

# Application-defined Multipath TCP Scheduling

Frühjahrstreffen der GI Fachgruppen BS und KuVS  
Amr Rizk



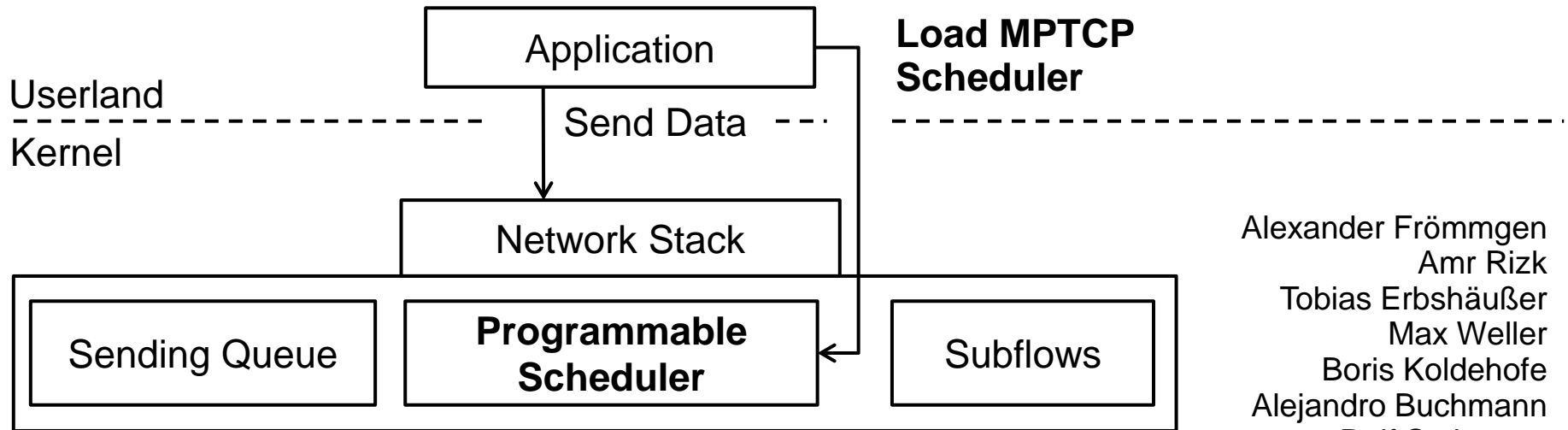
TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



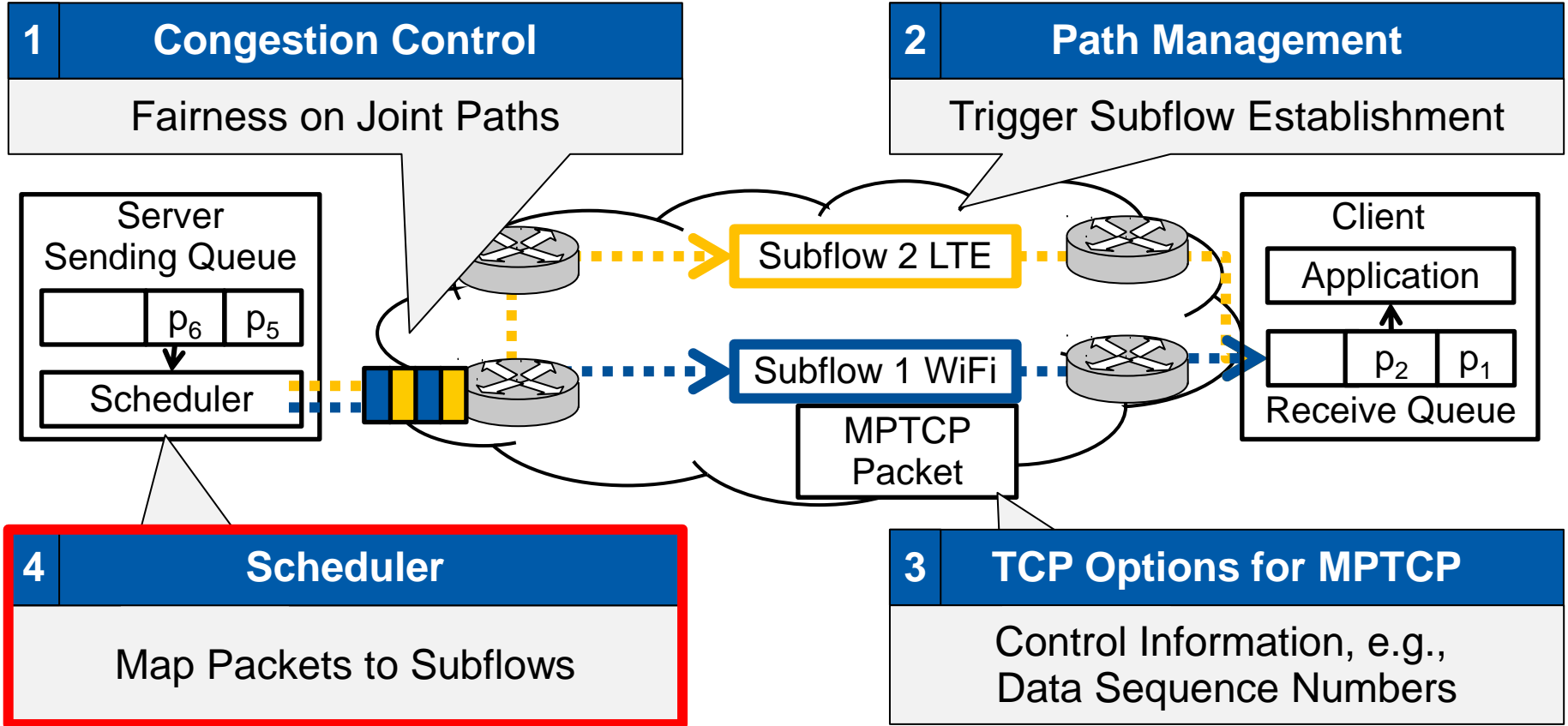
KOM - Multimedia  
Communications Lab



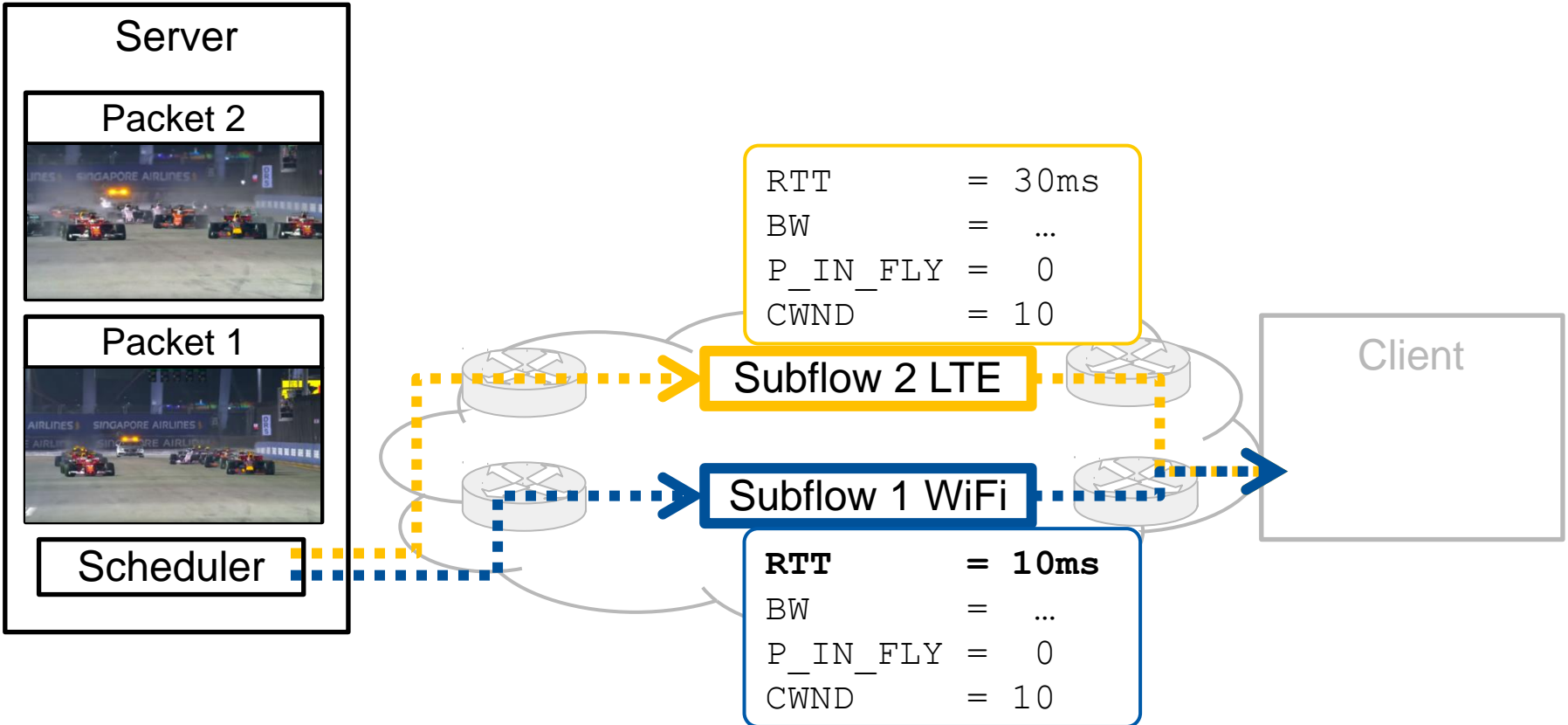
DFG Collaborative Research Centre 1053 - MAKI  
Multi Mechanism Adaptation for the Future Internet



# Multipath TCP in a Nutshell

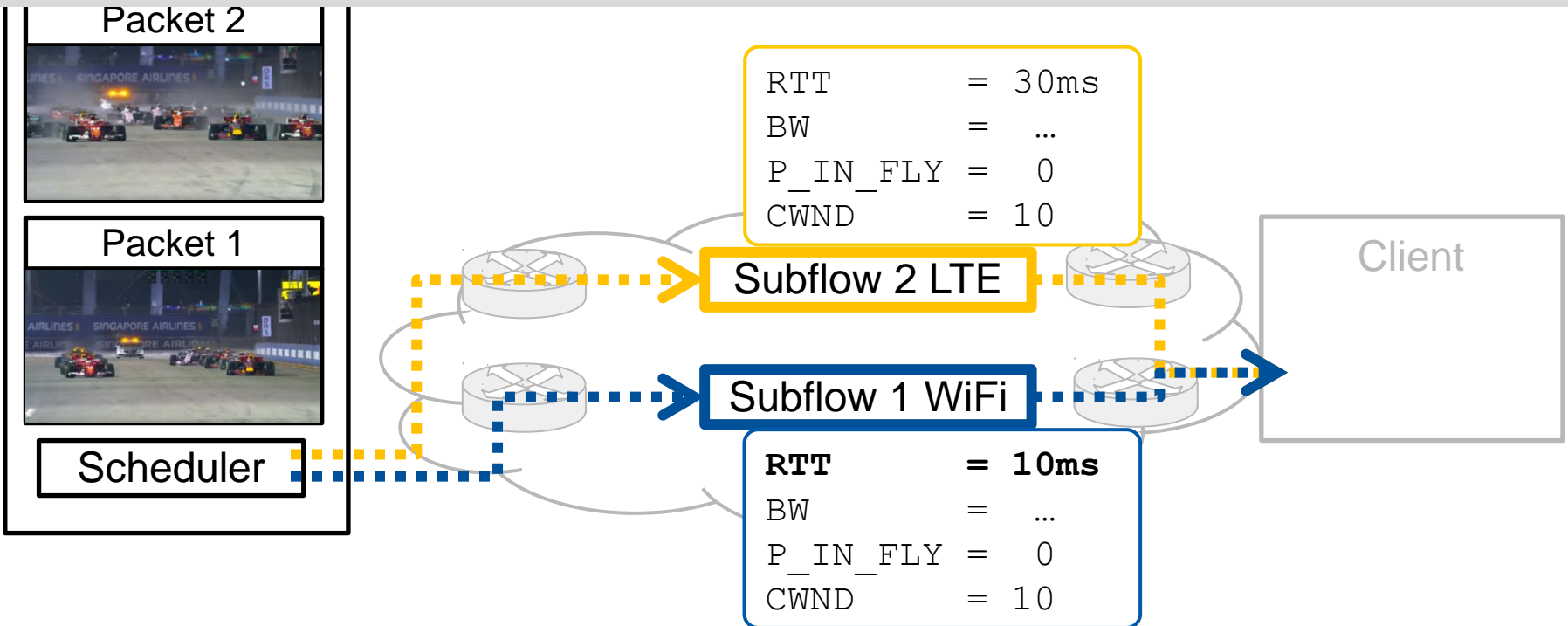


# Multipath TCP Scheduling



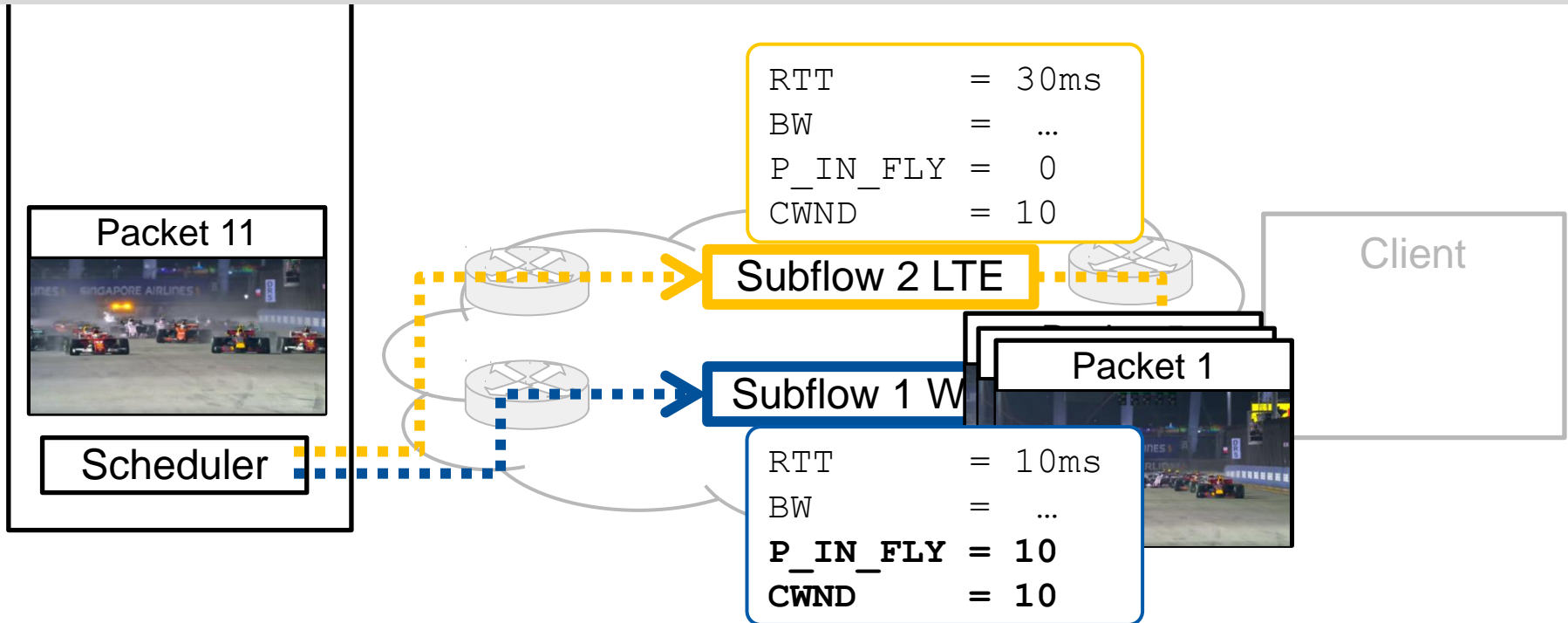
# Multipath TCP Scheduling

Intuition: Schedule on Subflow with minimum round-trip time (**minRTT**).



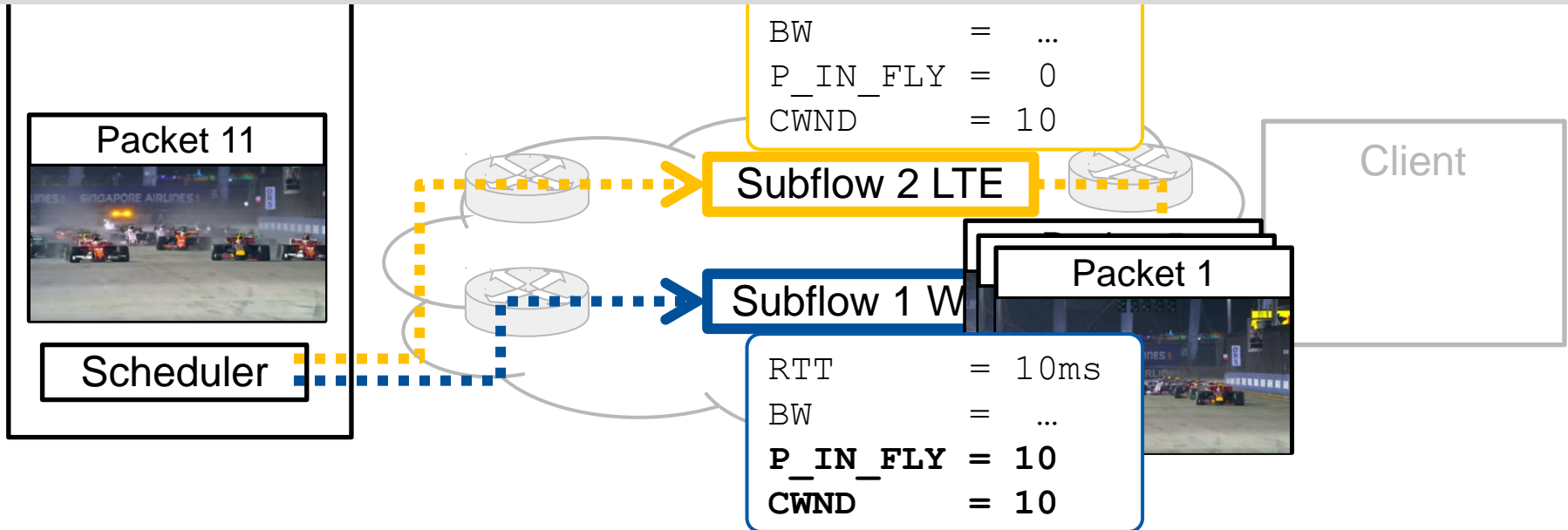
# Multipath TCP Scheduling

Schedule on Subflow with minimum round-trip time (minRTT).



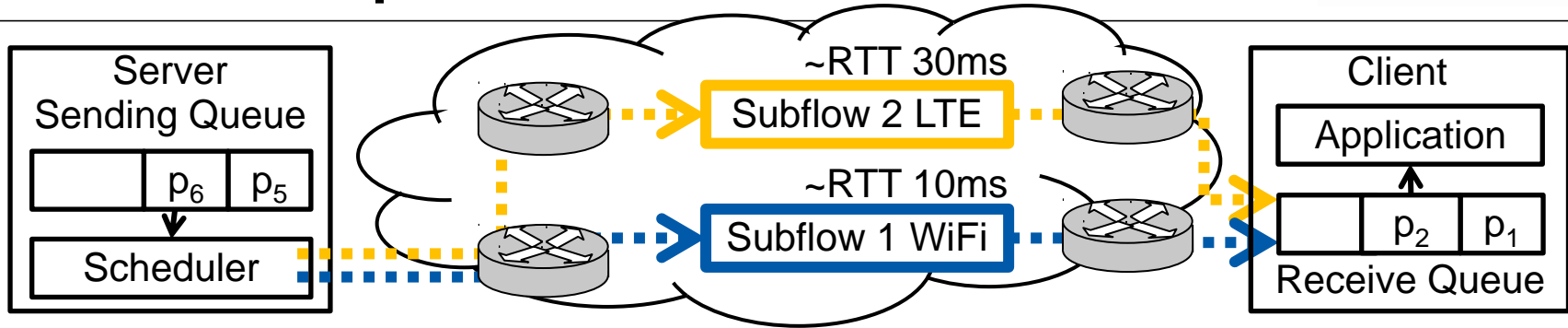
# Multipath TCP Scheduling

Schedule on Subflow with minimum round-trip time (**minRTT**),  
which is not saturated (Congestion window > packets in flight).



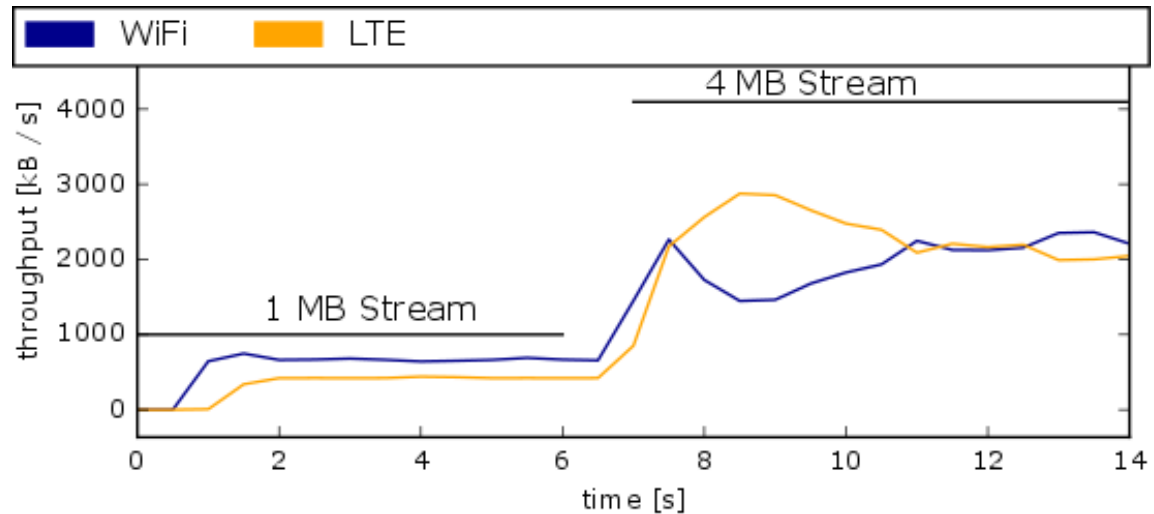
# Multipath TCP Scheduling

## Real-world Experiment



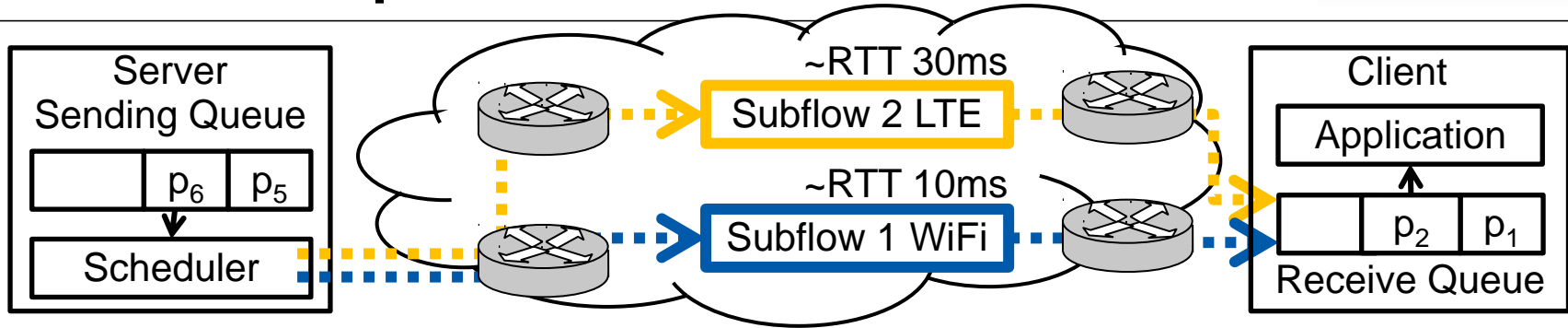
Constant Bitrate Stream

Home Network WiFi  
and LTE  
Germany



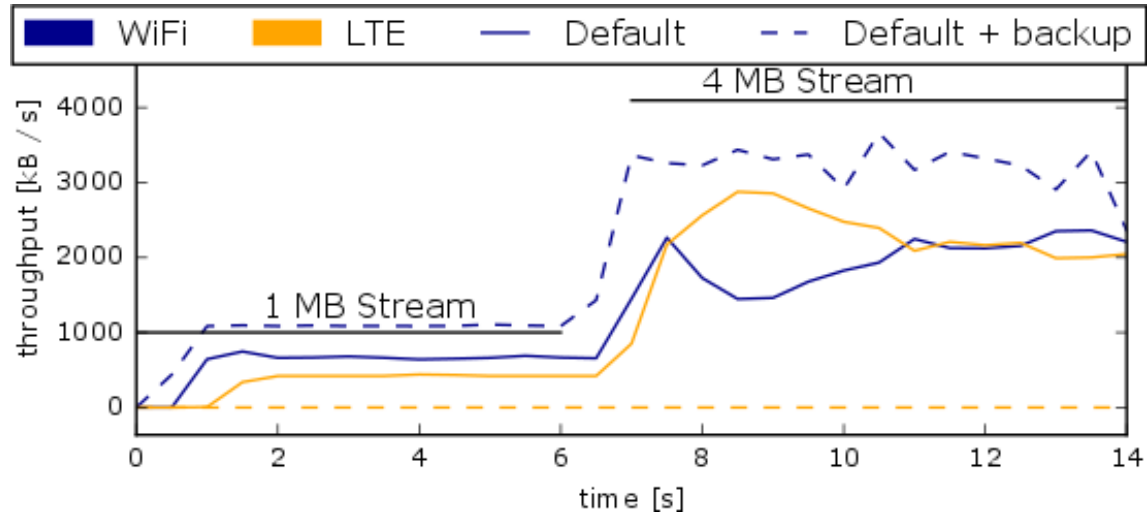
# Multipath TCP Scheduling

## Real-world Experiment



Constant Bitrate Stream

Home Network WiFi  
and LTE  
Germany





# Multipath TCP Scheduler Overview

Name	Domain	Preferences	Available in Linux (Implementation)
Min RTT ( <b>Default</b> )	General Purpose	Binary	✓
Round-Robin	General Purpose	Binary	✓
Redundant	Thin Flows	Binary	✓
Compensate Loss	Short Datacenter Flows	?	<b>No</b>
Energy preserving		?	<b>No</b>
DASH Video		Yes	<b>No</b>
...			

# Multipath TCP Scheduler Overview

## Great Concepts without MPTCP Implementation



The reference FIFO scheduler. We compare our cross-layer scheduler and the optimal solution to the *reference FIFO* scheduler, which mimics the current implementation of MPTCP without any cross-layer scheduler. The video

The cross-layer scheduler can be implemented either at the transport layer or at the application layer. In the following of the paper, we consider an implementation at the application level because it is more convenient. Here are some details.

X. Corbillon, R. Aparicio-Pardo, N. Kuhn, G. Texier, and G. Simon. Cross-layer scheduler for video streaming over MPTCP. In ACM MMSys, 2016

In future work, we plan to implement DAPS in FreeBSD's CMT-SCTP stack and in Linux implementation of MPTCP to evaluate the performance gain in realistic network conditions and address the related practical issues.

N. Kuhn, E. Lochin, A. Mifdaoui, G. Sarwar, O. Mehani, and R. Boreli. DAPS: intelligent delay-aware packet scheduling for multipath transport. In IEEE ICC, 2014.

independent subflows. This approach provides transparent applications, at the cost of requiring kernel upgrade space. RepFlow can be incorporated into MPTCP multi-path support.

It is evident that RepFlow lends itself to many implementation choices. No matter the details, it is crucial that path diversity is utilized, i.e. the five-tuples of original and replicated flow have to be different (assuming

H. Xu and B. Li. RepFlow: Minimizing flow completion times for replicated flows in data centers. In IEEE INFOCOM, 2014

# Goals (of this talk)

---

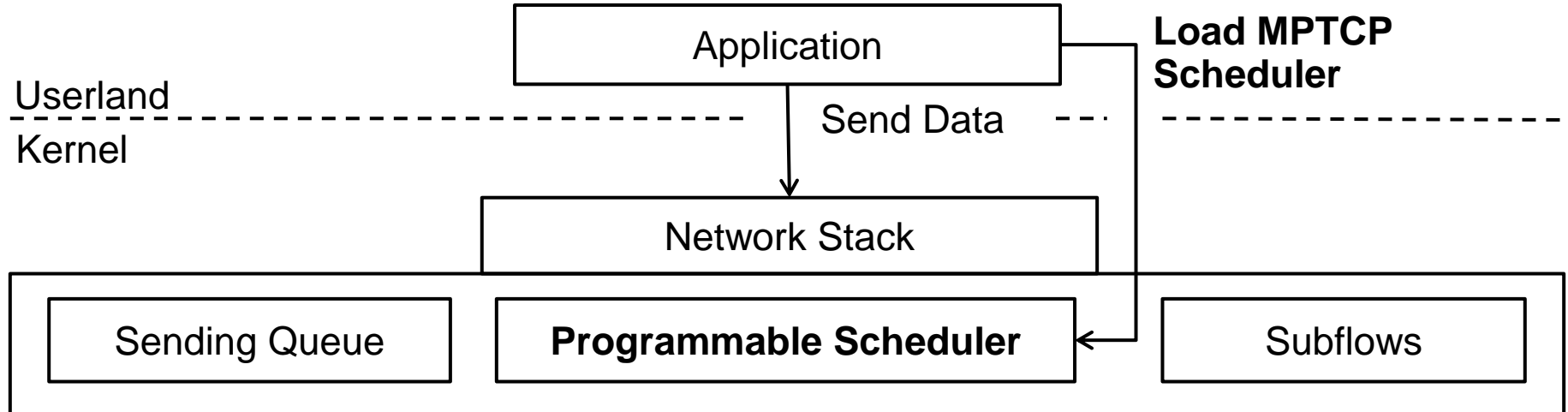


TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

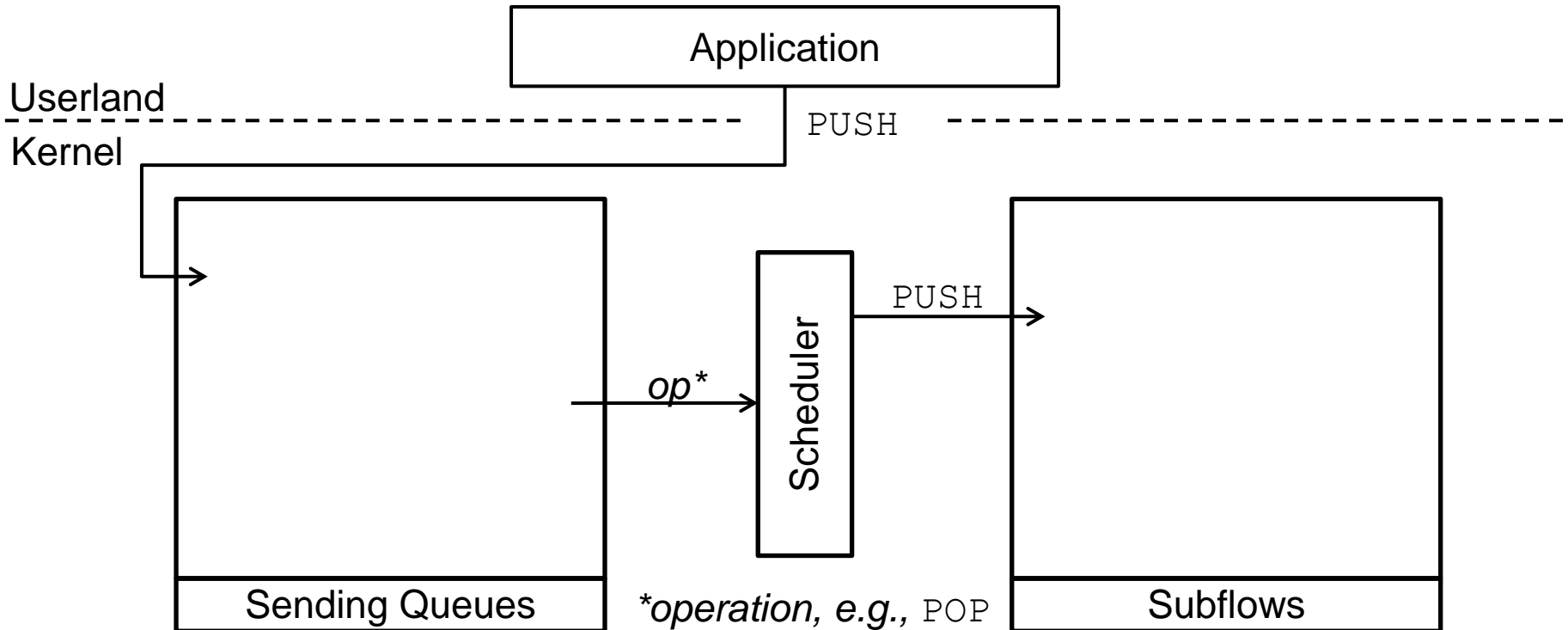
Systematically **specify** and **execute** MPTCP schedulers

**Enable application-defined** MPTCP scheduling

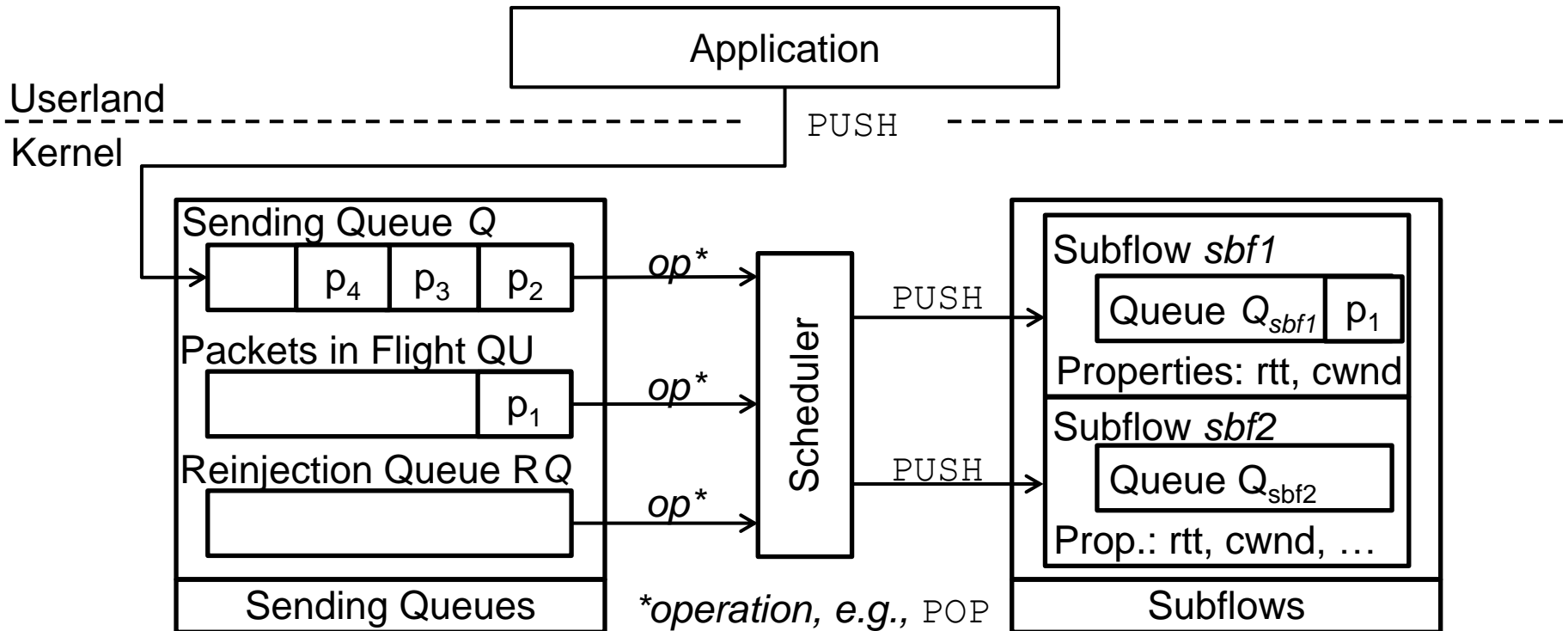
# Part I: Towards a Programmable Scheduler in the Network Stack



# Model of the Scheduler Environment



# Model of the Scheduler Environment

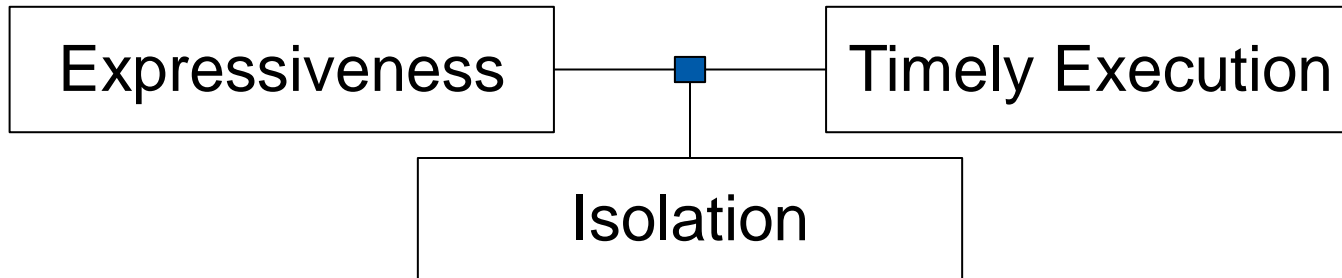


# Specifying Multipath TCP Schedulers



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

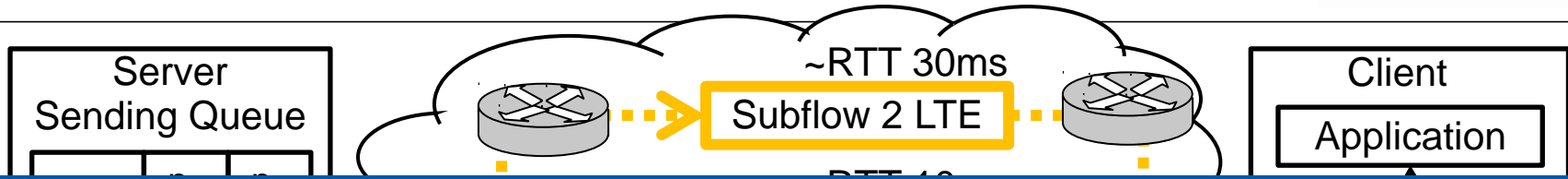
## Requirements



## Design Decisions

1. Modelled Elements as Entities: Set of Subflows with their properties, Queues of Packets with their Properties
2. Declarative Packet and Subflow Selection (Filter, Min, Max)
3. No Recursion, No Functions, Limited Loops
4. Variables with Single Assignment, Implicit and Static Type System
5. No, One, or Multiple Packets per Scheduler Execution

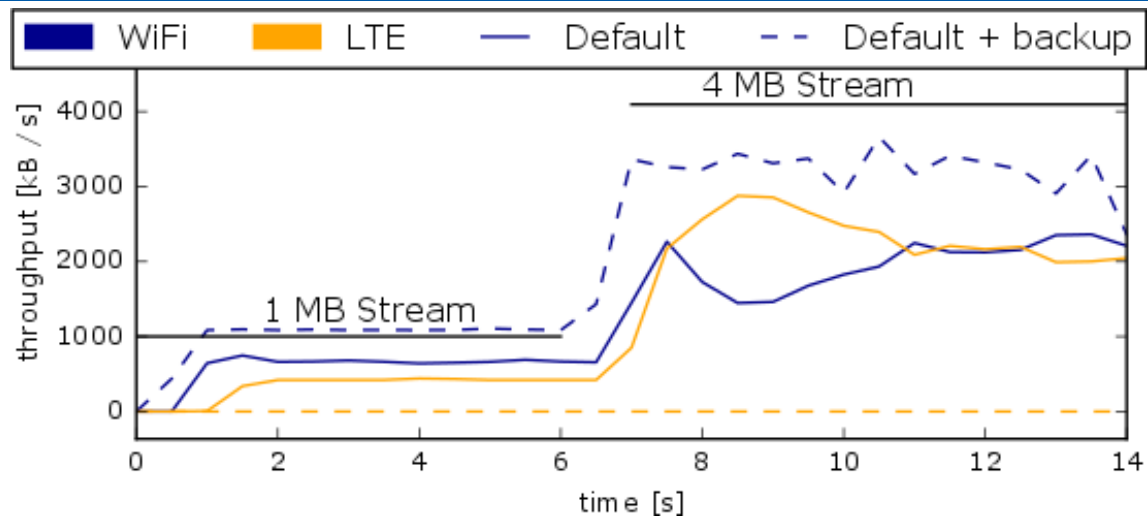
# Back to the Motivating Example



## How to use the Backup Flow?

### Constant Bitrate Stream

Home Network WiFi  
and LTE  
Germany





# Systematically Specify MPTCP Schedulers

## Domain Specific Specification Language



### Example: Preference-aware RTT-sensitive Scheduler

```
1  VAR sbfCandis = SUBFLOWS.FILTER(  
2      sbf => sbf.CWND > sbf.SKBS_IN_FLIGHT + sbf.QUEUED  
3      AND !sbf.TSQ_THROTTLED AND !sbf.LOSSY);  
4  
5  VAR backSbf = sbfCandis.FILTER(  
6      sbf => sbf.IS_BACKUP).MIN(sbf => sbf.RTT);  
7  VAR nonBackSbf = sbfCandis.FILTER(  
8      sbf => !sbf.IS_BACKUP).MIN(sbf => sbf.RTT);  
9  
10 IF (nonBackSbf.RTT_MS > 100 AND backSbf.RTT_MS < 80) {  
11     backSbf.PUSH(Q.POP());  
12 } ELSE {  
13     nonBackSbf.PUSH(Q.POP());  
14 }
```

# Systematically Specify MPTCP Schedulers

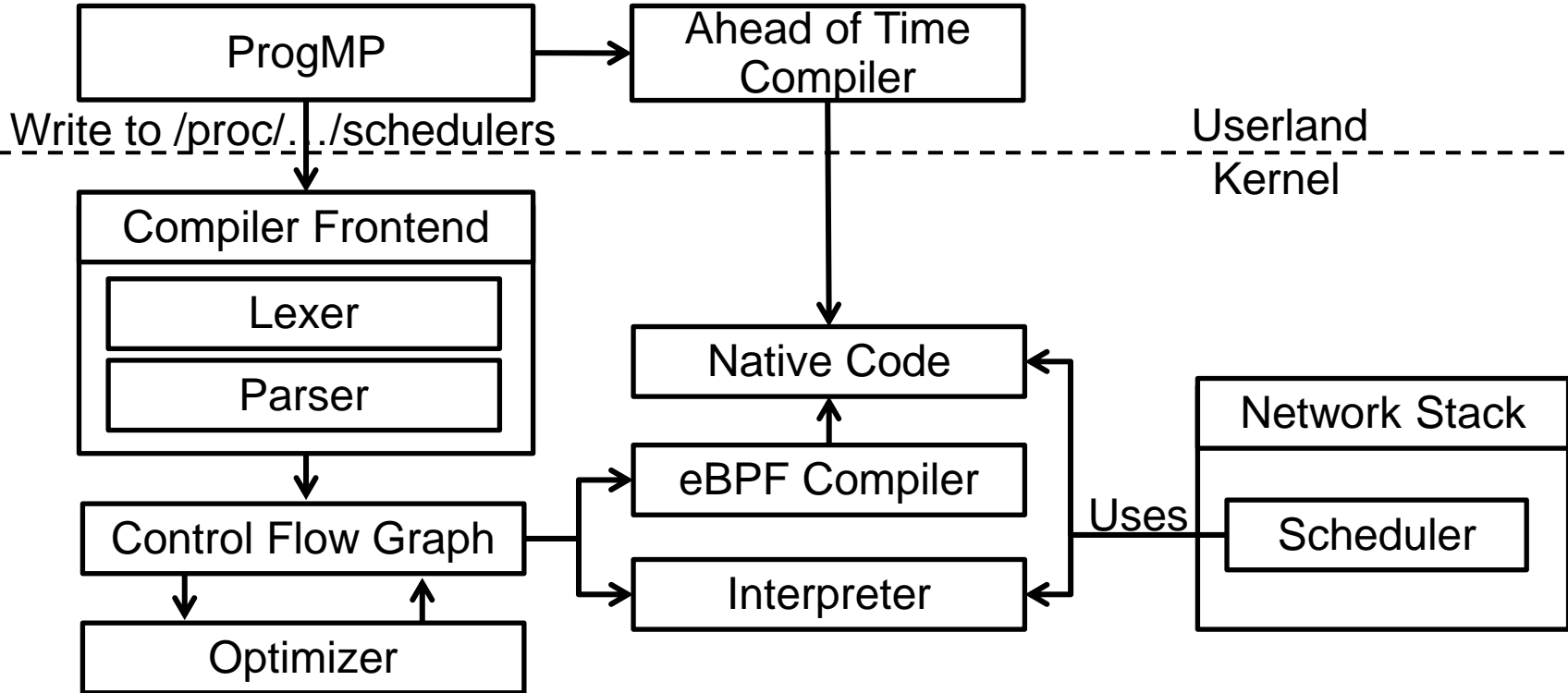
## Domain Specific Specification Language



### Example: Preference-aware RTT-sensitive Scheduler

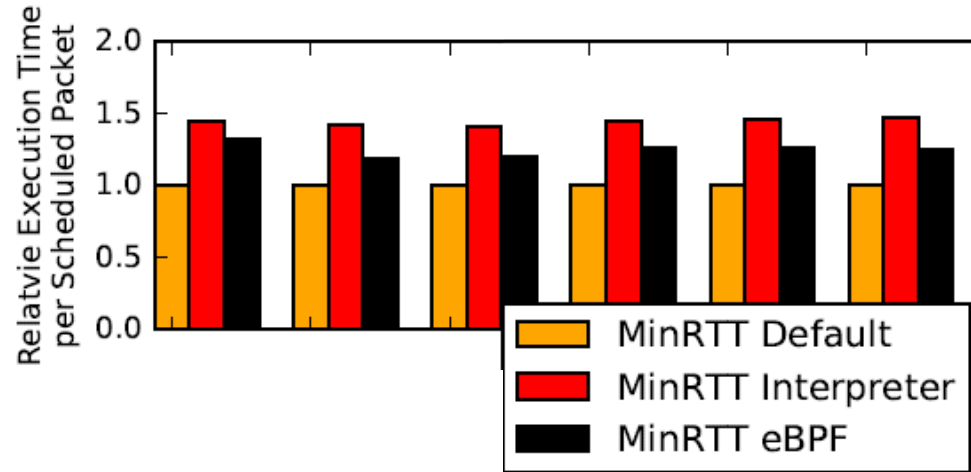
```
1  VAR sbfCandis = SUBFLOWS.FILTER(  
2      sbf => sbf.CWND > sbf.SKBS_IN_FLIGHT + sbf.QUEUED  
3          AND !sbf.TSQ_THROTTLED AND !sbf.LOSSY);  
4  
5  VAR backSbf = sbfCandis.FILTER(  
6      sbf => sbf.IS_BACKUP).MIN(sbf => sbf.RTT);  
7  VAR nonBackSbf = sbfCandis.FILTER(  
8      sbf => !sbf.IS_BACKUP).MIN(sbf => sbf.RTT);  
9  
10 IF (nonBackSbf.RTT_MS > R1 AND backSbf.RTT_MS < R2) {  
11     backSbf.PUSH(Q.POP());  
12 } ELSE {  
13     nonBackSbf.PUSH(Q.POP());  
14 }
```

# Systematically Specify and Execute MPTCP Schedulers

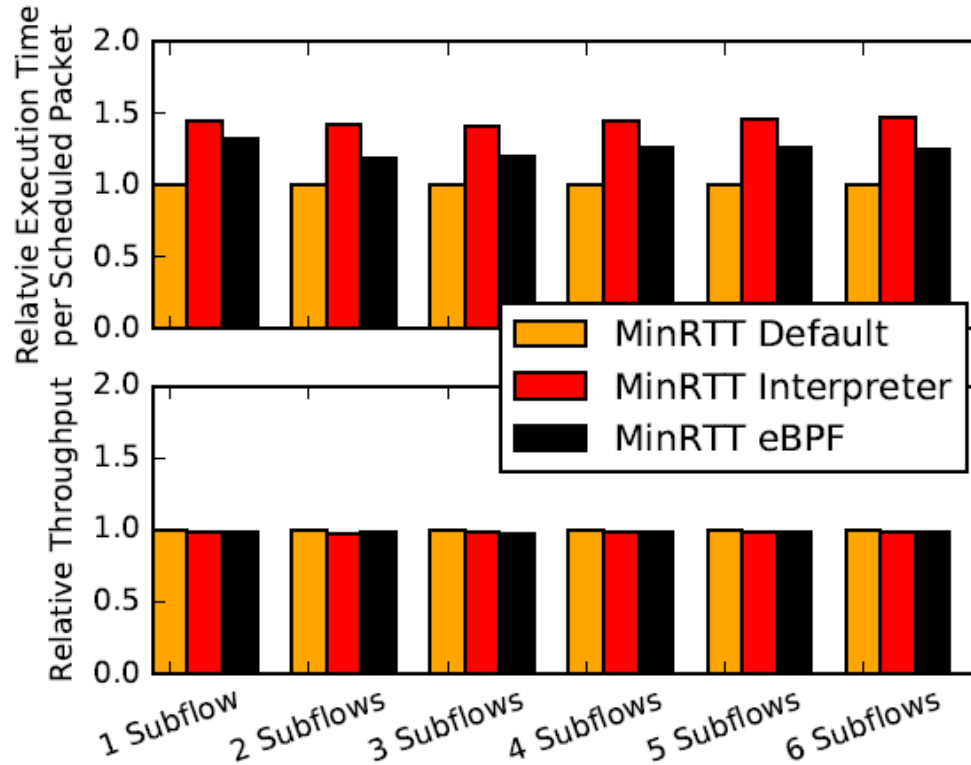


Specified schedulers are executable in the Linux Kernel

# Abstraction vs. Overhead

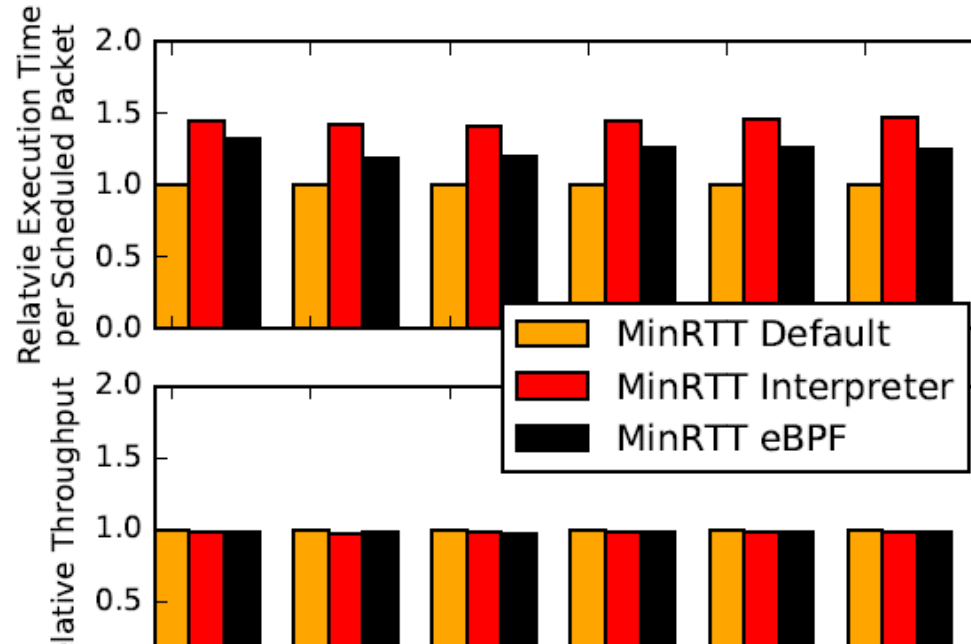


# Abstraction vs. Overhead



link saturation  
at 1 Gbit/s

# Abstraction vs. Overhead



link saturation  
at 1 Gbit/s

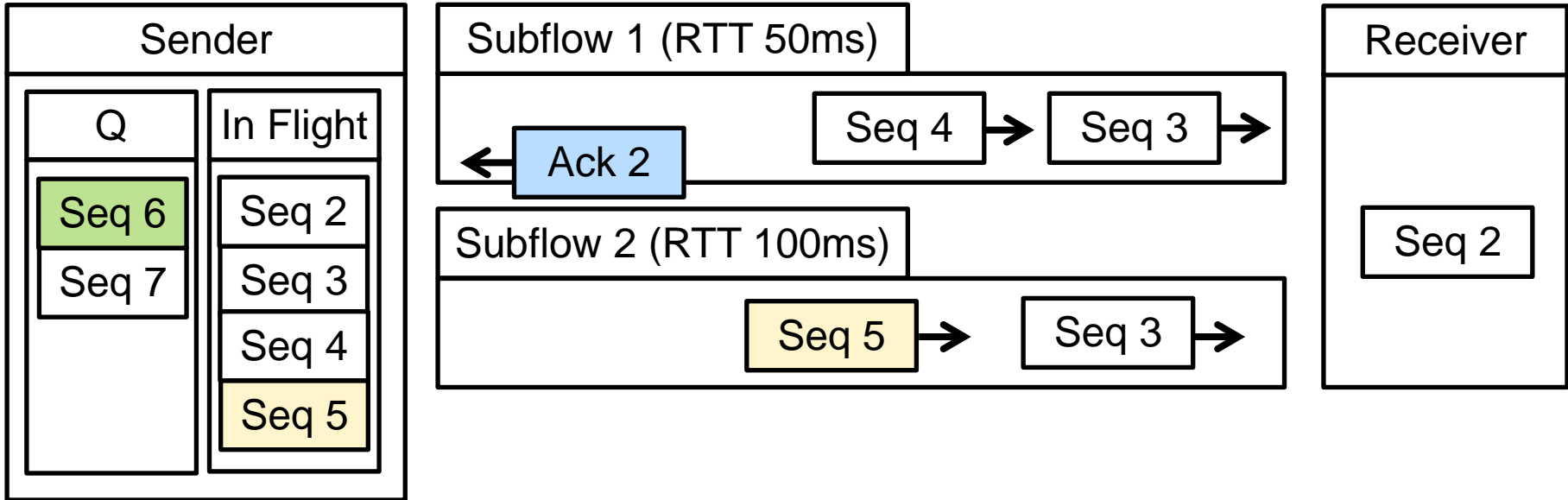
The runtime environment induces an overhead,  
which is acceptable for most application scenarios.

# Part II: Design of Novel Multipath TCP Schedulers



	Preference-aware	Application-aware	Executable
Round-trip Time-aware	✓	✓	✓
Constant Bitrate Stream Scheduling	✓	✓	✓
Redundant Scheduling			✓
HTTP/2-aware Scheduling	✓	✓	✓
More in Paper and Under Review			

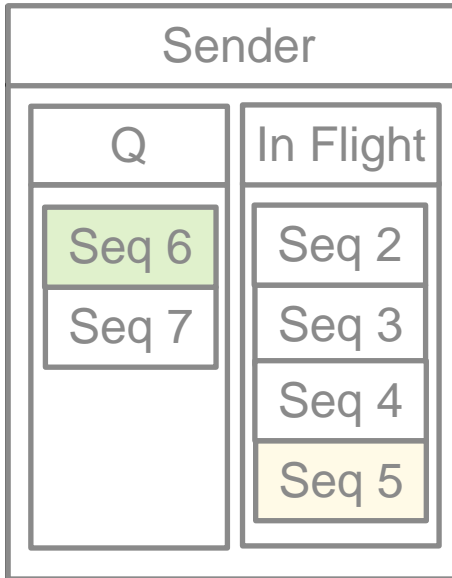
# A Close Look at Redundant Schedulers



When the acknowledgement **Ack 2** arrives at the sender:  
→ should we send the fresh packet **Seq 6** or the old packet **Seq 5**



# A Close Look at Redundant Schedulers



## Prefer Per Subflow Fresh Packets

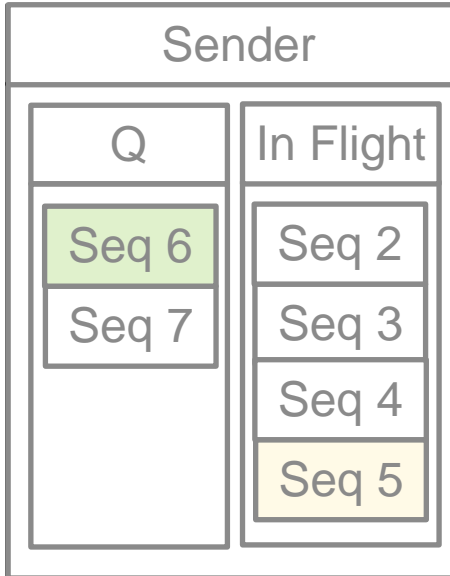
```
1 FOREACH(VAR sbf IN sbfCandidates) {
2   VAR skb = QU.FILTER(s => !s.SENT_ON(sbf).TOP);
3   IF(skb != NULL) {
4     sbf.PUSH(skb);
5   } ELSE {
6     sbf.PUSH(Q.POP());
7   }
8 }
```

## Prefer Global Fresh Packets

```
1 IF(!sbfCandidates.EMPTY) {
2   FOREACH(VAR sbf IN sbfCandidates) {
3     sbf.PUSH(Q.TOP);
4   }
5   DROP(Q.POP());
6 }
```



# A Close Look at Redundant Schedulers



## Prefer Per Subflow Fresh Packets

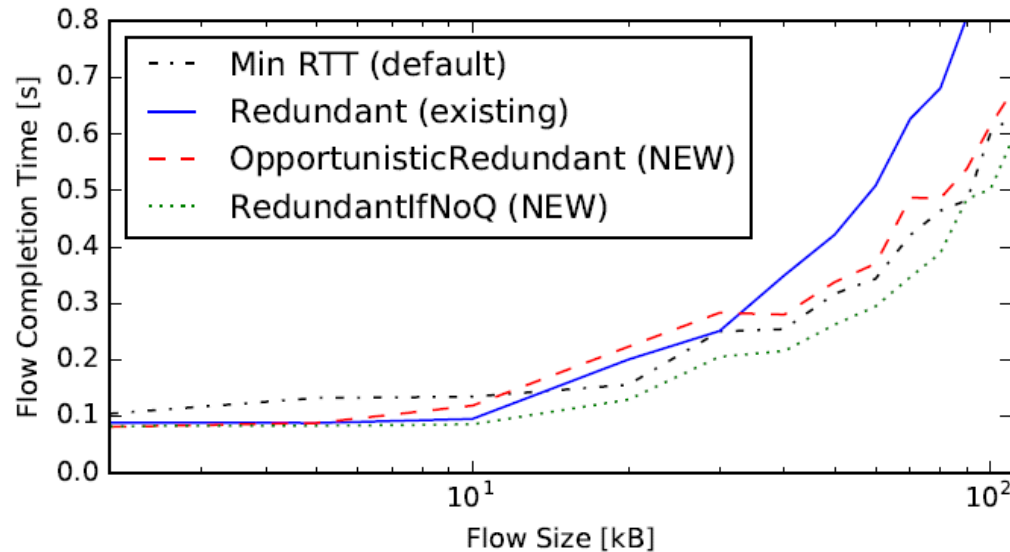
```
1 FOREACH(VAR sbf IN sbfCandidates) {
2   VAR skb = QU.FILTER(s => !s.SENT_ON(sbf).TOP);
3   IF(skb != NULL) {
4     sbf.PUSH(skb);
5   } ELSE {
6     sbf.PUSH(Q.POP());
7   }
8 }
```

## Prefer Global Fresh Packets

```
1 IF(!sbfCandidates.EMPTY) {
2   FOREACH(VAR sbf IN sbfCandidates) {
3     sbf.PUSH(Q.TOP);
```

ProgMP enables rapid specification and evaluation of schedulers

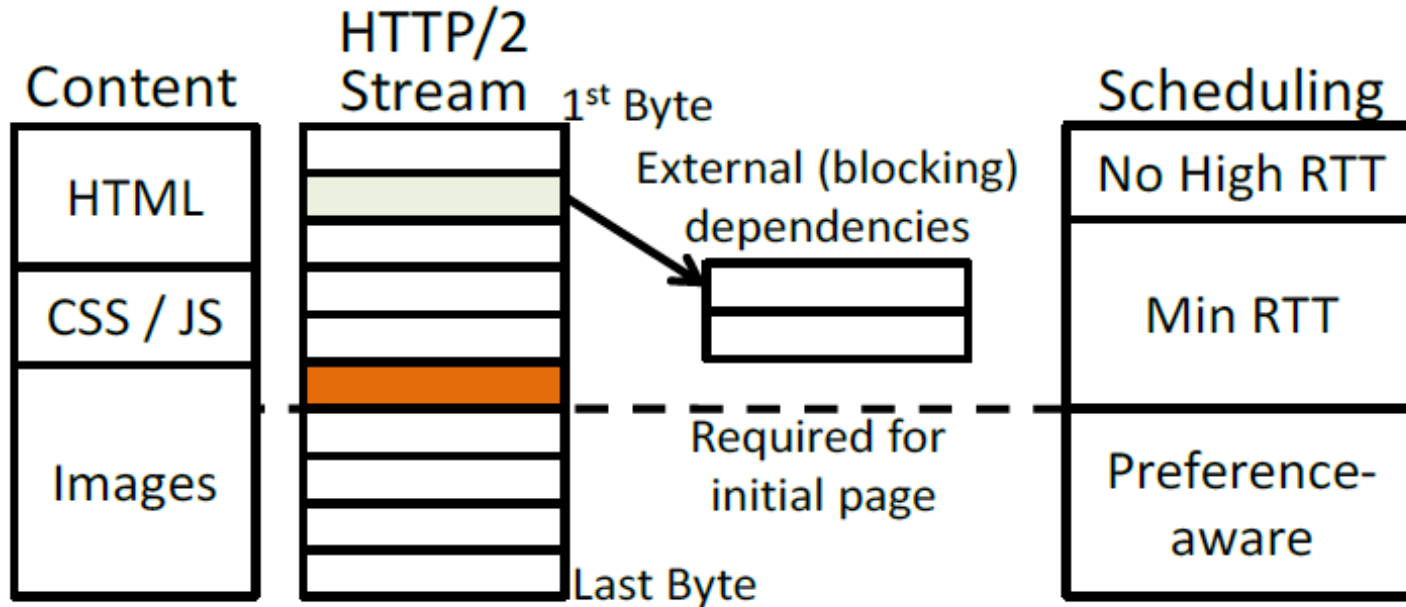
# A Close Look at Redundant Schedulers



Better

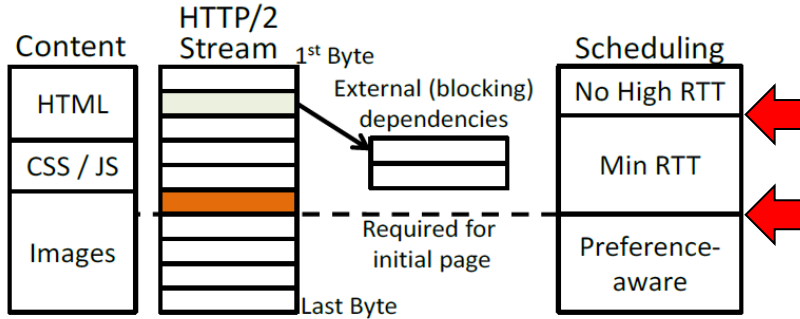
ProgMP enables novel redundant schedulers,  
which outperform established approaches.

# HTTP/2-aware Scheduling

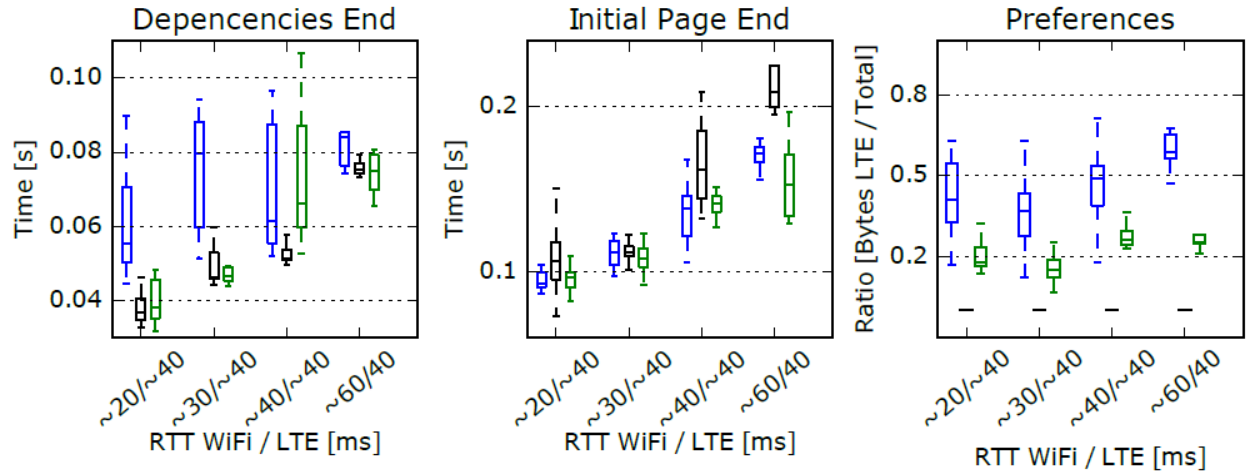


ProgMP enables HTTP/2-aware Scheduling.

# HTTP/2-aware Scheduling



Left to right boxplots: ■ Default ■ Single Path ■ HTTP/2-Aware





# Conclusion

We presented a **programming model for Multipath TCP scheduling**

- **Specification** and **execution** of MPTCP schedulers
- **Application-defined** MPTCP scheduling

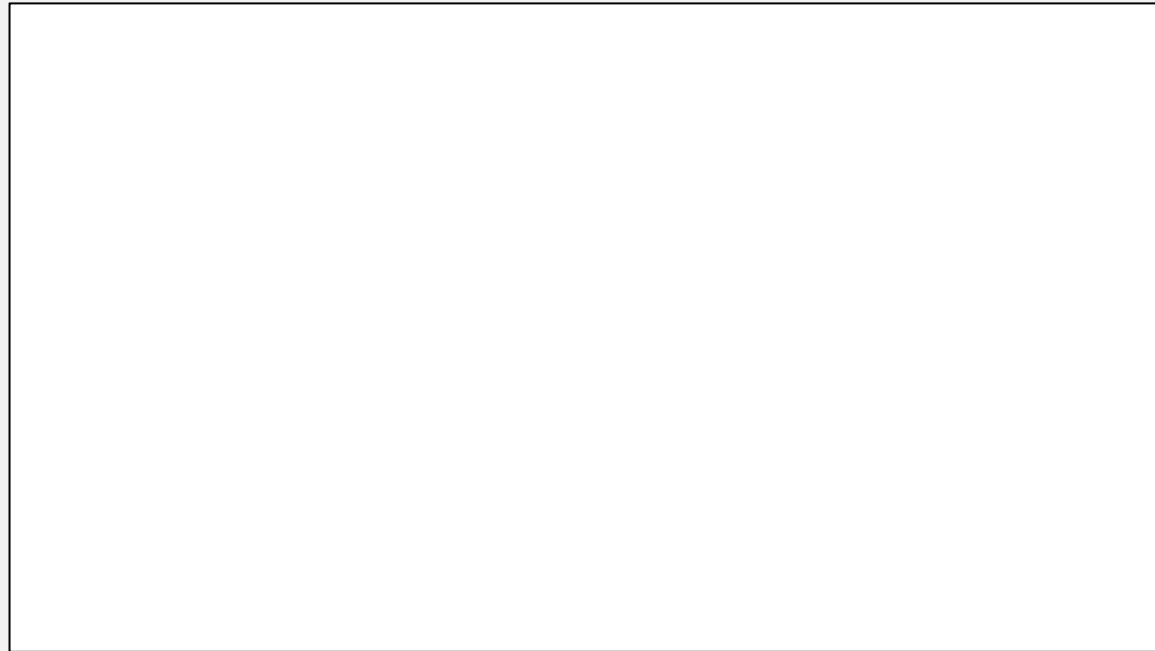
We **proposed** and **evaluated** novel MPTCP schedulers

- RTT-aware scheduler
- Constant bitrate schedulers
- Flavors of redundant schedulers
- HTTP/2-aware scheduler
- ...

# Conclusion



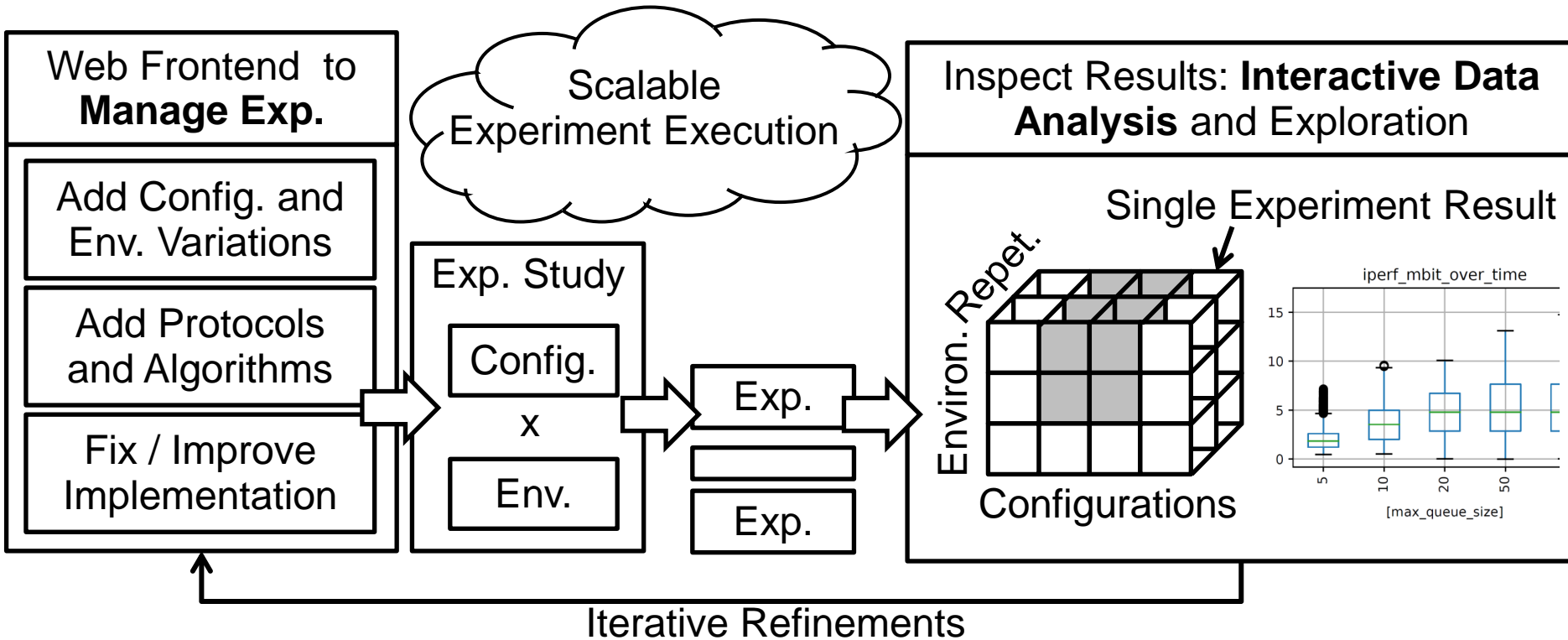
<https://progmp.net>



We pre

We pro

# How can we systematically compare and evaluate scheduler design decisions?





# How can we systematically compare and evaluate scheduler design decisions?

Web Frontend to  
**Manage Exp.**

Add Config. and  
Env. Variations

Add Protocols

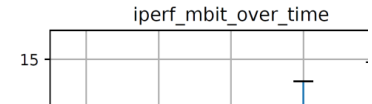
Scalable  
Experiment Execution

Exp. Study

Inspect Results: **Interactive Data  
Analysis** and Exploration

Single Experiment Result

Repet.



A Framework for the Management, the Scalable Execution  
and Interactive Analysis of Extensive Network Experiments

<https://maci-research.net>

Iterative Refinements