# MultiPath TCP : Linux Kernel implementation
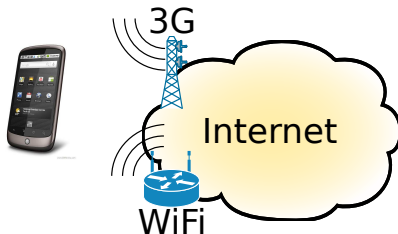
Presenter: Christoph Paasch
IP Networking Lab
UCLouvain, Belgium

August 28, 2012
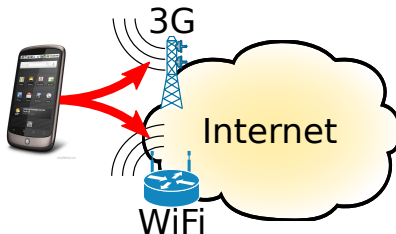
http://multipath-tcp.org

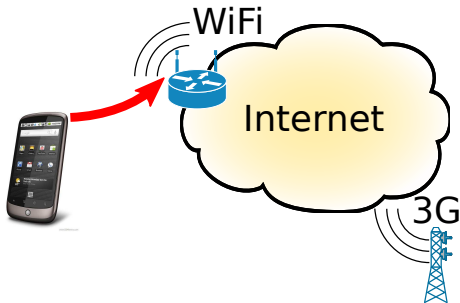Mobile devices can connect to the Internet via different interfaces

Mobile devices can connect to the Internet via different interfaces

Smartphones have to restart their data-transfer when moving
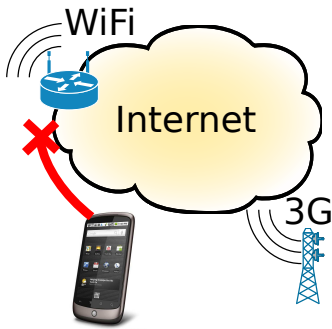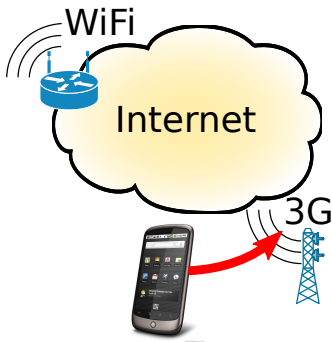away from the WiFi access-point.

Smartphones have to restart their data-transfer when moving away from the WiFi access-point.

Smartphones have to restart their data-transfer when moving away from the WiFi access-point.

Data-centers have a large redundant infrastructure
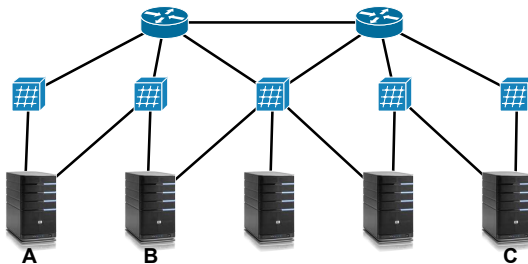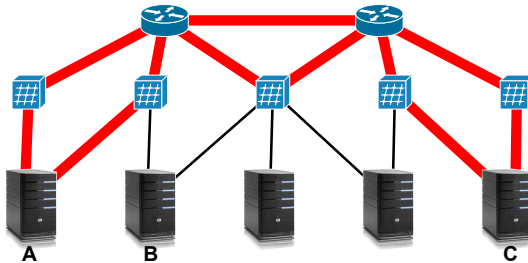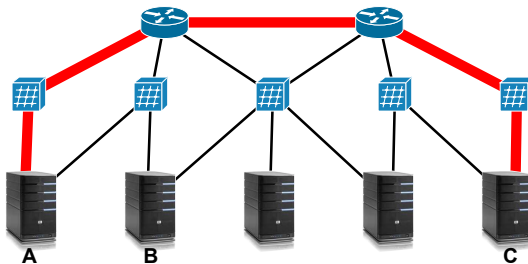
Data-centers have a large redundant infrastructure

Collisions in data-center reduce the bandwidth and result in suboptimal load-balancing

Collisions in data-center reduce the bandwidth and result in suboptimal load-balancing

Collisions in data-center reduce the bandwidth and result in suboptimal load-balancing

Collisions in data-center reduce the bandwidth and result in suboptimal load-balancing

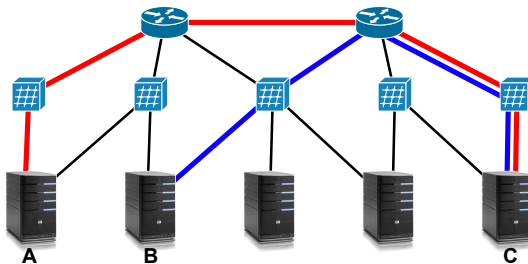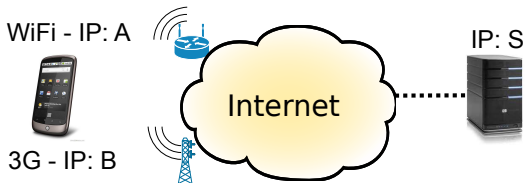Mismatch between the **multipath network** and the **single-path transport** protocol.

# MultiPath TCP

- Runs with **unmodified applications**
- Works over **today's Internet**
- **IPv4/IPv6** are both supported (even simultaneously)

# MultiPath TCP



WiFi - IP: A

3G - IP: B

IP: S

**Is the server MPTCP-capable?**

SrcIP: A
DstIP: S

SYN
MP_CAPABLE : ID12

**Is the server MPTCP-capable?**

SrcIP: S
DstIP: A

SYN + ACK
MP_CAPABLE: ID23

**Sending Data - Naiv approach**

**Sending Data - Naiv approach**

**Sending Data - Naiv approach**

**Sending Data - Naiv approach**

**Establish separate subflows**



SrcIP: B
DstIP: S

SYN
MP_JOIN : ID23

**Establish separate subflows**



SrcIP: S
DstIP: B

SYN + ACK
MP_JOIN

**Establish separate subflows**

**Establish separate subflows**



Subflow 1
write_seq
snd_cwnd
rcv_nxt

Subflow 2
write_seq
snd_cwnd
rcv_nxt

**Sending Data**

**Sending Data**

# MultiPath TCP

- Both subflows can be used simultaneously.
- Addresses are advertised with ADD_ADDR and removed by REMOVE_ADDR.
- Subflows can be dynamically added and removed during the lifetime of the connection.

# MultiPath TCP

## Linux Kernel Implementation

Available at http://multipath-tcp.org

# Exchanged Messages



SrcIP: A
DstIP: S

SYN
MP_CAPABLE : ID12

High-Level Kernel design - Client Side

## In-depth call-stack - Client Side

# In-depth call-stack - Server Side

# Exchanged Messages



SrcIP: S
DstIP: A

SYN + ACK
MP_CAPABLE: ID23

# High-Level Kernel design - Client Side

In-depth call-stack - Client Side

# Exchanged Messages

# In-depth call-stack - Server Side

# High-Level Kernel design

# Exchanged Messages



SrcIP: B
DstIP: S

SYN
MP_JOIN : ID23

## High-Level Kernel design - Client Side

## In-depth call-stack - Client Side



```
mptcp_init4_subsockets(...);
```

```
mptcp_add_sock(...);
```
Attach subflow to
MPTCP-level

Creates a new IPv4 subflow

```
connect(...);
```

```
tcp_v4_connect(...);
```

```
tcp_connect(...);
```
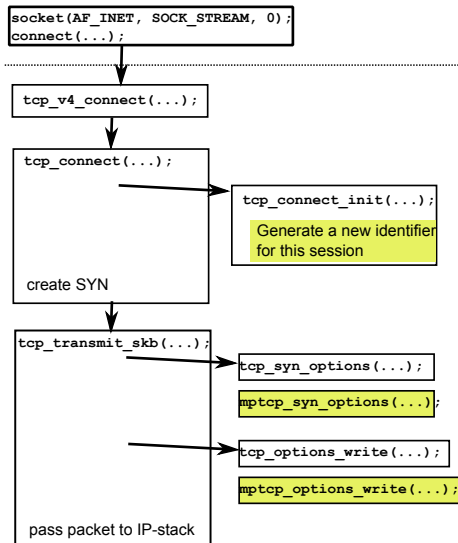Adds the MP_JOIN

SYN

## In-depth call-stack - Server Side

# Exchanged Messages



SrcIP: S
DstIP: B

SYN + ACK
MP_JOIN

## Exchanged Messages



ACK
MP_JOIN
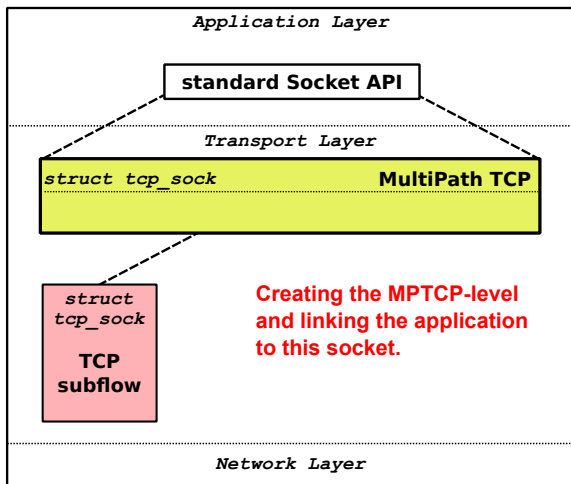
In-depth call-stack - Server Side

# Exchanged Messages

## High-Level Kernel design

# High-Level Kernel design

# Sending Data

## High-Level Kernel design

# Receiving Data

Packets can be reordered at the data-level due to delay-differences.



SubSeq: 2
DataSeq: 1002

SubSeq: 1
DataSeq: 1000

SubSeq: 200
DataSeq: 1001

Packets can be reordered at the data-level due to delay-differences.

Packets can be reordered at the data-level due to delay-differences.

A loss at the subflow-level (or network-reordering) can also cause reordering at the subflow-level

A loss at the subflow-level (or network-reordering) can also cause reordering at the subflow-level

# Receiving Data

Subflow-level out-of-order queues are necessary to handle the
retransmission at the subflow-level

# High-Level Kernel design

# MultiPath TCP

### Design Challenges

# Writing the Data-sequence number in the TCP-options

*Naive approach*

```
struct tcp_skb_cb {
        __u32              seq; /* Starting sequence number */
        [...]
#ifdef CONFIG_MPTCP
        __u32              data_seq;
        __u32              end_data_seq;
        __u32              data_ack;
        [...]
#endif
}
```

- Writing data-seq in *tcp_options_write*
- Increased *tcp_skb_cb* by 24 bytes.

*Our solution*

Inside the MPTCP-scheduler write the data-seq on top of the payload, **before** calling *tcp_transmit_skb*.



No more increase of *tcp_skb_cb*.

## Current Design

Current Design

## Current Design

Application-level

```
struct socket
struct sock *sk;
```

MPTCP-level

TCP subflow-level

```
struct tcp_sock
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;
```

## Current Design

# Current Design

Application-level

```
struct socket
struct sock *sk;
```

MPTCP-level

```
struct mptcp_cb
struct tcp_sock
```

TCP subflow-level

```
struct tcp_sock
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;
```

```
struct mptcp_tcp_sock
MPTCP-specific fields
```

## Current Design

Application-level

struct socket
struct sock *sk;

Change sk-pointer

MPTCP-level

**struct mptcp_cb**
**struct tcp_sock**

TCP subflow-level

**struct tcp_sock**
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;

**struct mptcp_tcp_sock**
MPTCP-specific fields

## Current Design

- Problems, if the application does a system-call on the socket, **before** the reception of the SYN+ACK
- Fix: Wait for the SYN+ACK. E.g., `tcp_sendmsg`:

```
/* Wait for a connection to finish. */
if ((1 << sk->sk_state) & ~(TCPF_ESTABLISHED |
   TCPF_CLOSE_WAIT))
       if ((err = sk_stream_wait_connect(sk, &timeo))
          != 0)
             goto do_error;
```

- We need to do this in **all** functions that take a lock on the socket! `tcp_recvmsg`, `tcp_splice_read`, `ip_setsockopt`, `ip_get_sockopt`, `tcp_ioctl`, ... and many more

*Upcoming Design*

*Upcoming Design*

## Upcoming Design

Application-level

```
struct socket
struct sock *sk;
```

```
struct tcp_sock
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;
```

MPTCP-level

TCP subflow-level

*Upcoming Design*

Application-level

**struct socket**
struct sock *sk;

**struct tcp_sock**
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;

MPTCP-level

sk_clone()

**struct tcp_sock**
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;

TCP subflow-level

# Creating the MPTCP-level on the client-side.



*Upcoming Design*

Application-level

**struct socket**
struct sock *sk;

MPTCP-level

**struct mptcp_cb**
MPTCP specific variables
Pointers to subsockets

**struct tcp_sock**
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;

TCP subflow-level

**struct tcp_sock**
struct mptcp_cb *mpcb;
struct mptcp_tcp_sock *mptcp;

**struct mptcp_tcp_sock**
Per socket MPTCP-specific fields

*Questions*

- Lots of socket options in the TCP/IP stack
- Some are for the MPTCP-level (*SO_SNDBUF*), some should get passed onto all other subflows (*IP_TTL*)
- This requires a lot of changes in TCP unrelated functions (e.g., `do_ip_setsockopt`)

How could we handle this?

# Still lots of changes to the TCP-stack

*Questions*

We have a lot of:

```
if (tcp_sk(sk)->mpc) {
        DO_SOME_MPTCP_STUFF
} else {
        DO_USUAL_TCP_STUFF
}
```

# Submitting MPTCP upstream???

- $\sim$ 10000 lines of code
- Tightly integrated in the TCP-stack
- More work to do:
  - Cleanup - better separate MPTCP from TCP
  - Some missing features
  - Support TSO
  - Support NET_DMA
  - . . .

- How to split the patch in small pieces?

# Conclusion

Freely available at **http://multipath-tcp.org**
Download it, try it out, contribute!

UCLouvain MPTCP-Team:
Sébastien Barré
Christoph Paasch
Gregory Detal
Fabien Duchene

Prof. Olivier Bonaventure

Thanks to our previous and present partners/contributors: