# Online Reprogrammable Multi Tenant Switches

Johannes Krude[1], Jaco Hofmann[2], Matthias Eichholz[2],
Klaus Wehrle[1], Andreas Koch[2], Mira Mezini[2]

[1]RWTH Aachen University, [2]Technische Universität Darmstadt

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

https://comsys.rwth-aachen.de/

ENCP '19, 2019-12-09

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS | RWTH AACHEN UNIVERSITY
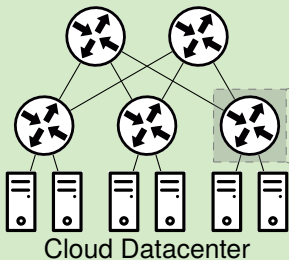
## Programmable Switch as a Service

- **On switch …**
    - ► …stateful load balancer **replaces hundreds of servers** [SilkRoad 2017]
    - ► …data aggregation **speeds up databases** [Lerner et.al. 2019, …]
    - ► …paxos **reduces coordination overhead** [NetChain 2018, …]
    - ► …key-value caching **improves throughput and latency** [NetCache 2017, …]

# Programmable Switch as a Service

- **On switch …**
  - ▸ …stateful load balancer **replaces hundreds of servers** [SilkRoad 2017]
  - ▸ …data aggregation **speeds up databases** [Lerner et.al. 2019, …]
  - ▸ …paxos **reduces coordination overhead** [NetChain 2018, …]
  - ▸ …key-value caching **improves throughput and latency** [NetCache 2017, …]

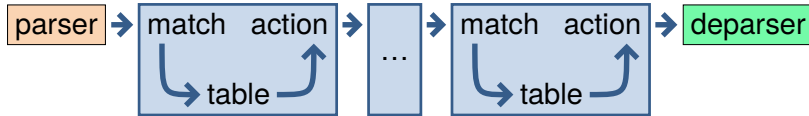**Programmable Switch as a Service**



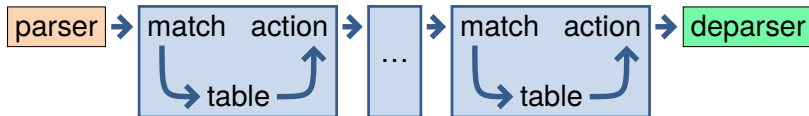| Tenant 1: **Load Balancer** | Tenant 2: **SQL group-by** |
|---|---|
| Tenant 3: **PAXOS node** | Tenant 4: **Key-Value Cache** |
| **Packet Forwarding** | |

Cloud Datacenter

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY
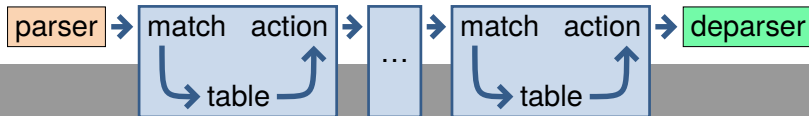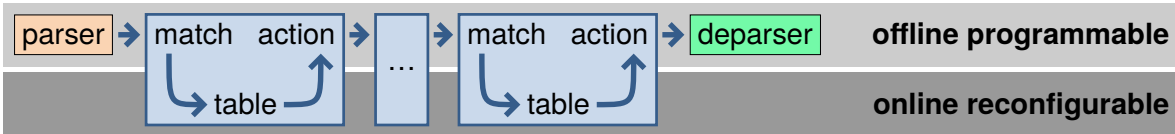
Krude et al.

- **Runs a single P4 program**

# Current Programmable Switch Pipeline [RMT 2013, Barefoot Tofino]



- **Runs a single P4 program**

Krude et al.

# Current Programmable Switch Pipeline [RMT 2013, Barefoot Tofino]



parser → match action / table → ... → match action / table → deparser **offline programmable**

**online reconfigurable**

- **Runs a single P4 program**
- **Reprogramming causes switch and network downtime**

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS | RWTH AACHEN UNIVERSITY
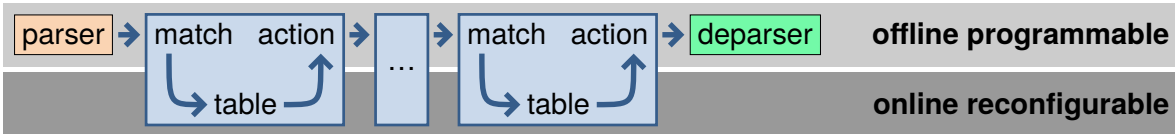
# Current Programmable Switch Pipeline [RMT 2013, Barefoot Tofino]



- **Runs a single P4 program**
- **Reprogramming causes switch and network downtime**

**We propose to modify the programmable switch architecture**
- To enable hot-pluggability of on-switch functions

Krude et al.

**Definition**

The ability to insert, modify, and remove on-switch functions without affecting other on-switch functions and packet forwarding.

Tenant 1:
**Load Balancer**
*needs high availability*

Tenant 2:
**SQL group-by**
*lifetime of seconds*

**Packet Forwarding**

Programmable Switch

## Hot-Pluggability

**Definition**

The ability to insert, modify, and remove on-switch functions without affecting other on-switch functions and packet forwarding.

**Related Work**

- **Use a dedicated switch for each application** [PPS 2019]
- **Put generalized functionality permanently onto switches** [NetAccel 2019, Ports et al. 2019]
- **Emulate P4 in Match-Action Tables** [Hyper4 2016, HyperVDP 2019]
  - ▶ Excessive Resource Consumption

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS | RWTH AACHEN UNIVERSITY

# Hot-Pluggability

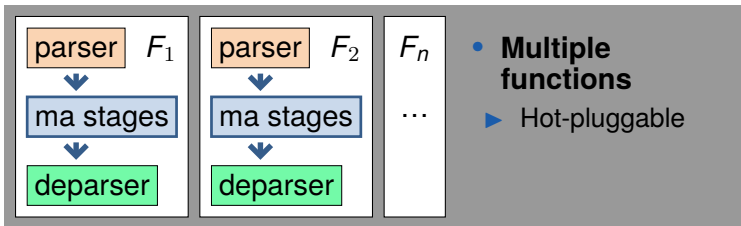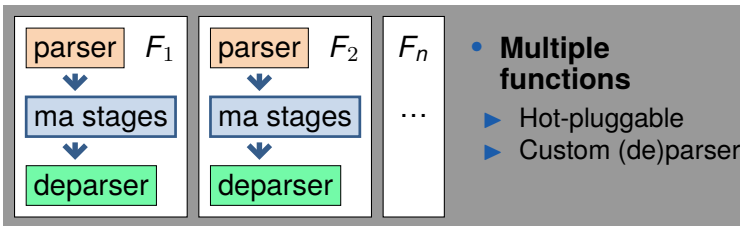| Definition | The ability to insert, modify, and remove on-switch functions without affecting other on-switch functions and packet forwarding. |
| --- | --- |

| Related Work | • **Use a dedicated switch for each application** [PPS 2019]<br>• **Put generalized functionality permanently onto switches** [NetAccel 2019, Ports et al. 2019]<br>• **Emulate P4 in Match-Action Tables** [Hyper4 2016, HyperVDP 2019]<br>  ▶ Excessive Resource Consumption |
| --- | --- |

We want: **Switch Sharing** & **On-Demand Instantiation** & **Individual Customization**

TECHNISCHE UNIVERSITÄT DARMSTADT & COMSYS RWTH AACHEN UNIVERSITY

Krude et al.

- **Multiple functions**
  - ▶ Hot-pluggable

Krude et al.

Packet Structure

| Inner Header |
|---|
| Payload |

parser $F_1$  parser $F_2$  $F_n$
↓          ↓
ma stages   ma stages   ⋯
↓          ↓
deparser    deparser

- **Multiple functions**
  - ▶ Hot-pluggable
  - ▶ Custom (de)parser

- **Function selection**
  - ▶ Max one function per packet
  - ▶ Based on VXLAN, IP dst, …

- **Multiple functions**
  - ▶ Hot-pluggable
  - ▶ Custom (de)parser

Packet Structure

| Outer Header |
|---|
| Inner Header |
| Payload |

Krude et al.

- **Function selection**
  - ▶ Max one function per packet
  - ▶ Based on VXLAN, IP dst, …

- **Multiple functions**
  - ▶ Hot-pluggable
  - ▶ Custom (de)parser

- **Post processing**
  - ▶ Updating dst/output port

Packet Structure

| Outer Header |
|---|
| Inner Header |
| Payload |

- **Function selection**
  - ▶ Max one function per packet
  - ▶ Based on VXLAN, IP dst, ...

- **Multiple functions**
  - ▶ Hot-pluggable
  - ▶ Custom (de)parser

- **Post processing**
  - ▶ Updating dst/output port

Packet Structure

| Outer Header |
| Inner Header |
| Payload |

- **Program isolation**
  - ▶ Access to only own packets
  - ▶ Limit access to outer header
  - ▶ Control plane virtualization

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY

- **We present three different possible implementations**
  - ▶ None of them yet implemented

| **Multiple Switching ASICs** | **Using FPGAs** | **An ASIC extension** |
|---|---|---|
| ✓ Easily realizable | ✓ Realizable with FPGA knowledge | ✗ To be done by switching ASIC vendors |
| ✗ No statefull functions | ✗ Reduced Throughput | ✓ High performance |

- **Alternate between two switching ASICs**



*active*   *standby*

$F_1$, $F_2$

*front-end*

- **Alternate between two switching ASICs**



*active*  *standby*

$F_1, F_2$

*front-end*

- **Alternate between two switching ASICs**

- **Alternate between two switching ASICs**



*standby*

*active*

$F_1, F_2$

$F_1, F_2, F_3$

*front-end*

- **Alternate between two switching ASICs**



*standby*     *active*

$F_1, F_2, F_3$

*front-end*

Krude et al.

- **Alternate between two switching ASICs**



*standby*

$F_1, F_3$

*active*

$F_1, F_2, F_3$

*front-end*

- **Alternate between two switching ASICs**

- **Alternate between two switching ASICs**
- **Merge functions into single program**



*active*

*standby*

$F_1, F_3$

$F_1, F_2, F_3$

*front-end*

Krude et al.

- **Alternate between two switching ASICs**
- **Merge functions into single program**
- **Compiler provides isolation**
  - ▶ Restricts access to outer headers
  - ▶ Control plane mapping from table memory to function



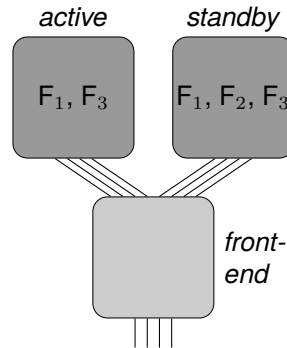*active*    *standby*

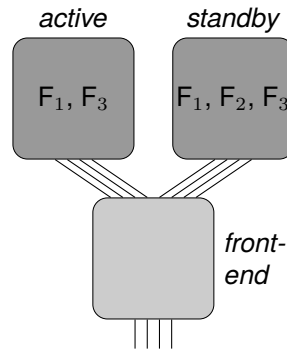$F_1, F_3$    $F_1, F_2, F_3$

*front-end*

# Multiple Programmable Switching ASICs

- **Alternate between two switching ASICs**
- **Merge functions into single program**
- **Compiler provides isolation**
  - ▶ Restricts access to outer headers
  - ▶ Control plane mapping from table memory to function



*active*   *standby*

$F_1, F_3$   $F_1, F_2, F_3$
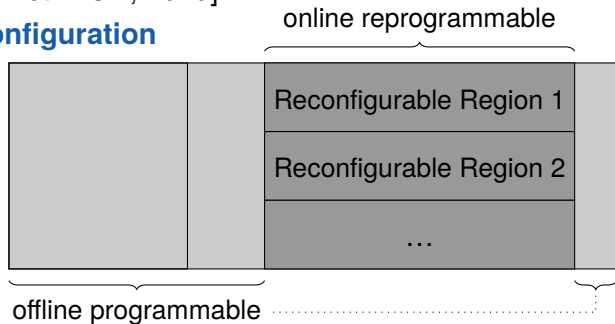
*front-end*

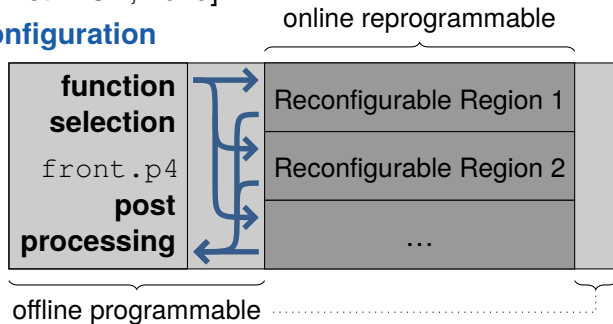| Advantages | Limitations |
|---|---|
| ✓ Based on available hardware | ✗ Problematic for statefull functions |

- **P4 can be executed on FPGAs** [P4$\rightarrow$ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**

Krude et al.

- **P4 can be executed on FPGAs** [P4→ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**

online reprogrammable

Reconfigurable Region 1

Reconfigurable Region 2

...

offline programmable

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY

- **P4 can be executed on FPGAs** [P4$\rightarrow$ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**

online reprogrammable

| **function selection** | Reconfigurable Region 1 |
| `front.p4` | Reconfigurable Region 2 |
| **post processing** | … |

offline programmable

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY

- **P4 can be executed on FPGAs** [P4→ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**
- **Isolation on FPGA level**
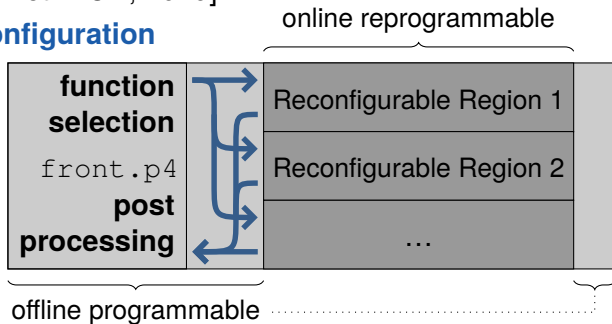  - ▶ Forward each packet to single region
  - ▶ Remove outer headers

online reprogrammable

| **function selection** | Reconfigurable Region 1 |
| `front.p4` | Reconfigurable Region 2 |
| **post processing** | ... |

offline programmable

Krude et al.

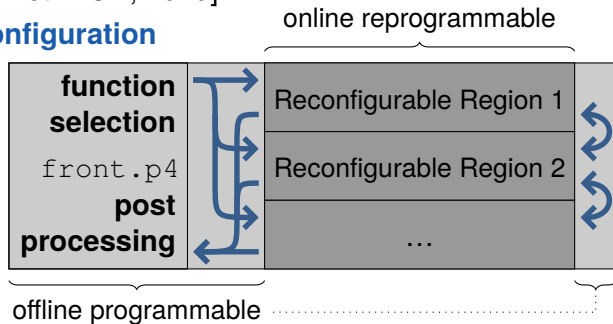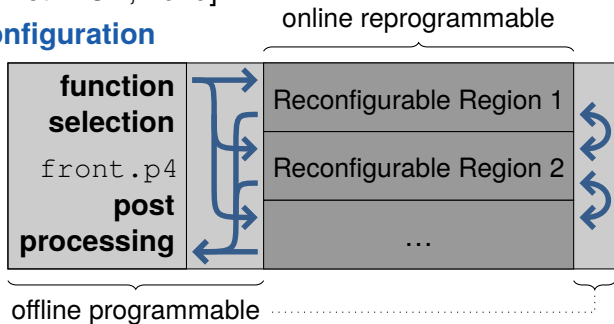TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY

- **P4 can be executed on FPGAs** [P4→ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**
- **Isolation on FPGA level**
  - ▶ Forward each packet to single region
  - ▶ Remove outer headers
- **Fixed sized reconfigurable regions**



online reprogrammable

| function selection | Reconfigurable Region 1 |
| front.p4 | Reconfigurable Region 2 |
| post processing | … |

offline programmable

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY

- **P4 can be executed on FPGAs** [P4$\rightarrow$ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**
- **Isolation on FPGA level**
  - ► Forward each packet to single region
  - ► Remove outer headers
- **Fixed sized reconfigurable regions**
  - ► Forward packets between regions
  - ► Split function across multiple regions



online reprogrammable

function selection

`front.p4`

post processing

Reconfigurable Region 1

Reconfigurable Region 2

…

offline programmable

- **P4 can be executed on FPGAs** [P4→ NetFPGA, 2019]
- **FPGAs support dynamic partial reconfiguration**
- **Isolation on FPGA level**
  - ▶ Forward each packet to single region
  - ▶ Remove outer headers
- **Fixed sized reconfigurable regions**
  - ▶ Forward packets between regions
  - ▶ Split function across multiple regions



online reprogrammable

**function selection**

`front.p4`

**post processing**

Reconfigurable Region 1

Reconfigurable Region 2

…

offline programmable

**Advantages**

- ✓ Readily available hardware
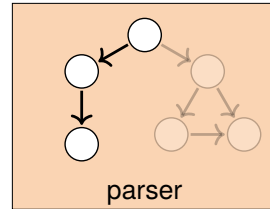- ✓ Non-reconfigured regions keep state

**Limitations**

- ✗ Limited throughput

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
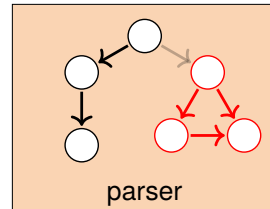  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]

Krude et al.

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
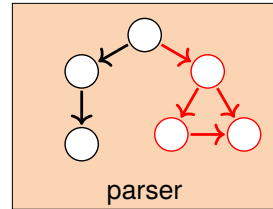- **Insertion/Removal order to avoid inconsistent states**



parser

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
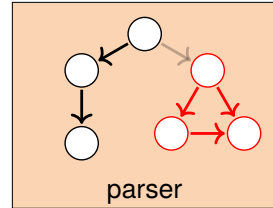- **Insertion/Removal order to avoid inconsistent states**



parser

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**
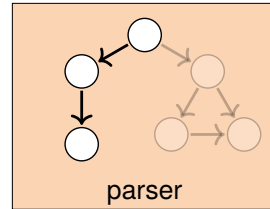


parser

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**



parser

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**
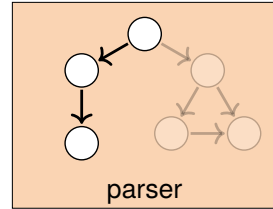


parser

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**
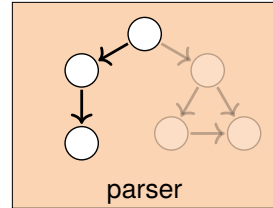- **P4 compiler tracks switch occupation**



parser

Krude et al.

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**
- **P4 compiler tracks switch occupation**
- **Isolation same as with multiple switching ASICs**



parser

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ► Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**
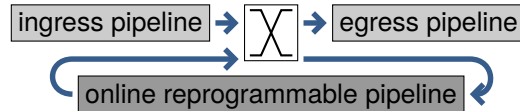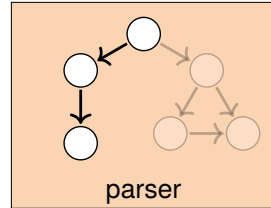- **P4 compiler tracks switch occupation**
- **Isolation same as with multiple switching ASICs**
- **Additional Pipeline for custom parser after function selection**
  - ► Some switches already have an additional pipeline
  - ► Pipelines can share resources [RMT 2013]



parser

ingress pipeline → ⋈ → egress pipeline

online reprogrammable pipeline

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS RWTH AACHEN UNIVERSITY

- **Parser, matching, actions, and deparser are stored in SRAM and TCAM**
  - ▶ Use per entry validity bit for atomic updating [CoPTUA 2004]
- **Insertion/Removal order to avoid inconsistent states**
- **P4 compiler tracks switch occupation**
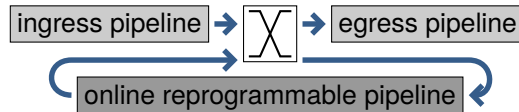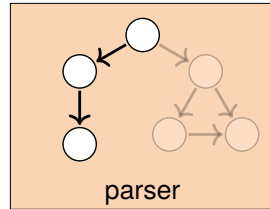- **Isolation same as with multiple switching ASICs**
- **Additional Pipeline for custom parser after function selection**
  - ▶ Some switches already have an additional pipeline
  - ▶ Pipelines can share resources [RMT 2013]



parser

ingress pipeline → ⊠ → egress pipeline

online reprogrammable pipeline

**Limitations**

✗ Needs to be done by ASIC vendors

**Advantages**

✓ Same performace as current ASICs

## Conclusion

- **Online reprogrammibility is needed for "Programmable Switches as a Service"**
- **We propose an architecture for online reprogrammibility**
  - ▶ No implementation yet

Krude et al.

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS | RWTHAACHEN UNIVERSITY

# Conclusion

- **Online reprogrammibility is needed for "Programmable Switches as a Service"**
- **We propose an architecture for online reprogrammibility**
  - ▶ No implementation yet

> **New interesting resource management questions**
> - Measuring & accounting resource usage
> - Resource allocation
> - Avoiding resource fragmentation
> - …

TECHNISCHE UNIVERSITÄT DARMSTADT & COM SYS | RWTH AACHEN UNIVERSITY