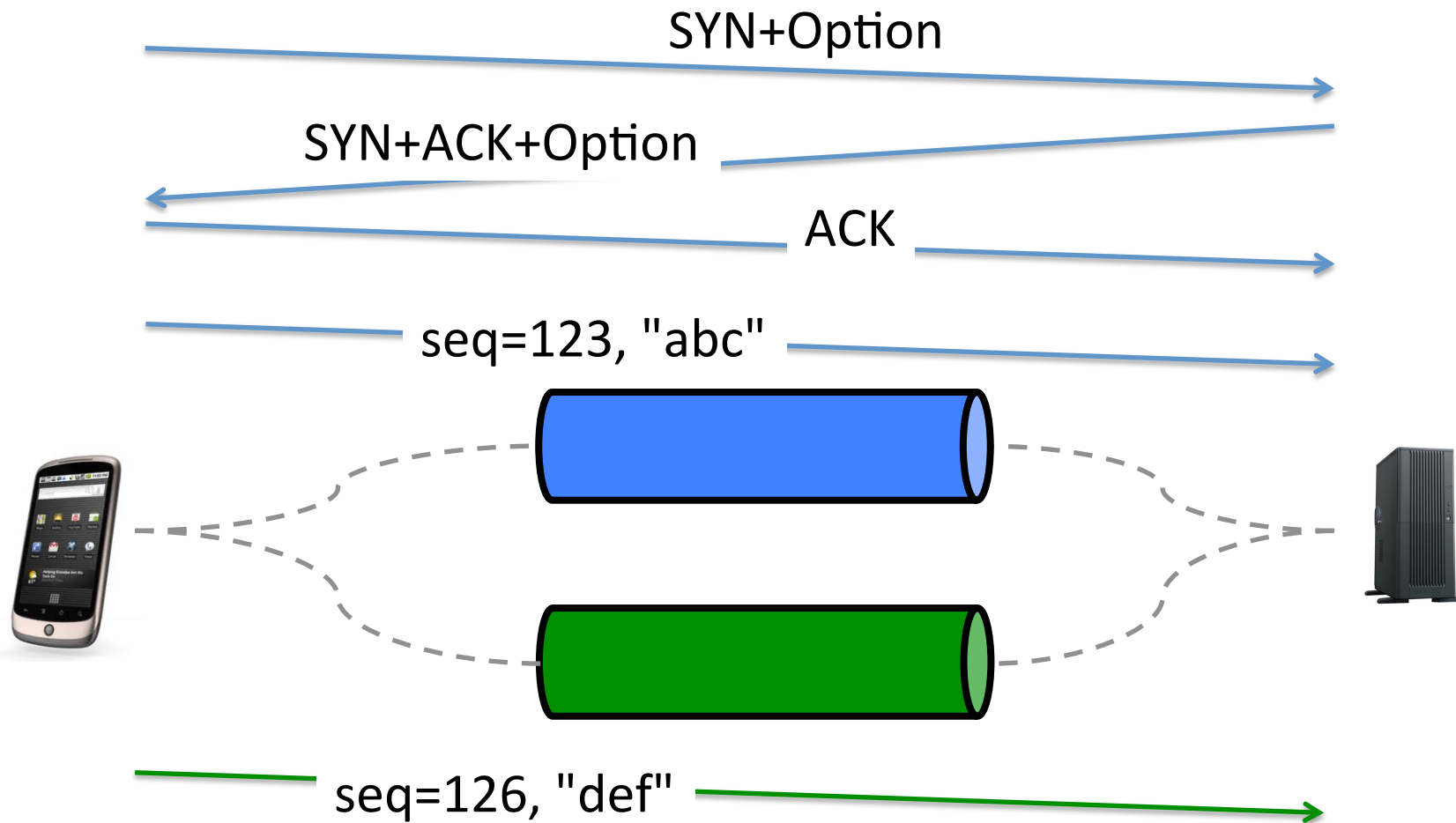# The *new* bytestream model

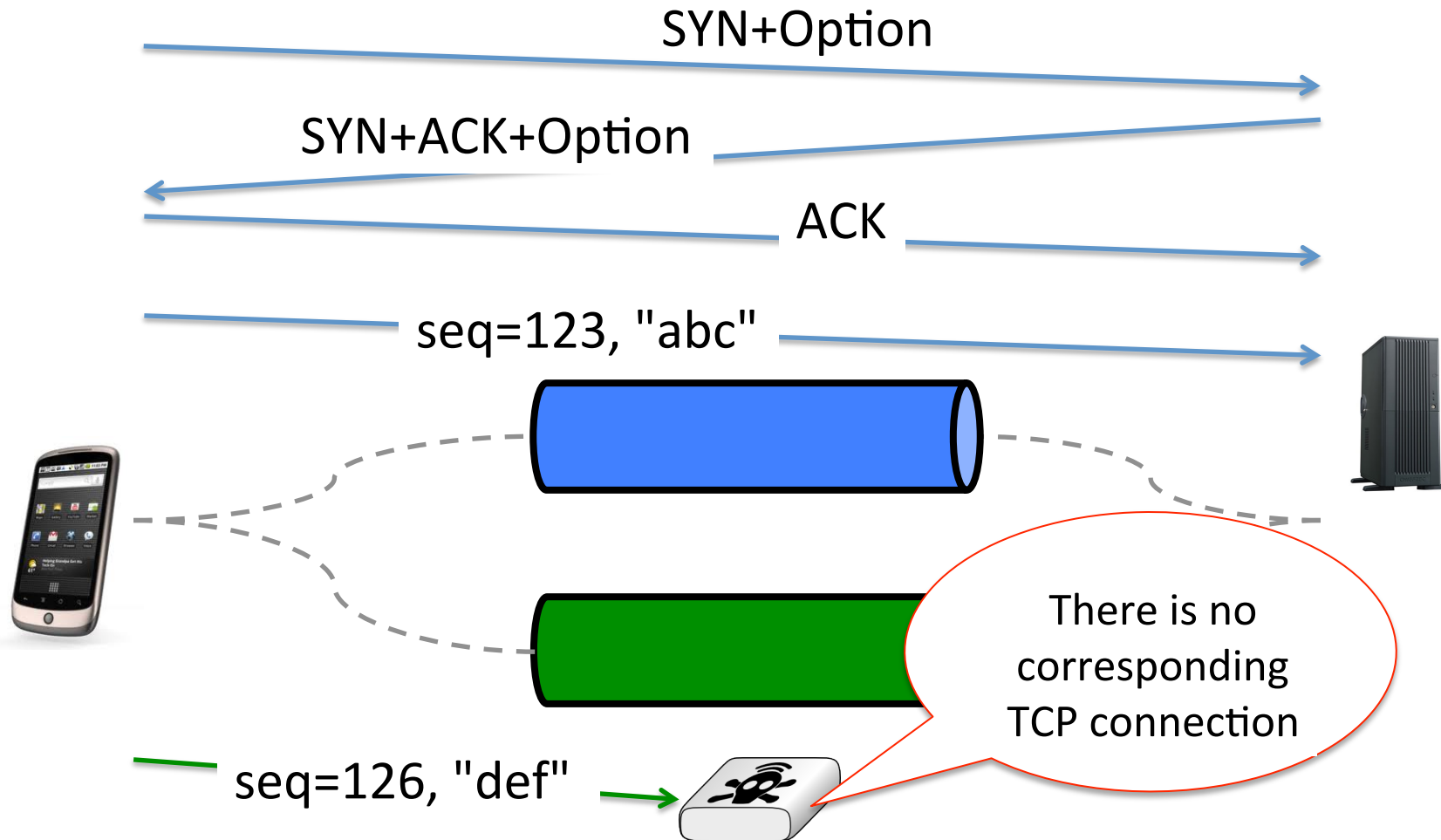# The Multipath TCP protocol

→ <span style="color:red">Control plane</span>
  – How to manage a Multipath TCP connection that uses several paths ?

• Data plane
  – How to transport data ?

• Congestion control
  – How to control congestion over multiple paths ?
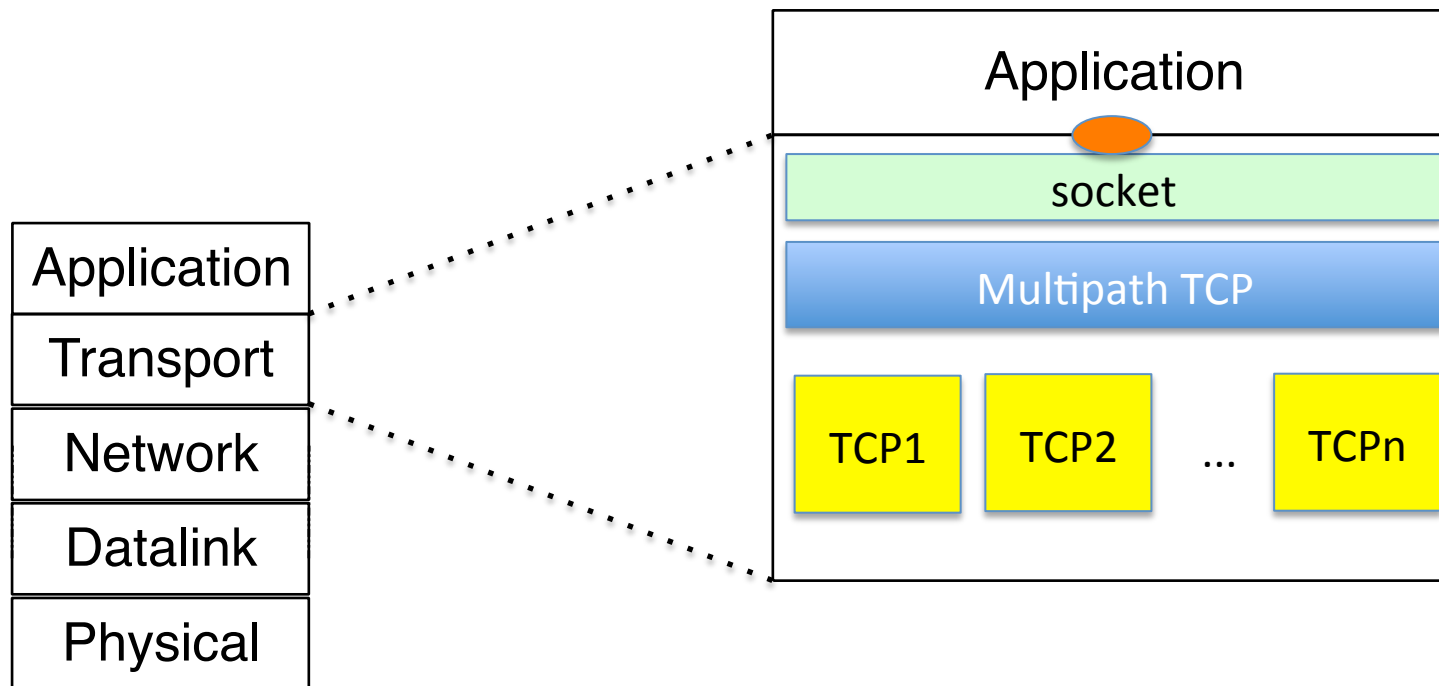
# A naïve Multipath TCP

SYN+Option

SYN+ACK+Option

ACK

seq=123, "abc"

seq=126, "def"

# A naïve Multipath TCP In today's Internet ?

SYN+Option

SYN+ACK+Option

ACK

seq=123, "abc"

seq=126, "def"

There is no corresponding TCP connection

# Design decision

- *A Multipath TCP connection is composed of one or more regular TCP subflows that are combined*

  - Each host maintains state that glues the TCP subflows that compose a Multipath TCP connection together

  - Each TCP subflow is sent over a single path and appears like a **regular TCP** connection along this path
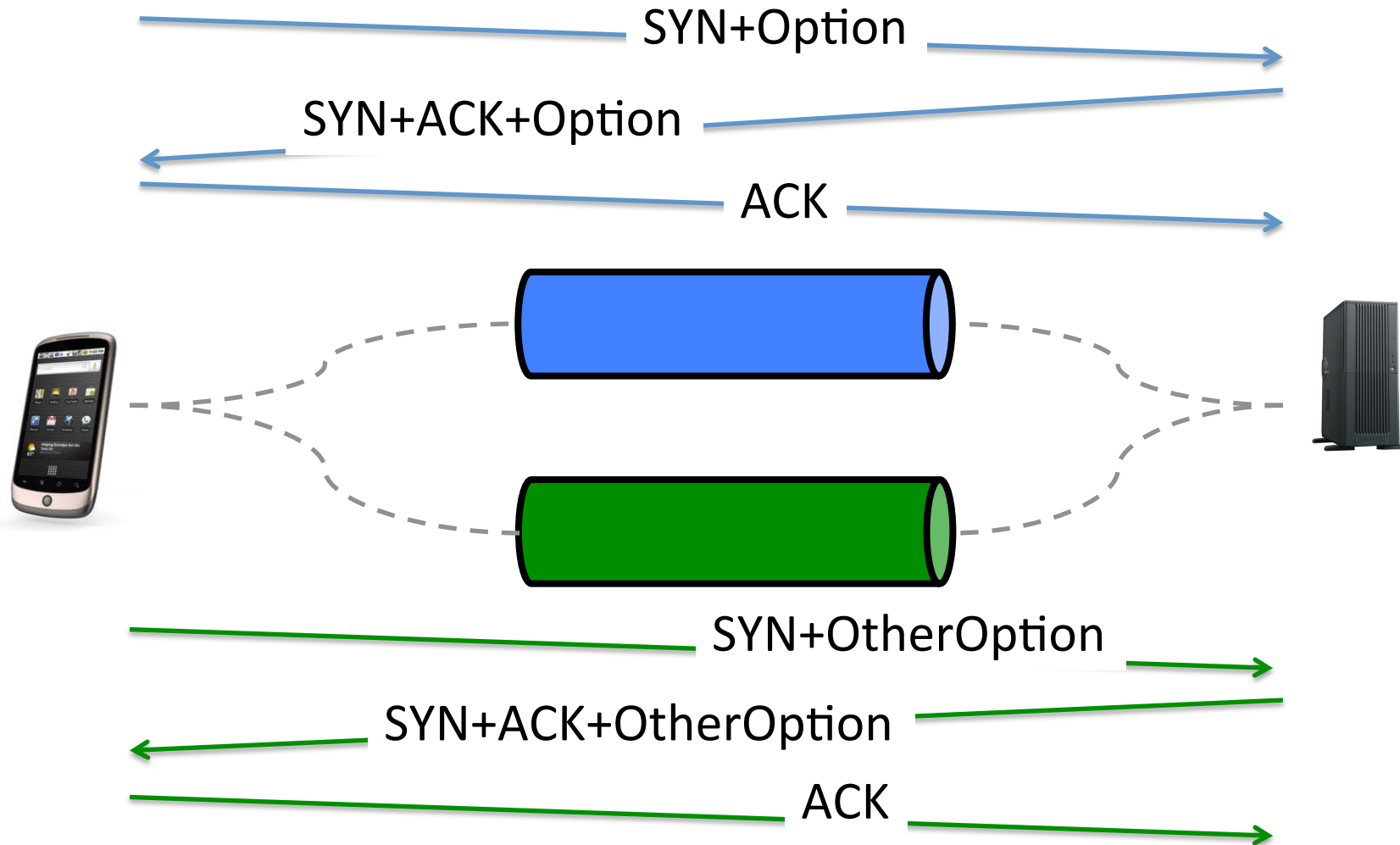
# Multipath TCP and the architecture



A. Ford, C. Raiciu, M. Handley, S. Barre, and J. Iyengar, "Architectural guidelines for multipath TCP development", RFC6182 2011.

# A *regular* TCP connection
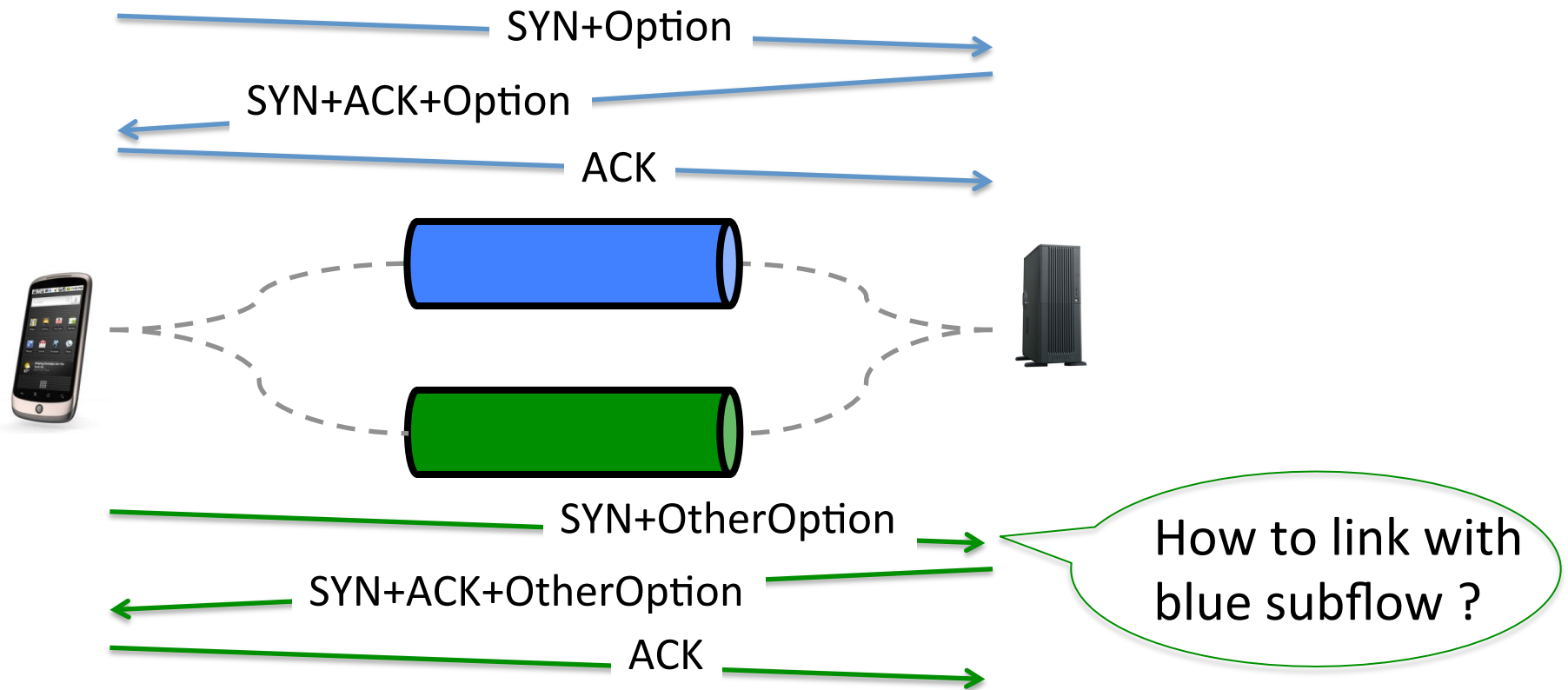
- What is a *regular* TCP connection ?

  - It starts with a three-way handshake
    - SYN segments may contain special options

  - All data segments are sent in sequence
    - There is no gap in the sequence numbers

  - It is terminated by using FIN or RST

# Multipath TCP

SYN+Option

SYN+ACK+Option

ACK

SYN+OtherOption

SYN+ACK+OtherOption

ACK

# How to combine two TCP subflows ?

# How to link TCP subflows ?

# How to link TCP subflows ?

SYN, Port$_{src}$=**1234**,Port$_{dst}$=**80**
+Option[Token=**5678**]

SYN+ACK+Option[Token=**6543**]

ACK

MyToken=5678
YourToken=6543

MyToken=6543
YourToken=5678

SYN, Port$_{src}$=1235,Port$_{dst}$=80
+Option[Token=**6543**]
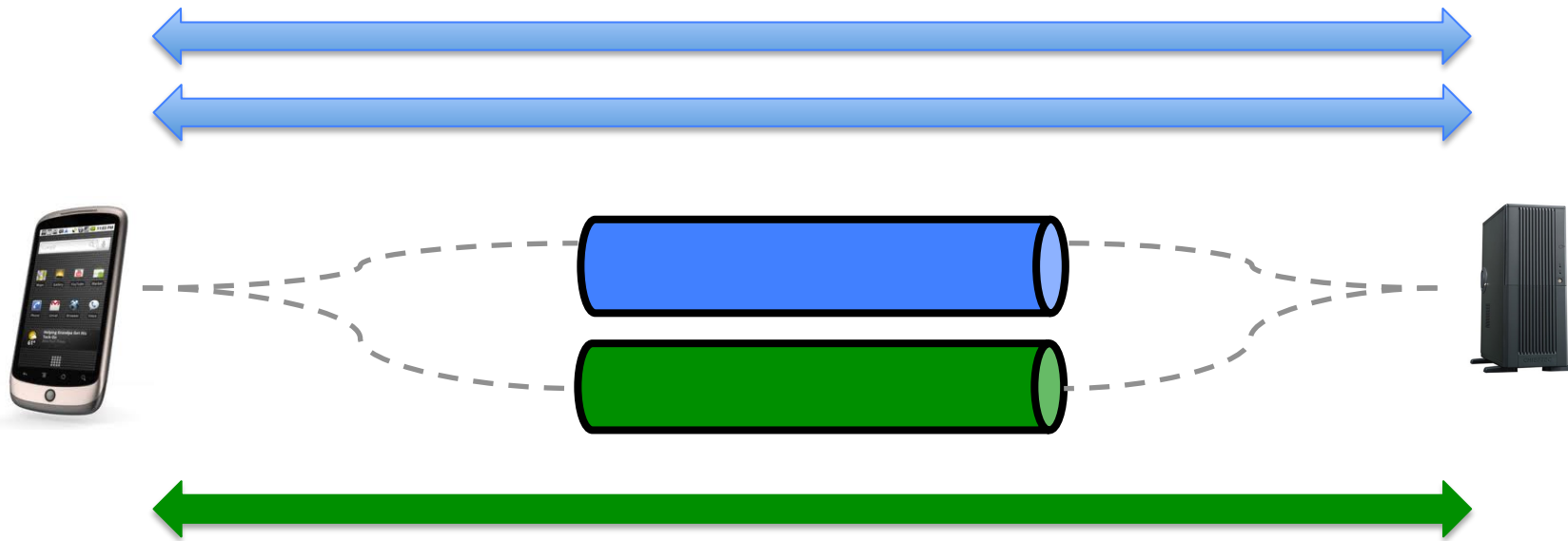
# TCP subflows

- Which subflows can be associated to a Multipath TCP connection ?

  - At least one of the elements of the four-tuple needs to differ between two subflows
    - Local IP address
    - Remote IP address
    - Local port
    - Remote port

# Subflow agility

- Multipath TCP supports
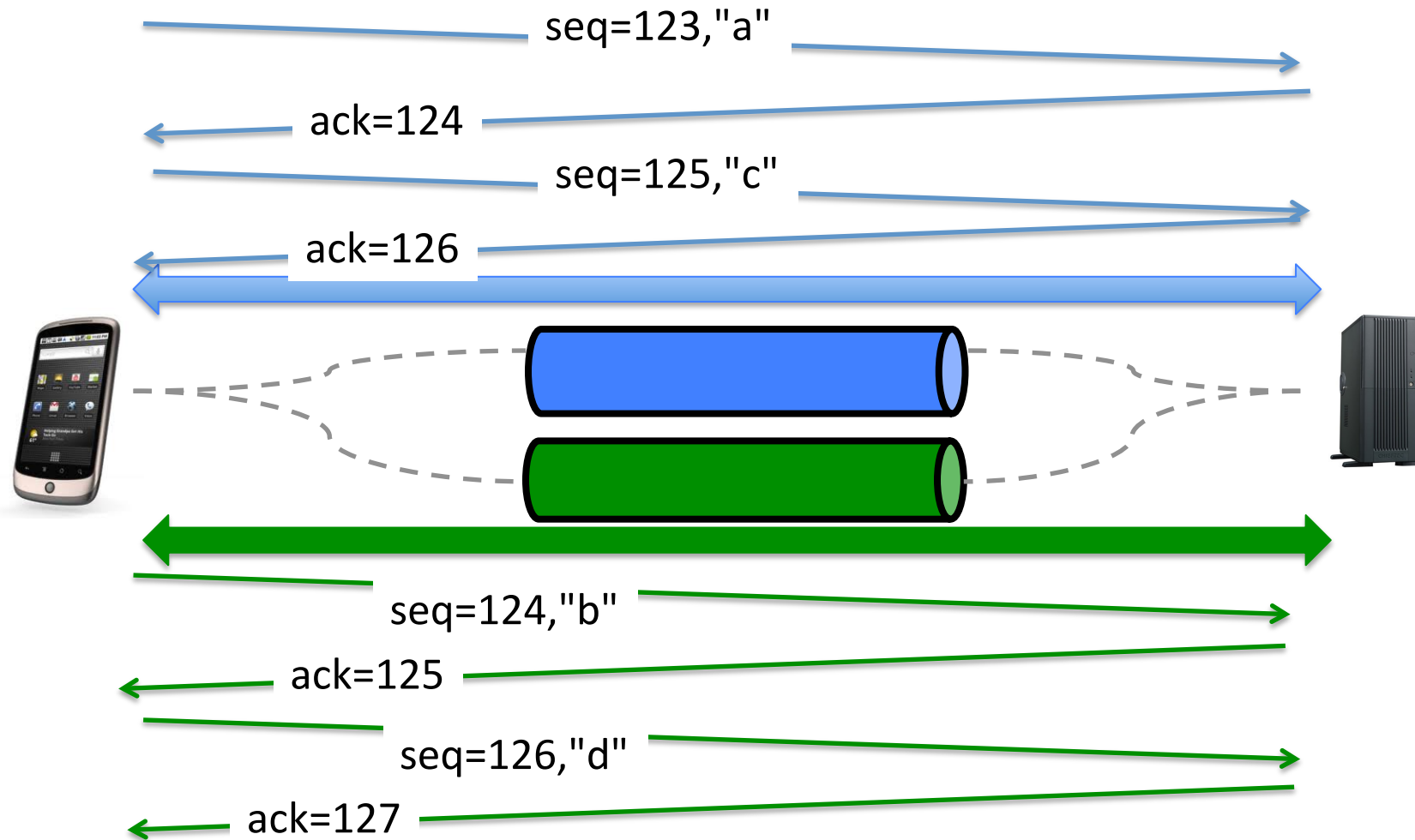  - addition of subflows
  - removal of subflows

# The Multipath TCP protocol

- Control plane
  - How to manage a Multipath TCP connection that uses several paths ?
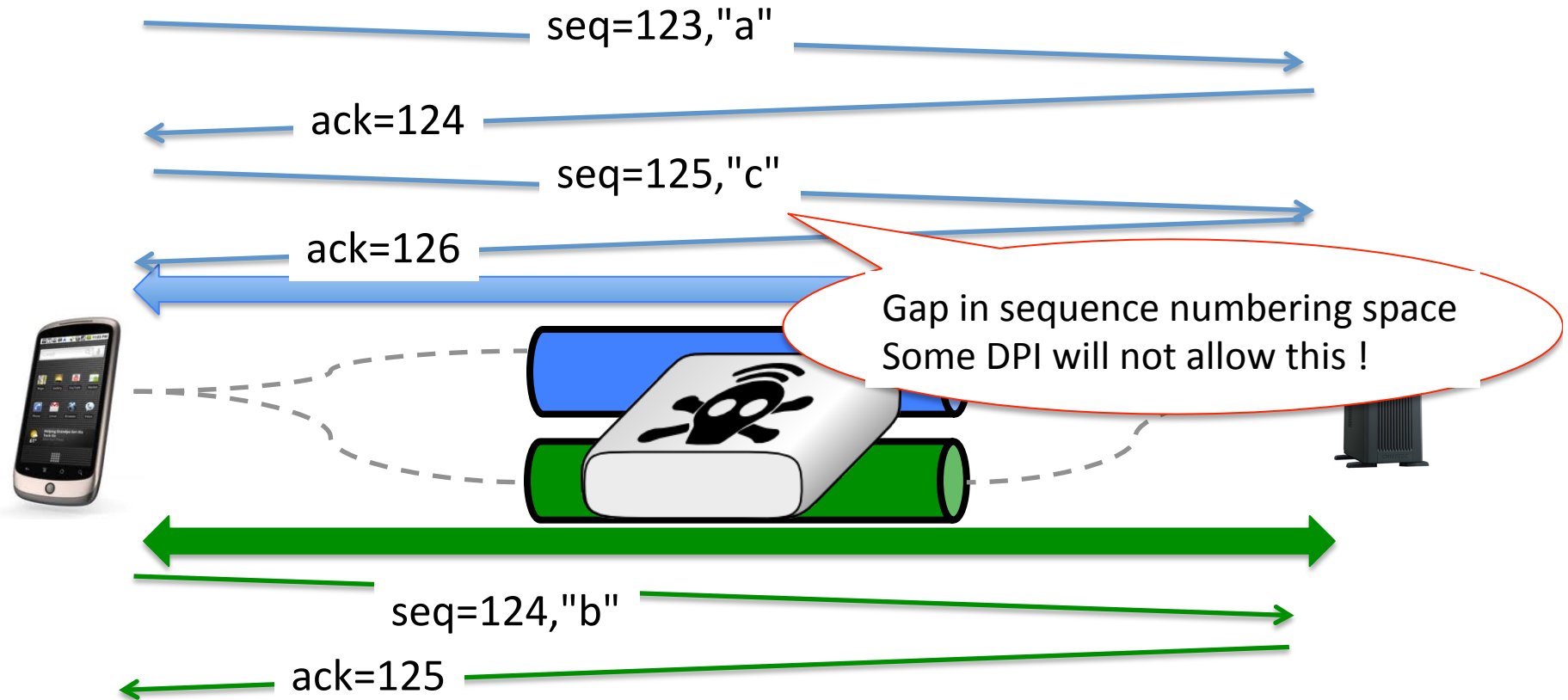
➡ **Data plane**

  - How to transport data ?


- Congestion control
  - How to control congestion over multiple paths ?

# How to transfer data ?

seq=123,"a"

ack=124

seq=125,"c"

ack=126

seq=124,"b"

ack=125

seq=126,"d"

ack=127

# How to transfer data
# in today's Internet ?

seq=123,"a"

ack=124

seq=125,"c"

ack=126

Gap in sequence numbering space
Some DPI will not allow this !

seq=124,"b"

ack=125

# Multipath TCP Data transfer

- Two levels of sequence numbers

# Multipath TCP
# Data transfer

Dseq=0,seq=123,"a"

DAck=1,ack=124

DSeq=2, seq=124,"c"

DAck=3, ack=125

DSeq=1, seq=456,"b"

DAck=2,ack=457
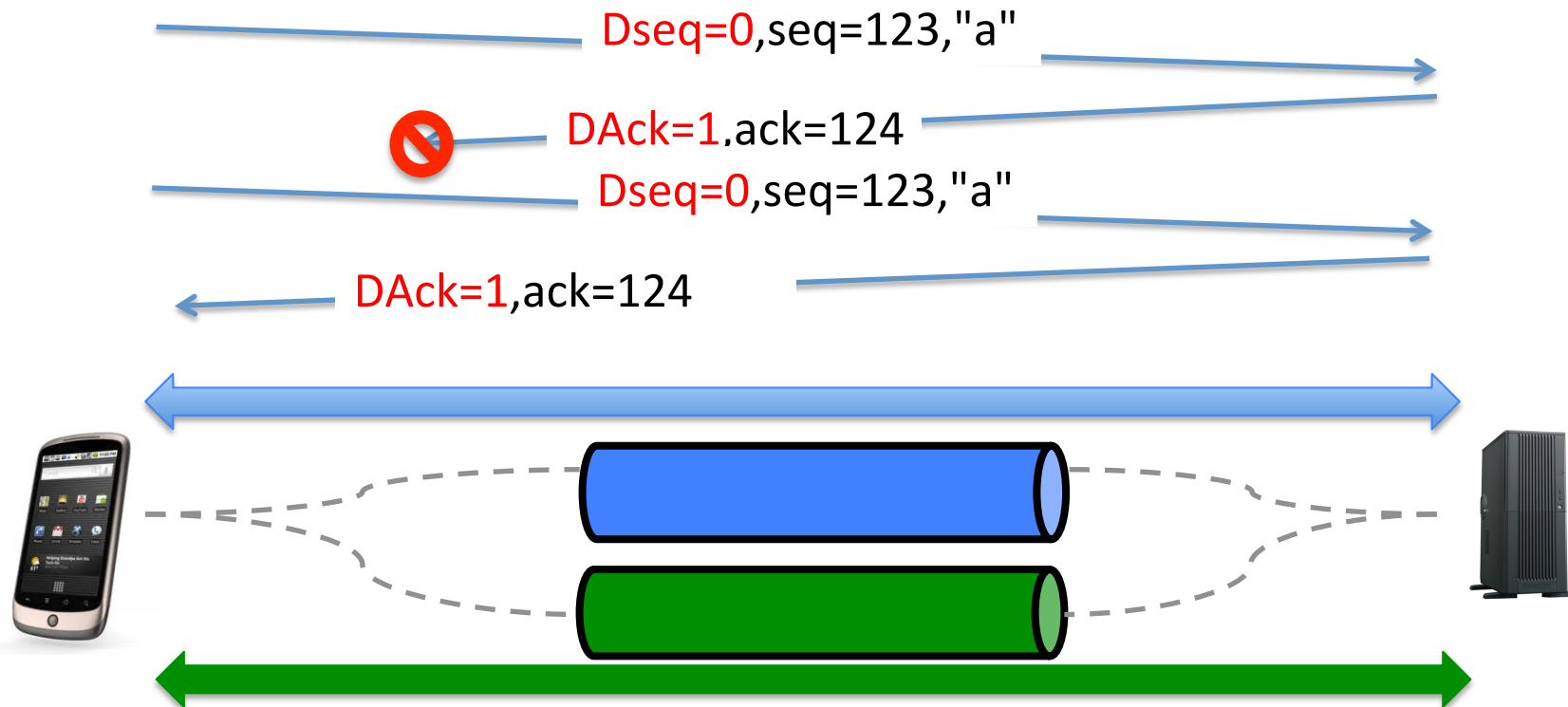
# Multipath TCP
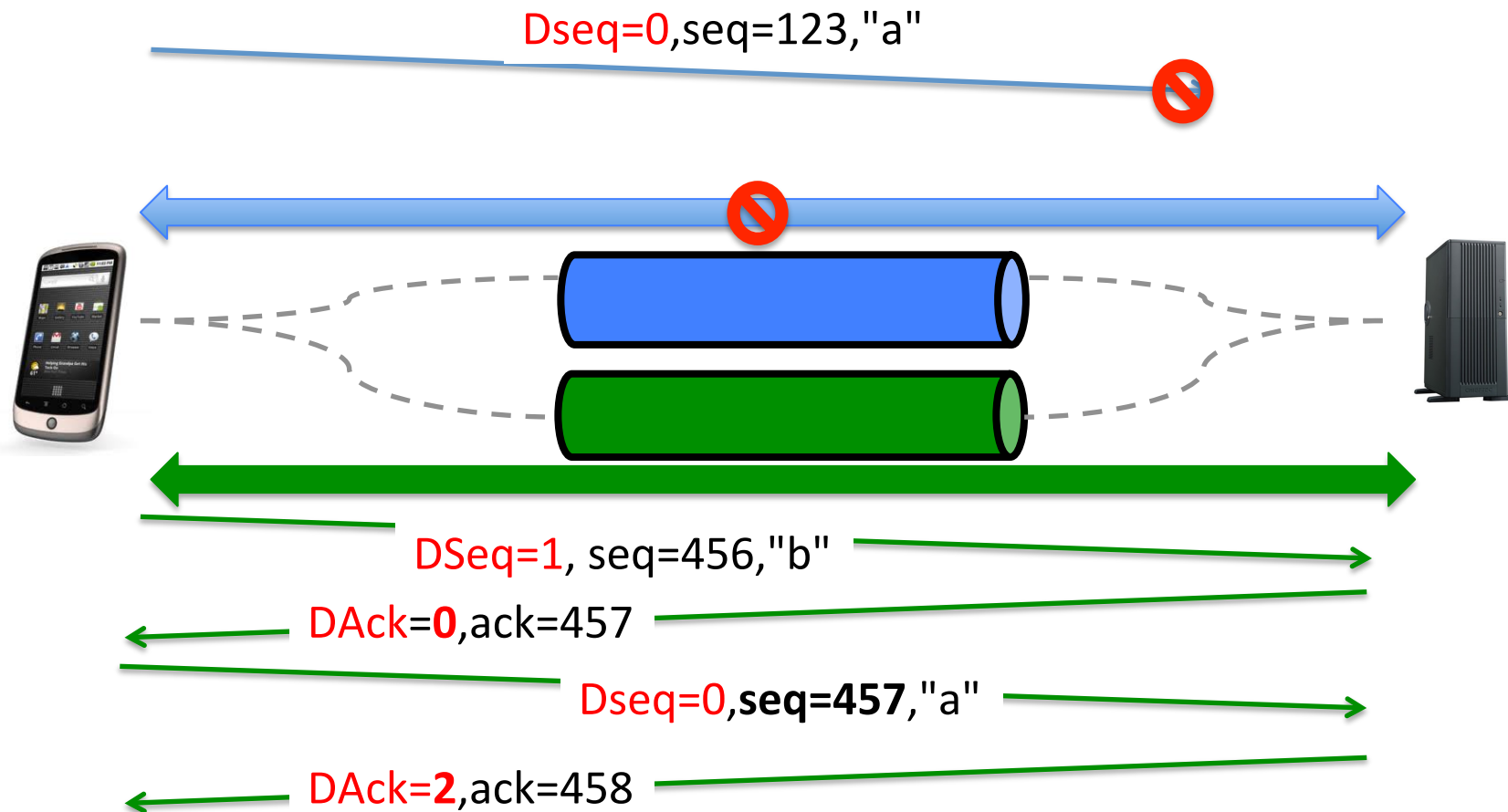# How to deal with losses ?

- Data losses over one TCP subflow
  - Fast retransmit and timeout as in regular TCP



Dseq=0,seq=123,"a"

DAck=1,ack=124

Dseq=0,seq=123,"a"

DAck=1,ack=124

# Multipath TCP

- What happens when a TCP subflow fails ?



Dseq=0,seq=123,"a"

DSeq=1, seq=456,"b"

DAck=**0**,ack=457

Dseq=0,**seq=457**,"a"
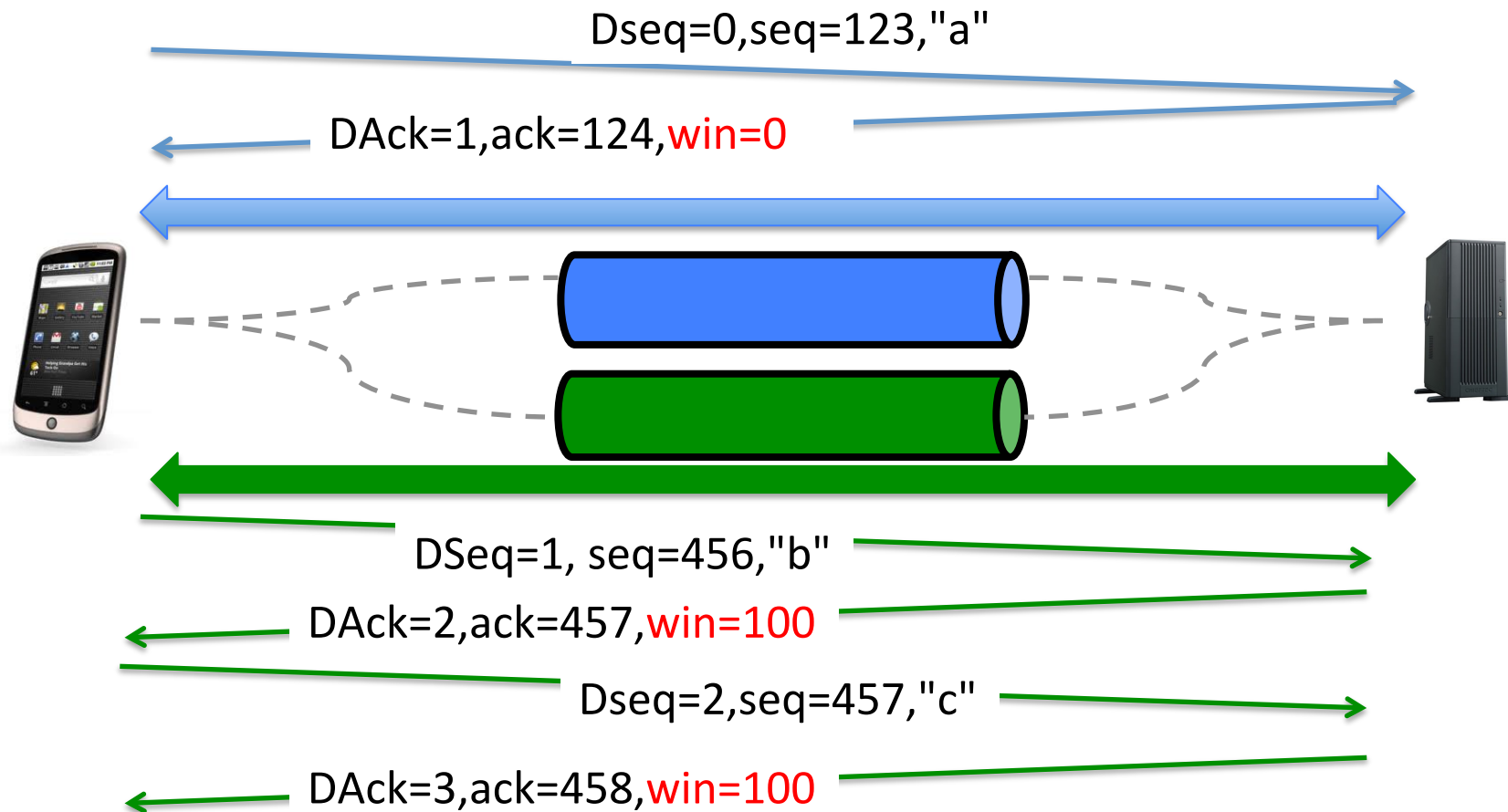
DAck=**2**,ack=458

# Retransmission heuristics

- Heuristics used by current Linux implementation

  - Fast retransmit is performed on the same subflow
    as the original transmission

  - Upon timeout expiration, reevaluate whether the
    segment could be retransmitted over another subflow

  - Upon loss of a subflow, all the unacknowledged data
    are retransmitted on other subflows
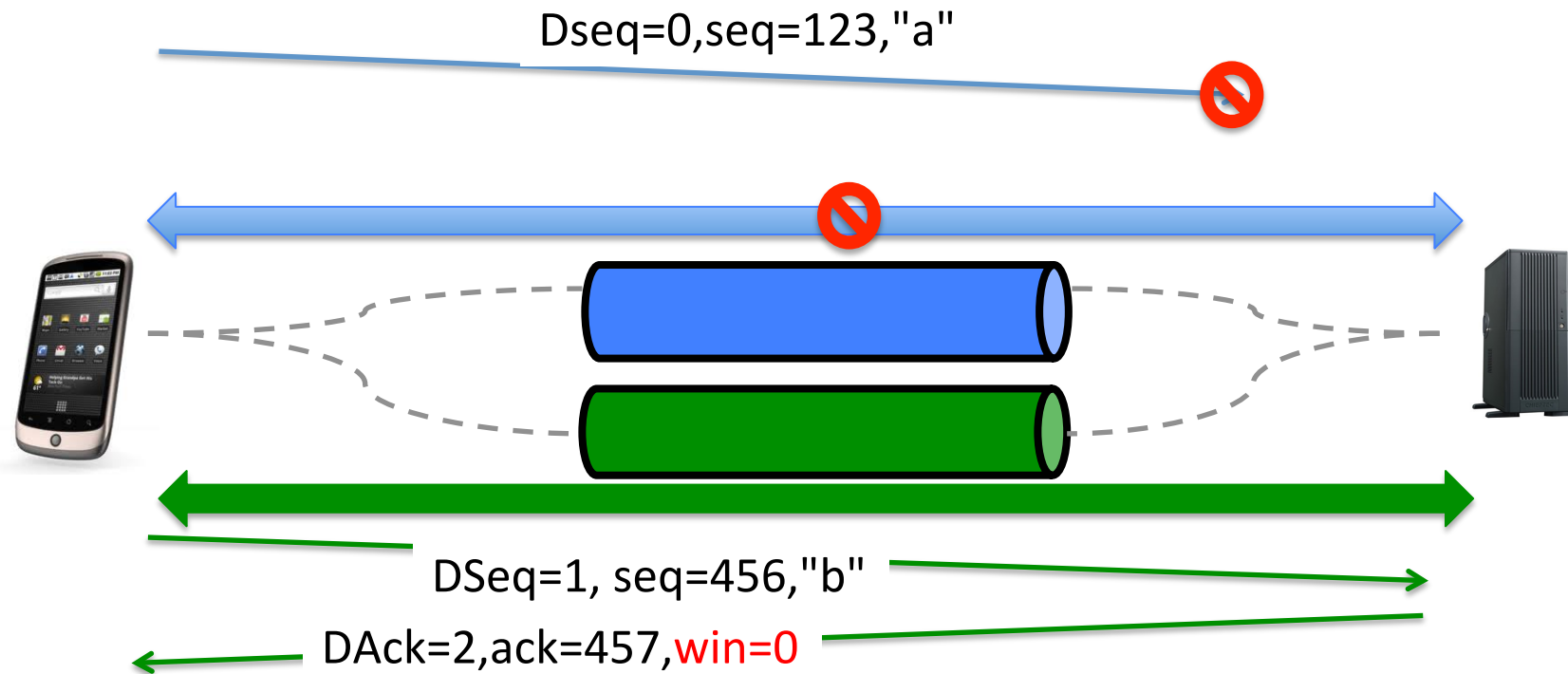
# Flow control

- How should the window-based flow control be performed ?

    – Independant windows on each TCP subflow

    – A single window that is shared among all TCP subflows

# Independant windows

# Independant windows possible problem



Dseq=0,seq=123,"a"

DSeq=1, seq=456,"b"

DAck=2,ack=457,win=0

- Impossible to retransmit, window is already full on green subflow

# A single window shared by all subflows



Dseq=0,seq=123,"a"

DAck=1,ack=124,win=10

DSeq=1, seq=456,"b"

DAck=2,ack=457,win=10

Dseq=2,seq=457,"c"

DAck=3,ack=458,win=10

# A single window shared by all subflows
# Impact of middleboxes



Dseq=0,seq=123,"a"

DAck=1,ack=124,win=100

DSeq=1, seq=456,"b"

DAck=2,ack=457,win=100

DAck=2,ack=457,*win=5*

# Multipath TCP Windows

- Multipath TCP maintains one window per Multipath TCP connection

    - Window is relative to the last acked data (<span style="color:red">Data Ack</span>)
    - Window is shared among all subflows
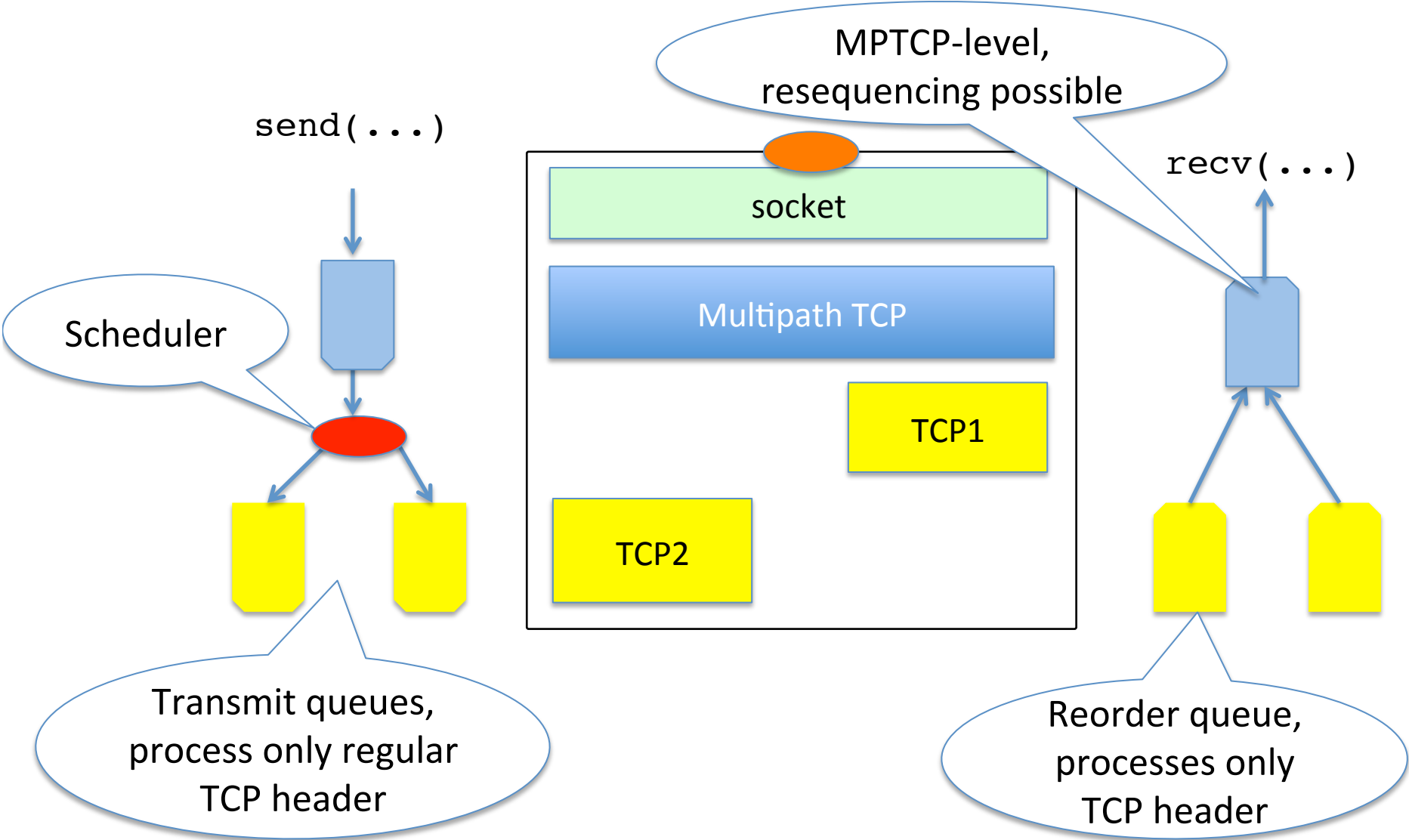        - It's up to the implementation to decide how the window is shared
    - Window is transmitted inside the `window` field of the regular TCP header
    - If middleboxes change `window field`,
        - use largest `window` received at MPTCP-level
        - use received `window` over each subflow to cope with the flow control imposed by the middlebox

# Multipath TCP buffers
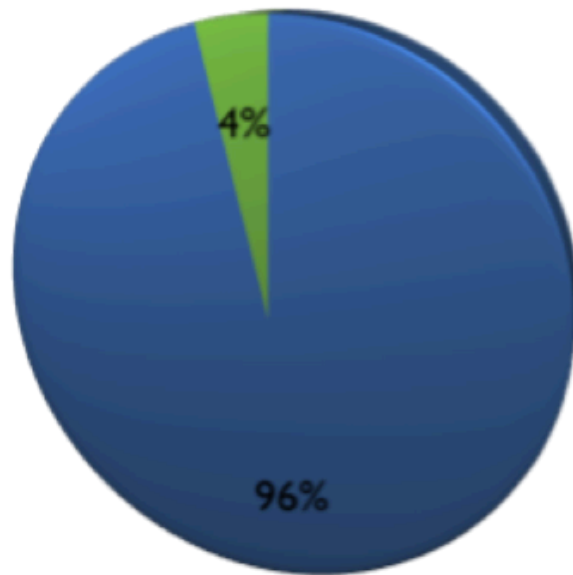
# Sending Multipath TCP information

- How to exchange the Multipath TCP specific information between two hosts ?

- Option 1
  – Use TLVs to encode data and control information inside payload of subflows

- **Option 2**
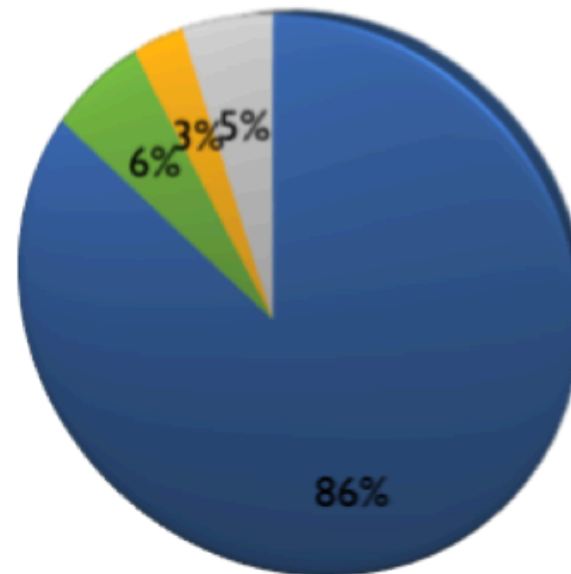  – Use TCP options to encode all Multipath TCP information

Option 1 : Michael Scharf, Thomas-Rolf Banniza , *MCTCP: A Multipath Transport Shim Layer*, GLOBECOM 2011

# Is it safe to use TCP options ?

- Known option (TS) in Data segments

Data segments, port 34443

4%

96%

- Pass
- Modify
- Remove

Data segments, port 80

6% 3% 5%

86%

- Pass
- Modify
- Remove
- Error

Honda, Michio, et al. "Is it still possible to extend TCP?." Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference. ACM, 2011.

© O. Bonaventure, 2011

# Is it safe to use TCP options ?

- Unknown option in Data segments

Data segments, port 34443

4%

96%

Pass    Remove
Modify

Data segments, port 80

9%    5%

86%

Pass    Remove
Modify    Error

Honda, Michio, et al. "Is it still possible to extend TCP?." Proceedings of the 2011 ACM
SIGCOMM conference on Internet measurement conference. ACM, 2011.

© O. Bonaventure, 2011

# Multipath TCP options

- TCP option format

| Kind | Length | Option-specific data |
|------|--------|----------------------|

- Initial design
  - One option kind for each purpose
    (e.g. Data Sequence number)

- Final design
  - A single variable-length Multipath TCP option

# Multipath TCP option

- A single option type
  - to minimise the risk of having one option accepted by middleboxes in SYN segments and rejected in segments carrying data

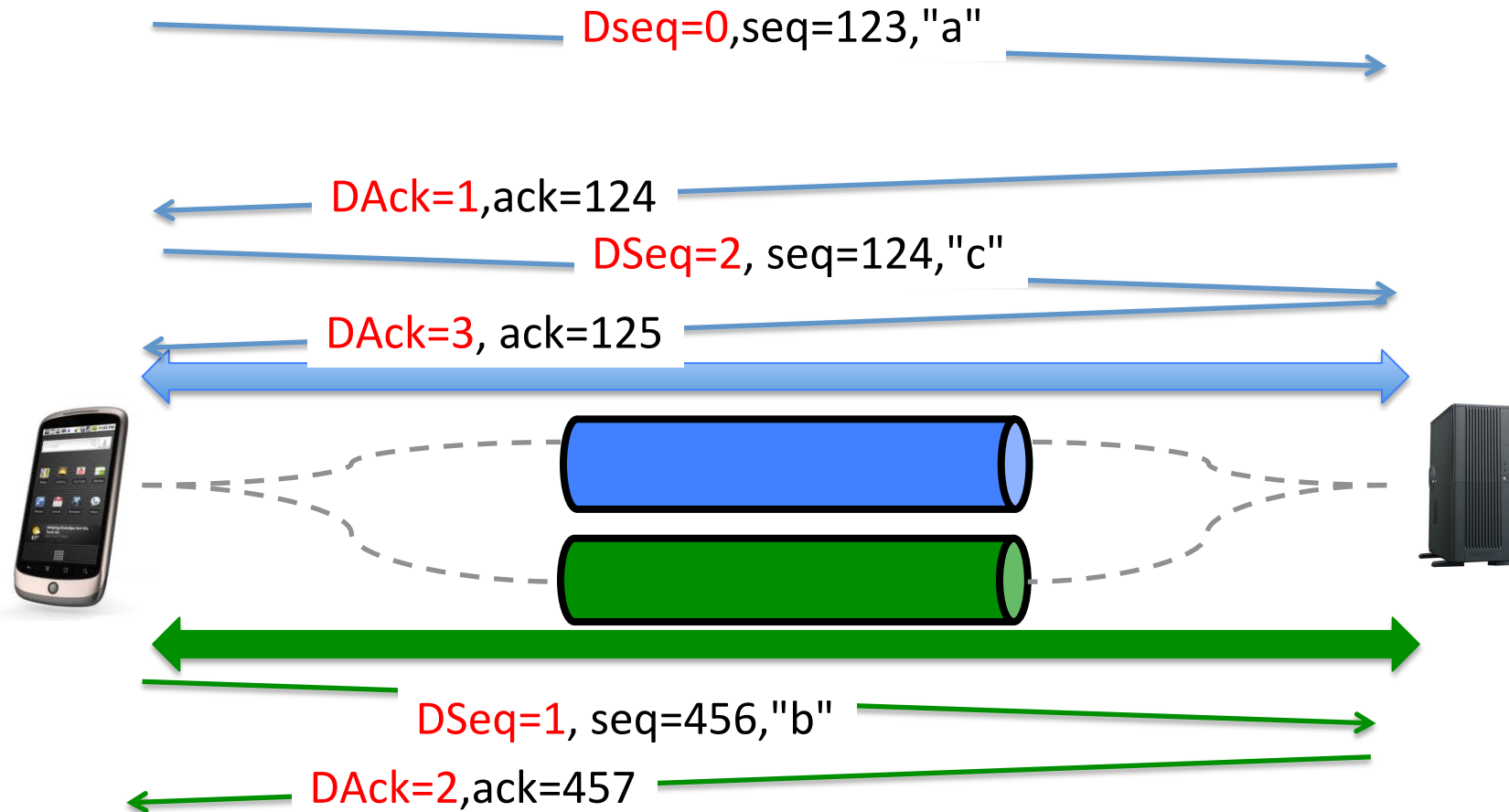| Kind | Length | Subtype | |
|------|--------|---------|---|
| Subtype specific data (variable length) | | | |

# Data sequence numbers and TCP segments

- How to transport Data sequence numbers ?
  - Same solution as for TCP
    - Data sequence number in TCP option is the Data sequence number of the first byte of the segment
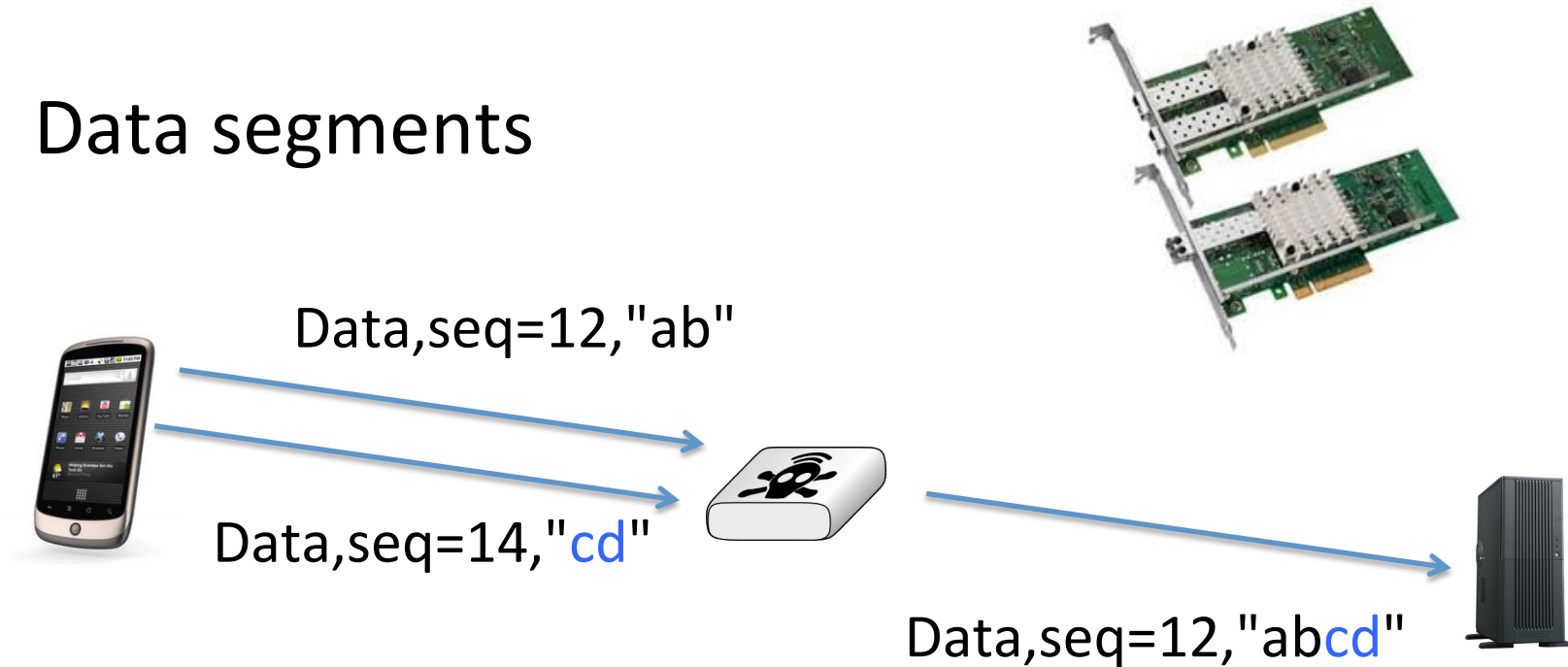
| Source port | Destination port |
|---|---|
| **Sequence number** | |
| Acknowledgment number | |

| THL | Reserved | Flags | Window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|
| ***Datasequence number*** | |
| Payload | |

# Multipath TCP
# Data transfer

Dseq=0,seq=123,"a"

DAck=1,ack=124

DSeq=2, seq=124,"c"

DAck=3, ack=125

DSeq=1, seq=456,"b"

DAck=2,ack=457

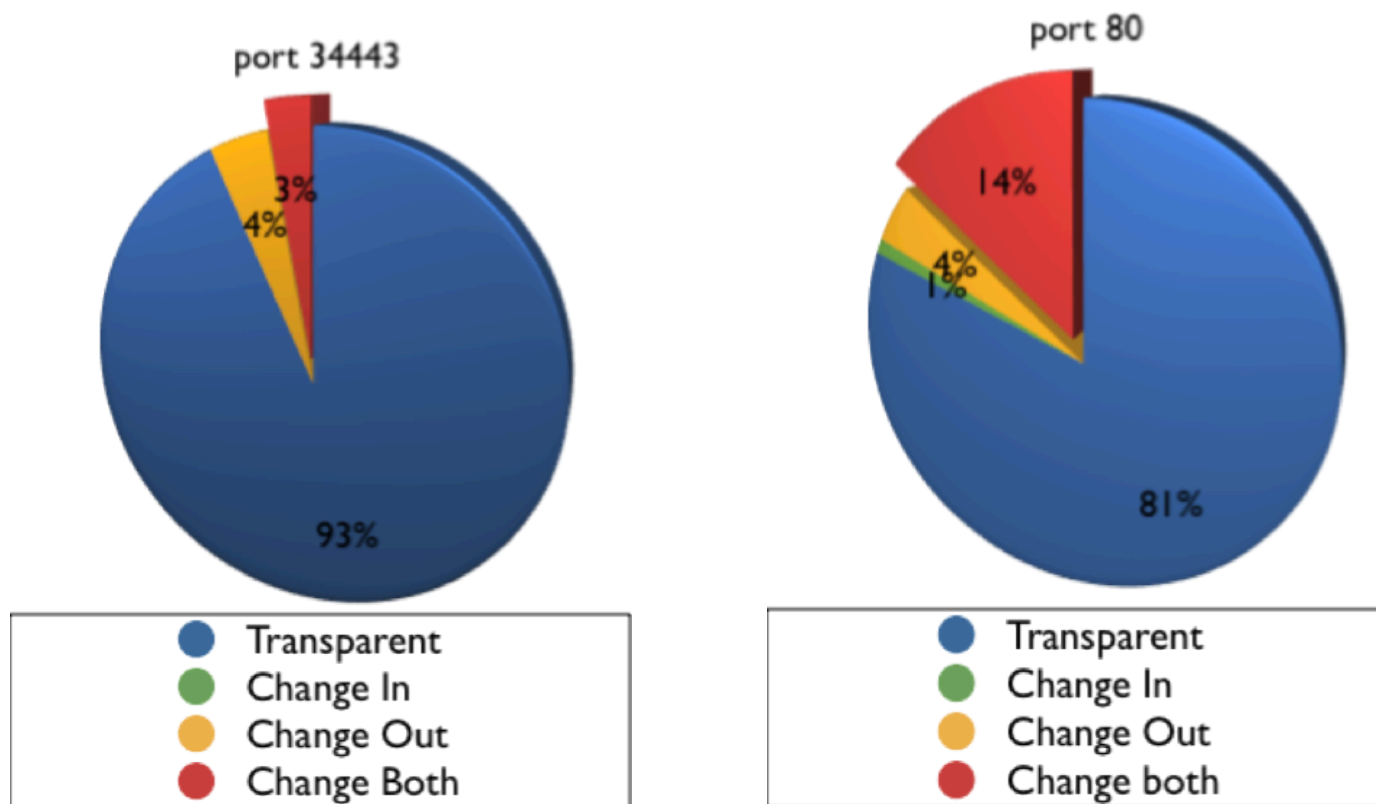# Middlebox interference

- Data segments

Data,seq=12,"ab"

Data,seq=14,"cd"

Data,seq=12,"abcd"

Such a middlebox could also be the network adapter of the
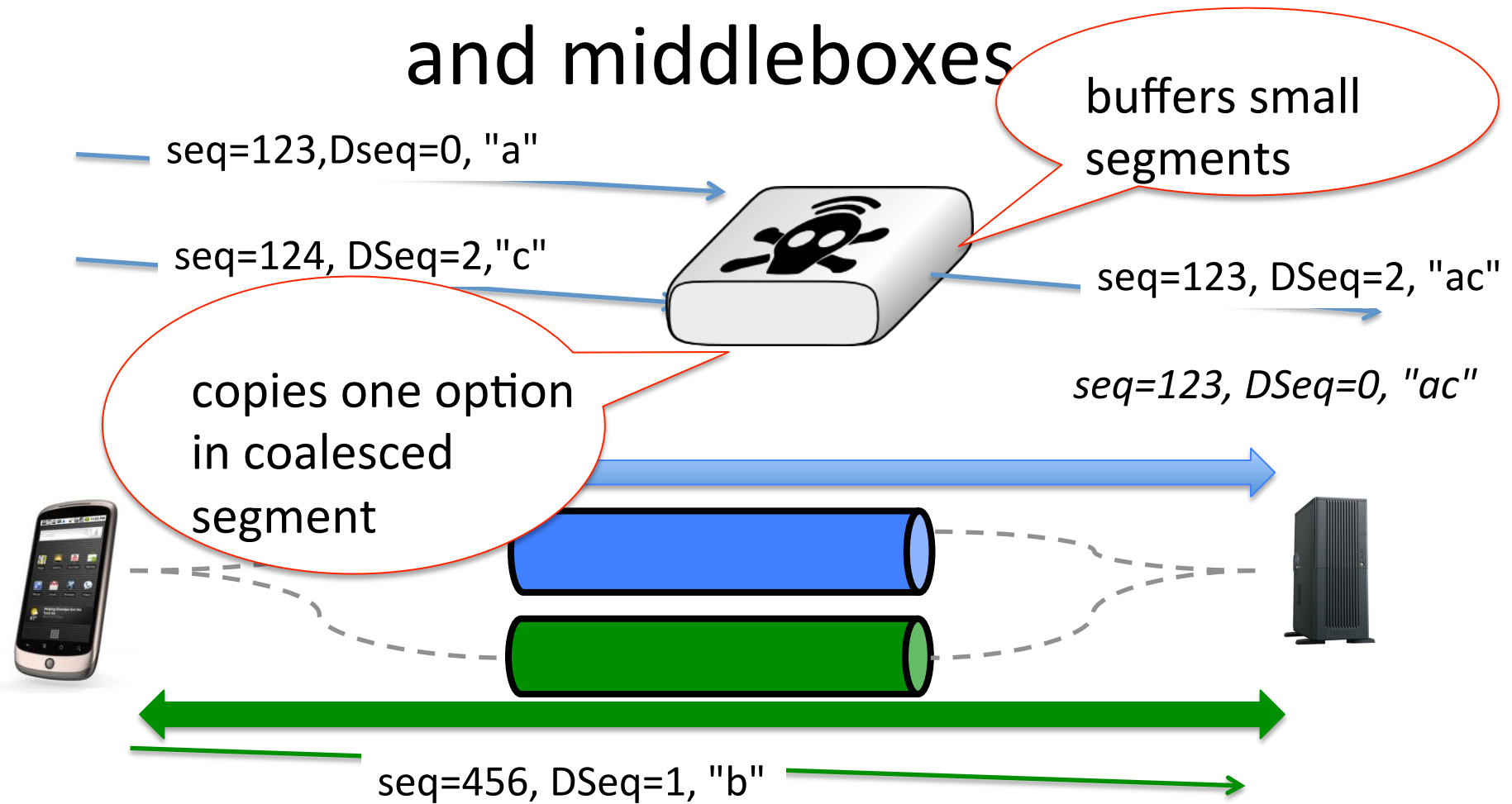server that uses LRO to improve performance.
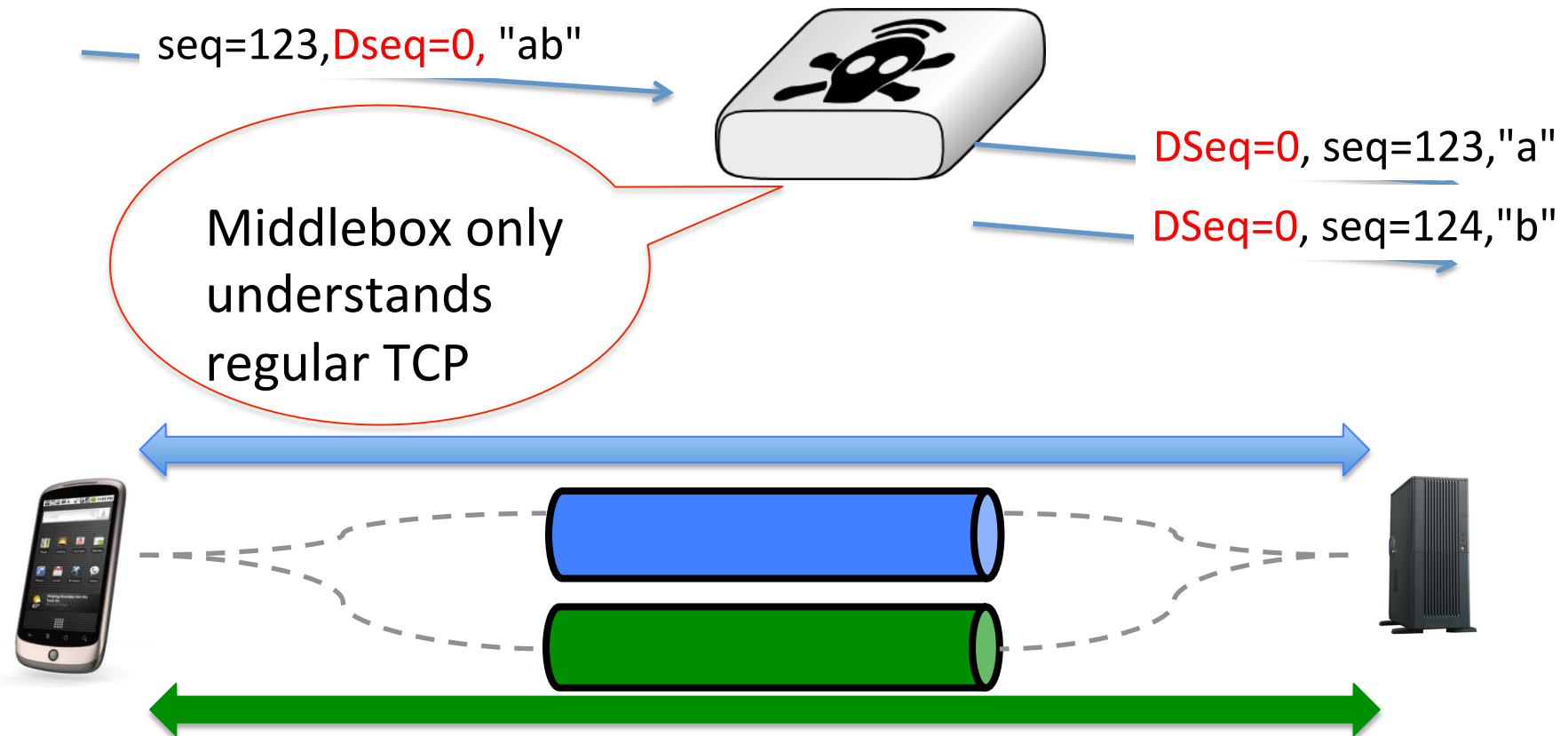
# Segment coalescing



Honda, Michio, et al. "Is it still possible to extend TCP?." Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference. ACM, 2011.

© O. Bonaventure, 2011

# Data sequence numbers and middleboxes

seq=123,Dseq=0, "a"

seq=124, DSeq=2,"c"

buffers small segments

seq=123, DSeq=2, "ac"

seq=123, DSeq=0, "ac"

copies one option in coalesced segment

seq=456, DSeq=1, "b"

# Data sequence numbers and middleboxes

seq=123,Dseq=0, "ab"



DSeq=0, seq=123,"a"

DSeq=0, seq=124,"b"

Middlebox only understands regular TCP
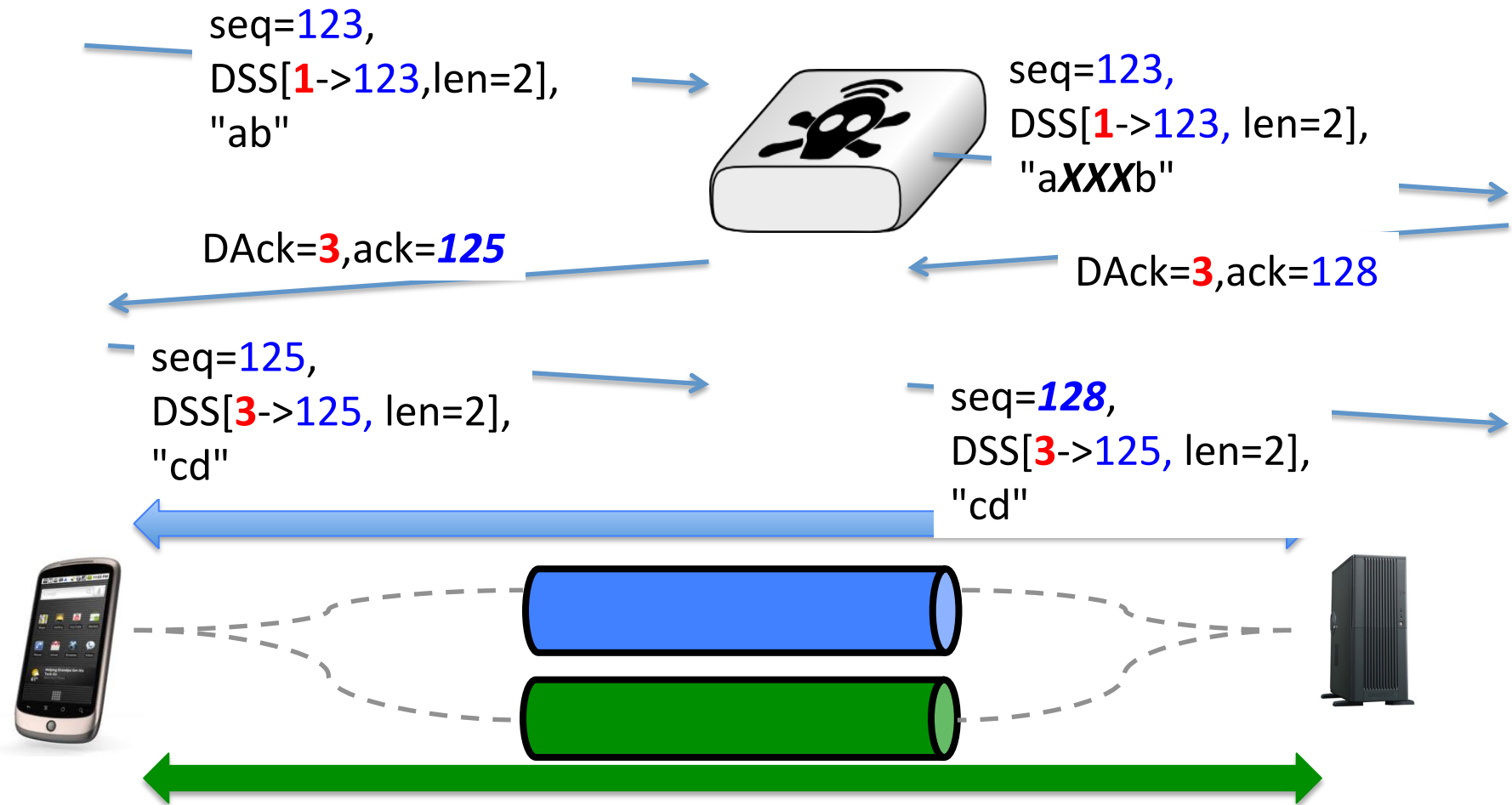
# Data sequence numbers and middleboxes

- How to avoid desynchronisation between the bytestream and data sequence numbers ?

- Solution
  - Multipath TCP option carries **mapping** between Data sequence numbers and *(difference between initial and current)* subflow sequence numbers
    - mapping covers a part of the bytestream (length)

# Multipath TCP and middleboxes

- With the DSS mapping, Multipath TCP can cope with middleboxes that
  - combine segments
  - split segments

- Are they the most annoying middleboxes for Multipath TCP ?

  - Unfortunately not

# The worst middlebox

seq=123,
DSS[**1**->123,len=2],
"ab"

seq=123,
DSS[**1**->123, len=2],
 "a***XXX***b"

DAck=**3**,ack=***125***

DAck=**3**,ack=128

seq=125,
DSS[**3**->125, len=2],
"cd"

seq=***128***,
DSS[**3**->125, len=2],
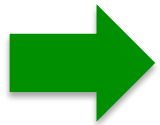"cd"

- Is this an academic exercise or reality ?

# The worst middlebox

- Is unfortunately very old...
  - Any ALG for a NAT

220 ProFTPD 1.3.3d Server (BELNET FTPD Server) [193.190.67.15]
ftp_login: user `<null>' pass `<null>' host `ftp.belnet.be'
Name (ftp.belnet.be:obo): anonymous
---> USER anonymous
331 Anonymous login ok, send your complete email address as your password
Password:
---> PASS XXXX
---> **PORT 192,168,0,7,195,120**
200 PORT command successful
---> LIST
150 Opening ASCII mode data connection for file list
lrw-r--r--   1 ftp     ftp          6 Jun  1  2011 pub -> mirror
226 Transfer complete

# Coping with the worst middlebox

- What should Multipath TCP do in the presence of such a worst middlebox ?
  - Do nothing and ignore the middlebox
    - but then the bytestream and the application would be broken and this problem will be difficult to debug by network administrators

  - Detect the presence of the middlebox
    - and fallback to regular TCP (i.e. use a single path and nothing fancy)

  Multipath TCP **MUST** work in all networks where regular TCP works.

# Detecting the worst middlebox ?

- How can Multipath TCP detect a middlebox that modifies the bytestream and inserts/removes bytes ?

  – Various solutions were explored

  – In the end, Multipath TCP chose to include its own checksum to detect insertion/deletion of bytes
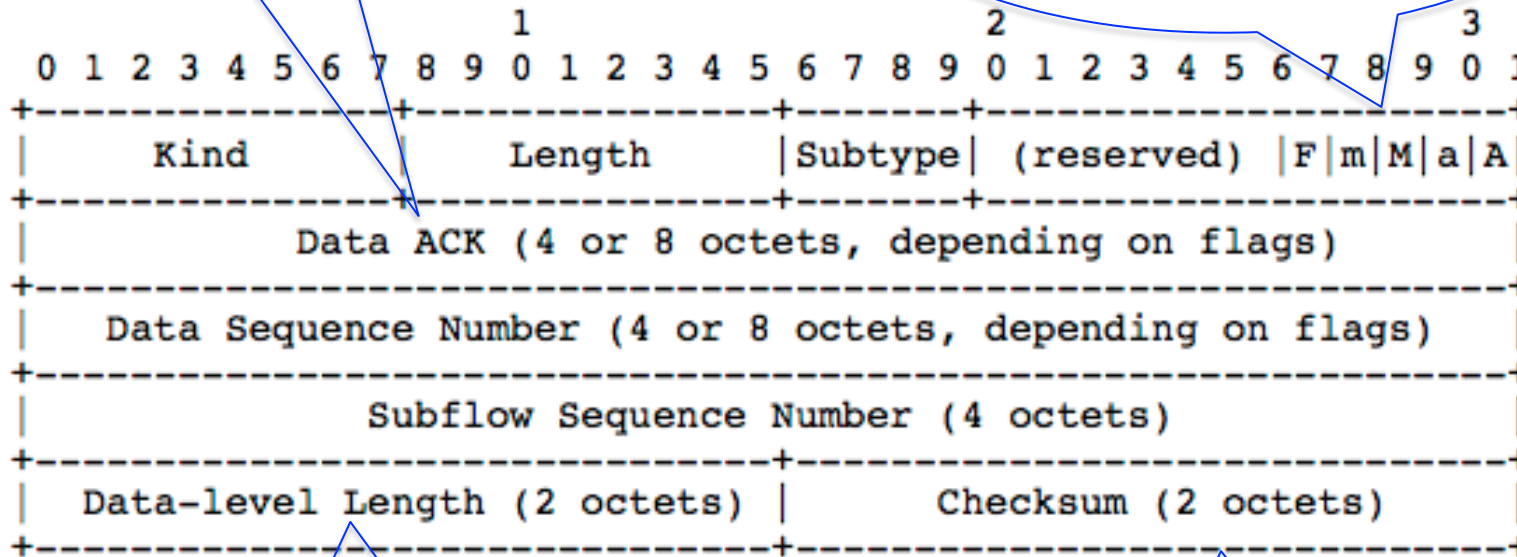
# Multipath TCP
# Data sequence numbers

- Data sequence numbers and Data acknowledgements

    - Maintained inside implementation as 64 bits field

    - Implementations can, as an optimisation, only transmit the lower 32 bits of the data sequence and acknowledgements

# Data Sequence Signal option

A = Data ACK present
a = Data ACK is 8 octets
M = mapping present
m = DSN is 8

Cumulative Data ack

```
                    1                   2                   3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---------------+---------------+-------+-------+---------------+
|      Kind     |     Length    |Subtype| (reserved) |F|m|M|a|A|
+---------------+---------------+-------+-------+---------------+
|              Data ACK (4 or 8 octets, depending on flags)     |
+---------------------------------------------------------------+
|       Data Sequence Number (4 or 8 octets, depending on flags)|
+---------------------------------------------------------------+
|              Subflow Sequence Number (4 octets)               |
+-------------------------------+-------------------------------+
|  Data-level Length (2 octets) |      Checksum (2 octets)      |
+-------------------------------+-------------------------------+
```

Length of mapping, can extend beyond this segment

Computed over data covered by entire mapping + pseudo header

# The Multipath TCP protocol

- Control plane
  - How to manage a Multipath TCP connection that uses several paths ?

- Data plane
  - How to transport data ?

➡ **Congestion control**
  - How to control congestion over multiple paths ?

# AIMD in TCP

- Congestion control mechanism
  - Each host maintains a *congestion window (cwnd)*
  - No congestion
    - Congestion avoidance (**additive increase**)
      - increase *cwnd* by one segment every round-trip-time
  - Congestion
    - TCP detects congestion by detecting losses
    - Mild congestion (fast retransmit – **multiplicative decrease**)
      - *cwnd=cwnd*/2 and restart congestion avoidance
    - Severe congestion (timeout)
      - *cwnd*=1, set slow-start-threshold and restart slow-start

# Congestion control for Multipath TCP

- Simple approach
  - independant congestion windows

Threshold

Threshold

Threshold

# Independant congestion windows

- Problem

12Mbps

# Coupled congestion control

- Congestion windows are coupled
  - congestion window growth cannot be faster than TCP with a single flow
  - Coupled congestion control aims at **moving traffic away from congested path**

# Linked increases congestion control

- Algorithm
  - For each loss on path r, $cwin_r = cwin_r / 2$

  - Additive increase

$$cwin_r = cwin_r + \min\left(\frac{\max\left(\frac{cwnd_i}{(rtt_i)^2}\right)}{\left(\sum_i \frac{cwnd_i}{rtt_i}\right)^2}, \frac{1}{cwnd_r}\right)$$

D. Wischik, C. Raiciu, A. Greenhalgh, and M. Handley, "Design, implementation and evaluation of congestion control for multipath TCP," NSDI'11: Proceedings of the 8th USENIX conference on Networked systems design and implementation, 2011.

# Other Multipath-aware congestion control schemes

R. Khalili, N. Gast, M. Popovic, U. Upadhyay, J.-Y. Le Boudec , MPTCP is not Pareto-optimal: Performance issues and a possible solution, Proc. ACM Conext 2012

Y. Cao, X. Mingwei, and X. Fu, "Delay-based Congestion Control for Multipath TCP," ICNP2012, 2012.

T. A. Le, C. S. Hong, and E.-N. Huh, "Coordinated TCP Westwood congestion control for multiple paths over wireless networks," ICOIN '12: Proceedings of the The International Conference on Information Network 2012, 2012, pp. 92–96.

T. A. Le, H. Rim, and C. S. Hong, "A Multipath Cubic TCP Congestion Control with Multipath Fast Recovery over High Bandwidth-Delay Product Networks," *IEICE Transactions*, 2012.

T. Dreibholz, M. Becke, J. Pulinthanath, and E. P. Rathgeb, "Applying TCP-Friendly Congestion Control to Concurrent Multipath Transfer," Advanced Information Networking and Applications (AINA), 2010 24th IEEE International Conference on, 2010, pp. 312–319.

# The Multipath TCP protocol

→ **Control plane**
  - How to manage a Multipath TCP connection that uses several paths ?

- Data plane
  - How to transport data ?

- Congestion control
  - How to control congestion over multiple paths ?

# The Multipath TCP control plane

- Connection establishment
  - Beware of middleboxes that remove TCP options
  - Limited space inside TCP option in SYN

- Closing a Multipath TCP connection
  - Decouple closing the Multipath TCP connection from closing the subflows

- Address dynamics

# Multipath TCP
# Connection establishment

- Principle

# Roles of the initial TCP handshake

- Check willingness to open TCP connection
  - Propose initial sequence number
  - Negotiate Maximum Segment Size
- TCP options
  - negotiate Timestamps, SACK, Window scale
- Multipath TCP
  - check that server supports Multipath TCP
  - propose Token in each direction
  - propose initial Data sequence number in each direction
  - Exchange keys to authenticate subflows

# Putting everything inside the SYN

- How can we place inside SYN segment ?

    – Initial Data Sequence Number (64 bits)

    – Token (32 bits)

    – Authentication Key (the longer the better)

# Constraint on TCP options

| Ver | IHL | ToS | Total length | |
|-----|-----|-----|--------------|---|
| Identification | | | Flags | Frag. Offset |
| TTL | | Protocol | Checksum | |
| Source IP address | | | | |
| Destination IP address | | | | |
| Source port | | Destination port | | |
| Sequence number | | | | |
| Acknowledgment number | | | | |
| THL | Reserved | Flags | Window | |
| Checksum | | Urgent pointer | | |
| *Options* | | | | |
| Payload | | | | |

- Total length of TCP header : max 64 bytes

  - max 40 bytes for TCP options
  - *Options* length must be multiple of 4 bytes

# Key exchange

SYN, [MyKey="keyABC"]

SYN+ACK, [MyKey="keyDEF"]

ACK[MyKey="keyABC", YourKey="keyDEF"]

*MyKey="keyABC"*
*YourKey="keyDEF"*

*MyKey="keyDEF"*
*YourKey="keyABC"*

SYN,[NonceA=123]

SYN+ACK[NonceB=456,
HMAC(123||456,"keyDEF||keyABC")]

ACK,[HMAC(456||123,"keyABC||keyDEF")]

# The Multipath TCP control plane

- Connection establishment in details

- Closing a Multipath TCP connection

- Address dynamics

# Multipath TCP
# Address dynamics

- How to learn the addresses of a host ?

    IP=2.3.4.5

    IP=3.4.5.6
    IP6=2a00:1450:400c:c05::69

- How to deal with address changes ?

    IP=1.2.3.4

    IP=4.5.6.7

# Address dynamics

- Basic solution : multihomed server

SYN, [...]

SYN+ACK, [...]

ACK[...]

ADD_ADDR[3.4.5.6]

IP=2.3.4.5

ADD_ADDR[2a00:1450:400c:c05::69]

IP=3.4.5.6

IP6=2a00:1450:400c:c05::69

SYN,[...]

SYN+ACK[...]

ACK[..]

# Address dynamics

- Basic solution : mobile client

SYN, [...]

SYN+ACK, [...]

ACK[...]

ADD_ADDR [**4.5.6.7**]

IP=**2.3.4.5**

SYN,[...]

IP=1.2.3.4

SYN+ACK[...]

IP=**4.5.6.7**

ACK[..]

REMOVE_ADDR[1.2.3.4]

# Address dynamics
# in today's Internet

SYN, [...]

SYN+ACK, [...]

ACK[...]

ADD_ADDR [10.0.0.2]

ADD_ADDR [10.0.0.2]

IP=2.3.4.5

IP=1.2.3.4

IP=10.0.0.2

**?**

SYN [...]

# Address dynamics with NATs

- Solution

  - Each address has one identifier
    - Subflow is established between id=0 addresses
  - Each host maintains a list of <address,id> pairs of the addresses associated to an MPTCP endpoint
  - MPTCP options refer to the address identifier
    - ADD_ADDR contains  <address,id>
    - REMOVE_ADDR contains <id>

# Address dynamics

SYN, [...]

SYN+ACK, [...]

ACK[...]

ADD_ADDR [4.5.6.7,id=1]

IP=2.3.4.5

SYN,[id=1...]

IP=1.2.3.4

SYN+ACK[...]

IP=4.5.6.7

ACK[..]

REMOVE_ADDR[id=0]

# Agenda

- The motivations for Multipath TCP

- The changing Internet

- The Multipath TCP Protocol

- Multipath TCP use cases
  - Datacenters
  - Smartphones
  - IPv4/IPv6 coexistence

# TCP on servers

- How to increase server bandwidth ?



- Load balancing techniques
  - packet per packet
  - per flow load balancing
    - each TCP connection is mapped onto one interface

# Increasing server bandwidth
# with Multipath TCP



- Load balancing with Multipath TCP
  - Congestion control efficiently uses the two links
    for each MPTCP connection
  - Automatic failover in case of failures

# How fast can Multipath TCP go ?

# How fast can Multipath TCP go ?

# Datacenters evolve

- Traditional Topologies are tree-based
  - Poor performance
  - Not fault tolerant



- Shift towards multipath topologies: FatTree, BCube, VL2,

  Cisco, EC2



C. Raiciu, et al. "Improving datacenter performance and robustness with multipath TCP," *ACM SIGCOMM* 2011.

# TCP in data centers

# TCP in FAT tree networks
# Cost of collissions



C. Raiciu, et al. "Improving datacenter performance and robustness with multipath TCP," *ACM SIGCOMM* 2011.

# How to get rid of these collisions ?

- Consider TCP performance as an optimisation problem

# The Multipath TCP way



ECMP balances the subflows over different paths

Two subflows differ by their source port

Improving datacenter performance and robustness with multipath TCP," *ACM* 11.

# MPTCP better utilizes the FatTree network



C. Raiciu, et al. "Improving datacenter performance and robustness with multipath TCP," *ACM SIGCOMM* 2011.

See also G. Detal, et al. *, Revisiting Flow-Based Load Balancing: Stateless Path Selection in Data Center Networks*, Computer Networks, April 2013 for extensions to ECMP for MPTCP

# Multipath TCP on EC2

- Amazon EC2: infrastructure as a service
  - We can borrow virtual machines by the hour
  - These run in Amazon data centers worldwide
  - We can boot our own kernel
- A few availability zones have multipath topologies
  - 2-8 paths available between hosts not on the same machine or in the same rack
  - Available via ECMP

# Amazon EC2 Experiment

- 40 medium CPU instances running MPTCP
- During 12 hours, we sequentially ran all-to-all `iperf` cycling through:
  - TCP
  - MPTCP (2 and 4 subflows)

# MPTCP improves performance on EC2



C. Raiciu, et al. "Improving datacenter performance and robustness with multipath TCP," *ACM SIGCOMM* 2011.

# Agenda

- The motivations for Multipath TCP

- The changing Internet

- The Multipath TCP Protocol

- Multipath TCP use cases
  - Datacenters
  → Smartphones
  - IPv4/IPv6 coexistence

# Motivation

- One device, many IP-enabled interfaces

# MPTCP over WiFi/3G



8Mbps, 20ms

2Mbps, 150ms

# TCP over WiFi/3G



C. Raiciu, et al. "How hard can it be? designing and implementing a deployable multipath TCP," NSDI'12: Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation, 2012.

# MPTCP over WiFi/3G



C. Raiciu, et al. "How hard can it be? designing and implementing a deployable multipath TCP," NSDI'12: Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation, 2012.

# MPTCP over WiFi/3G

# MPTCP over WiFi/3G

# Understanding the performance issue

**Window full !**
**No new data can be sent on WiFi path**

D C B

8Mbps, 20ms

2Mbps, 150ms

A

Window

**Reinject segment on fast path**

**Halve congestion window on slow subflow**

# MPTCP over WiFi/3G

# Usage of 3G and WiFI

- How should Multipath TCP use 3G and WiFi ?

    - Full mode
        - Both wireless networks are used at the same time

    - Backup mode
        - Prefer WiFi when available, open subflows on 3G and use them as backup

    - Single path mode
        - Only one path is used at a time, WiFi preferred over 3G

# Evaluation scenario

WiFi:
Belgacom ADSL2+
(~8 Mbps, ~30 ms)

3G: Mobistar
(~2 Mbps, ~80ms)

# Recovery after failure



C. Paasch, et al. , "Exploring mobile/WiFi handover with multipath TCP," presented at the CellNet '12: Proceedings of the 2012 ACM SIGCOMM workshop on Cellular networks: operations, challenges, and future design, 2012.

# Recovery after failure



C. Paasch, et al. , "Exploring mobile/WiFi handover with multipath TCP," presented at the CellNet '12: Proceedings of the 2012 ACM SIGCOMM workshop on Cellular networks: operations, challenges, and future design, 2012.

# Agenda

- The motivations for Multipath TCP

- The changing Internet

- The Multipath TCP Protocol

- Multipath TCP use cases
  - Datacenters
  - Smartphones
  - IPv4/IPv6 coexistence

# IPv6 is coming …



Source http://6lab.cisco.com/stats/cible.php?country=world

# But IPv4 and IPv6 perf. may differ



Distribution of IPv4/IPv6 relative performance

E. Aben, *Measuring World IPv6 Day - Comparing IPv4 and IPv6 Performance*,
https://labs.ripe.net/Members/emileaben/measuring-world-ipv6-day-comparing-ipv4-and-ipv6-performance

# Happy eyeballs

# How to get best of IPv4 and IPv6 ?

# Conclusion

- Multipath TCP is becoming a reality
  - Due to the middleboxes, the protocol is more complex than initially expected
  - RFC has been published
  - there is running code !
  - Multipath TCP works over today's Internet !
- What's next ?
  - More use cases
    - Anycast, VM migration, storage, ...
  - Measurements and improvements to the protocol
    - Time to revisit 20+ years of heuristics added to TCP

?

# Try it by yourself !
# http://multipath-tcp.org

# References

- The Multipath TCP protocol
  - http://www.multipath-tcp.org
  - http://tools.ietf.org/wg/mptcp/

A. Ford, C. Raiciu, M. Handley, S. Barre, and J. Iyengar, "Architectural guidelines for multipath TCP development", RFC6182 2011.

A. Ford, C. Raiciu, M. J. Handley, and O. Bonaventure, "TCP Extensions for Multipath Operation with Multiple Addresses," RFC6824, 2013

C. Raiciu, C. Paasch, S. Barre, A. Ford, M. Honda, F. Duchene, O. Bonaventure, and M. Handley, "How hard can it be? designing and implementing a deployable multipath TCP," NSDI'12: Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation, 2012.

# Implementations

- Linux
  - http://www.multipath-tcp.org

    S. Barre, C. Paasch, and O. Bonaventure, "Multipath tcp: From theory to practice," *NETWORKING 2011*, 2011.

    Sébastien Barré. Implementation and assessment of Modern Host-based Multipath Solutions. PhD thesis. UCL, 2011

- FreeBSD
  - http://caia.swin.edu.au/urp/newtcp/mptcp/

- Simulators
  - http://nrg.cs.ucl.ac.uk/mptcp/implementation.html
  - http://code.google.com/p/mptcp-ns3/

# Middleboxes

M. Honda, Y. Nishida, C. Raiciu, A. Greenhalgh, M. Handley, and H. Tokuda, "Is it still possible to extend TCP?,"  IMC '11: Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference, 2011.

V. Sekar, N. Egi, S. Ratnasamy, M. K. Reiter, and G. Shi, "Design and implementation of a consolidated middlebox architecture," *USENIX NSDI*, 2012.

J. Sherry, S. Hasan, C. Scott, A. Krishnamurthy, S. Ratnasamy, and V. Sekar, "Making middleboxes someone else's problem: network processing as a cloud service," SIGCOMM '12: Proceedings of the ACM SIGCOMM 2012 conference on Applications, technologies, architectures, and protocols for computer communication, 2012.

# Multipath congestion control

- Background

  D. Wischik, M. Handley, and M. B. Braun, "The resource pooling principle," *ACM SIGCOMM Computer …*, vol. 38, no. 5, 2008.

  F. Kelly and T. Voice. Stability of end-to-end algorithms for joint routing and rate control. ACM SIGCOMM CCR, 35, 2005.

  P. Key, L. Massoulie, and P. D. Towsley, "Path Selection and Multipath Congestion Control," INFOCOM 2007. 2007, pp. 143–151.
- Coupled congestion control

  C. Raiciu, M. J. Handley, and D. Wischik, "Coupled Congestion Control for Multipath Transport Protocols," *RFC*, vol. 6356, Oct. 2011.

  D. Wischik, C. Raiciu, A. Greenhalgh, and M. Handley, "Design, implementation and evaluation of congestion control for multipath TCP," NSDI'11: Proceedings of the 8th USENIX conference on Networked systems design and implementation, 2011.

# Multipath congestion control

– More

R. Khalili, N. Gast, M. Popovic, U. Upadhyay, J.-Y. Le Boudec , MPTCP is not Pareto-optimal: Performance issues and a possible solution, Proc. ACM Conext 2012

Y. Cao, X. Mingwei, and X. Fu, "Delay-based Congestion Control for Multipath TCP," ICNP2012, 2012.

T. A. Le, C. S. Hong, and E.-N. Huh, "Coordinated TCP Westwood congestion control for multiple paths over wireless networks," ICOIN '12: Proceedings of the The International Conference on Information Network 2012, 2012, pp. 92–96.

T. A. Le, H. Rim, and C. S. Hong, "A Multipath Cubic TCP Congestion Control with Multipath Fast Recovery over High Bandwidth-Delay Product Networks," *IEICE Transactions*, 2012.

T. Dreibholz, M. Becke, J. Pulinthanath, and E. P. Rathgeb, "Applying TCP-Friendly Congestion Control to Concurrent Multipath Transfer," Advanced Information Networking and Applications (AINA), 2010 24th IEEE International Conference on, 2010, pp. 312–319.

# Use cases

- Datacenter

  C. Raiciu, S. Barre, C. Pluntke, A. Greenhalgh, D. Wischik, and M. J. Handley, "Improving datacenter performance and robustness with multipath TCP," *ACM SIGCOMM* 2011.

  G. Detal, Ch. Paasch, S. van der Linden, P. Mérindol, G. Avoine, O. Bonaventure, *Revisiting Flow-Based Load Balancing: Stateless Path Selection in Data Center Networks*, Computer Networks, April 2013

- Mobile

  C. Pluntke, L. Eggert, and N. Kiukkonen, "Saving mobile device energy with multipath TCP," *MobiArch '11: Proceedings of the sixth international workshop on MobiArch*, 2011.

  C. Paasch, G. Detal, F. Duchene, C. Raiciu, and O. Bonaventure, "Exploring mobile/WiFi handover with multipath TCP," CellNet '12: Proceedings of the 2012 ACM SIGCOMM workshop on Cellular networks: operations, challenges, and future design, 2012.