# Using P4-INT on Tofino for measuring device performance characteristics in a network lab

Speakers              Marcel Beausencourt (B.Eng), Max Julius Bode (B.Eng.)

Course of Studies    MA Information and Communication Engineering

FB 1 – Energy and Information

University of Applied Sciences Berlin (HTW Berlin)

Prof. Dr. Thomas Scheffler

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

# Agenda

1) Short introduction to P4

2) Inband Network Telemetry (INT)

3) Future prospects

# What's P4?

- Hardware-based programming of algorithms for packet processing

  - Create own packet processing behavior/algorithms

  - Manipulate every bit individually ($\rightarrow$ headers)

- Successor of Open-Flow

- P4 Goals

  I.   Reconfigurability

  II.  Protocol independence

  III. Portability (architecture independence)

htw.
Hochschule für Technik
und Wirtschaft Berlin
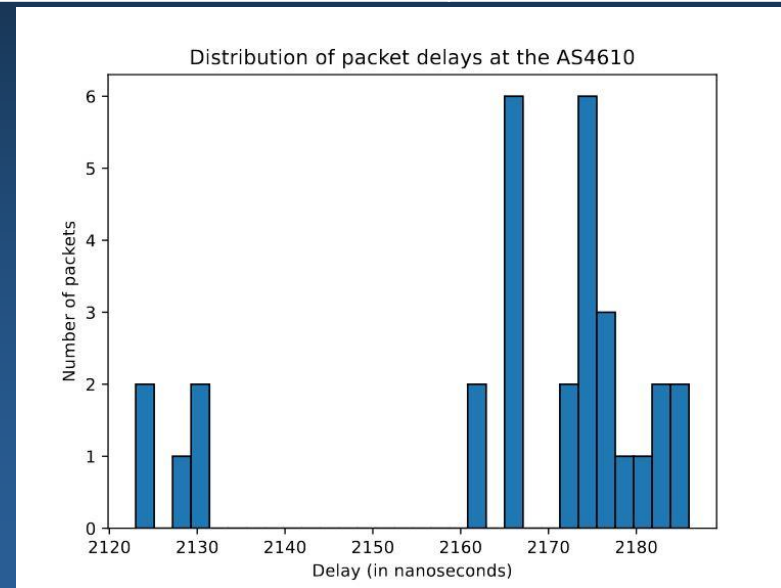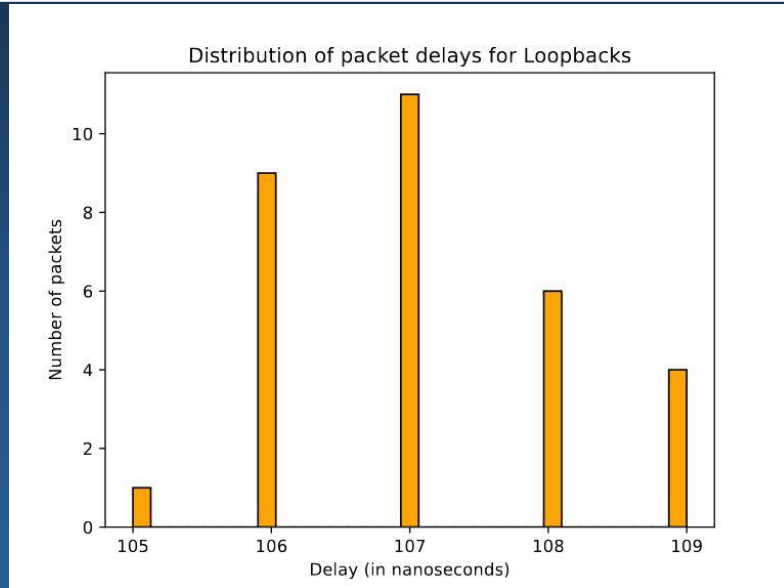University of Applied Sciences

# What can P4 do?

- Design and implementation of our own protocols

    - INT Header

    - Clone Header

- Pull unused protocols, optimize limited resources (TCAMs, ...)

- Building networks for customers that are using proprietary protocols
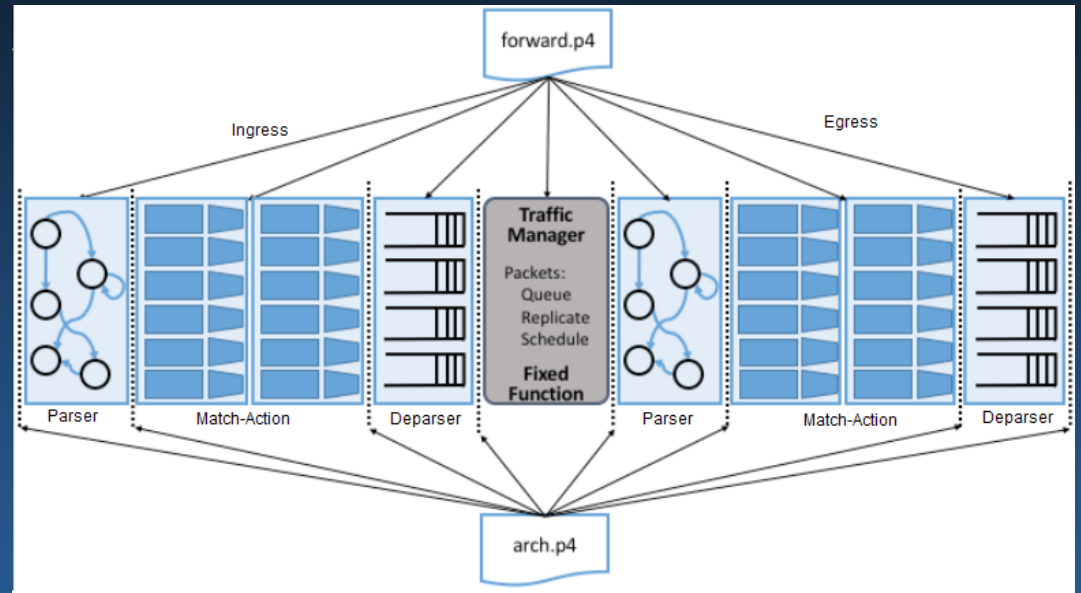
- Research applications

# First measurements with INT

| | Mean Delay in ns | Std. Dev. in ns |
|---|---|---|
| Loopback | 106.96 | 1.06 |
| AS4610 | 2128.96 | 16.58 |



Distribution of packet delays for Loopbacks



Distribution of packet delays at the AS4610

# How does P4 in Tofino work?

- Implement needed protocols
  - Parser (= State Machine)

- Implement your algorithms and tables
  - Match-Action

- Write the changes into the headers
  - Deparser (Checksums, MAC change, added headers → e.g. INT)
  - Send the packet ;)

- Payload isn't modified



Original Source: https://sdn.systemsapproach.org/

Hochschule für Technik
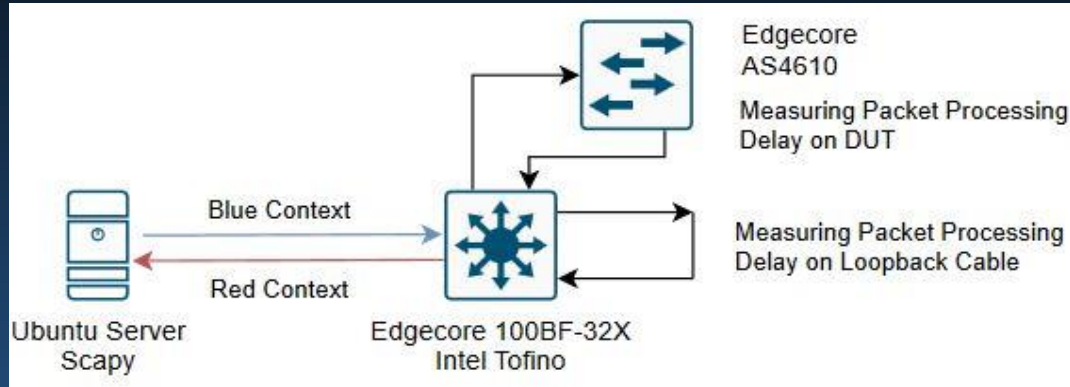und Wirtschaft Berlin
University of Applied Sciences
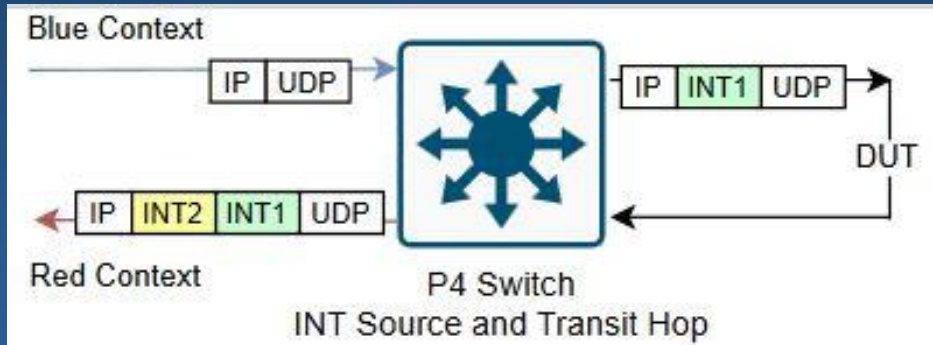
# INT Framework

- Framework made by the P4-Language-Consrtium → Current version: 2.1

- Collect & report network state by the data plane (no interaction with the control plane is needed)

- INT-XD

  – No packet modifications

- INT-MX

  – Packets carry only instructions for the switches

- INT-MD

  – Supports instructions and metadata inside of network packets

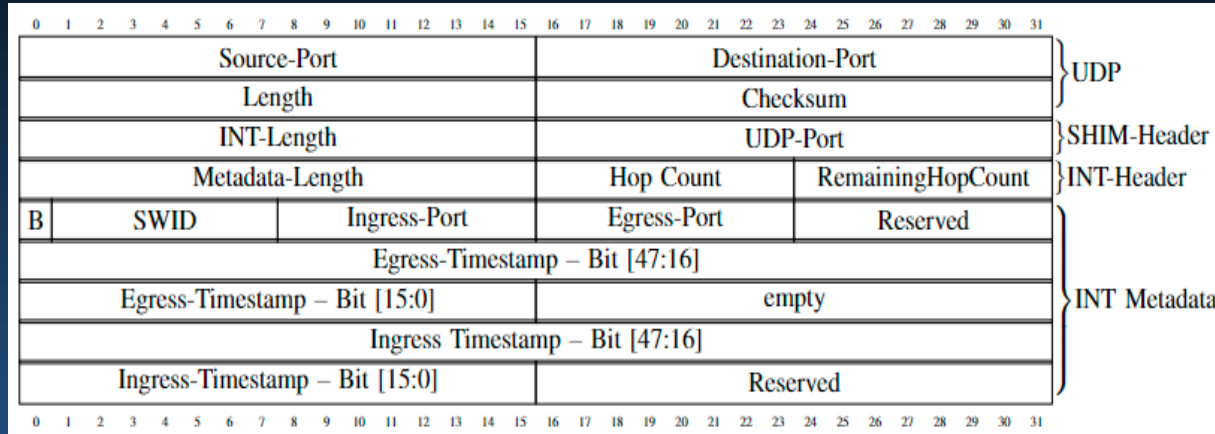# Measuring delays with P4-INT



Measuring topology



Adding INT headers into the packets

# INT Header



INT Header format with 48-bit timestamps and P4 code implementation

```
121  /***********INT Headers*********************************************/
122  header int_shim_h {
123          bit<16> len;
124          bit<16> port;
125  }
126
127  header int_header_h {
128          bit<16> hop_ml;
129          bit<8> nhop;
130          bit<8> remaining;
131  }
132
133  #ifdef MAC_STAMPING
134  header int_metadata_h {
135      bit<1> bos;
136      bit<7> swid;
137      bit<8> ig_port;
138      bit<8> eg_port;
139      bit<8> reserved_1;
140      bit<48> ig_mac_tstamp;
141      bit<16> reserved_2;
142  }
143  #else
144  header int_metadata_h {
145      bit<1> bos;
146      bit<7> swid;
147      bit<8> ig_port;
148      bit<8> eg_port;
149      bit<8> reserved_1;
150      bit<48> eg_mac_tstamp;
151      bit<16> empty;
152      bit<48> ig_mac_tstamp;
153      bit<16> reserved_2;
154  }
155  #endif
```

# What can be measured?

- Any networking device

    - Switches, Router

    - WDMs

- Effects/delays of different features (ACLs, Rate-Limits, …)

- Edge Cases → Load-Testing of networks running very unusual protocol stacks

# Traffic Analyzation I

# Traffic Analyzation II

```
> Frame 1: 120 bytes on wire (960 bits), 120 bytes captured (960 bits)
> Ethernet II, Src: HewlettP_ec:f4:40 (00:26:55:ec:f4:40), Dst: Broadcast (ff:ff:ff:ff:ff:ff)
> Internet Protocol Version 4, Src: 10.1.1.1, Dst: 10.1.1.2
> User Datagram Protocol, Src Port: 4321, Dst Port: 54321
∨ INT-SHIM Protocol Data
      Telemetry Length: 42
      UDP Port: 7
∨ INT-HEADER Protocol Data
      Hop ML: 34
      Number of Hops: 2
      Remaining Hops: 0
∨ INT-METADATA Protocol Data - Packet No. 2
      0... .... = BOS bit: 0
      .000 0000 = Switch-ID: 0
      Ingress Port: 133
      Egress Port: 135
      Egress Mac Timestamp: 479459370260
      Empty: 0
      Ingress Mac Timestamp: 479459369874
∨ INT-METADATA Protocol Data - Packet No. 1
      1... .... = BOS bit: 1
      .000 0000 = Switch-ID: 0
      Ingress Port: 134
      Egress Port: 132
      Egress Mac Timestamp: 479459367498
      Empty: 0
      Ingress Mac Timestamp: 479459367133
```

# Packet-Duplication / Clone-Header

| 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 |
|---|
| Anzahl der gewünschten Clones |
| # of Clones |
| Secure String |
| MC/BC Option · Length of Payload · Version # · Future Use |
| IP for MC/BC |
| Netmask for MC/BC |
| Standard Test String |
| 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 |

- Makes X clones of a packet → 1  field

- Compares wanted number of packets with real number of packet → "# of Clones" = Counter

- Relies on Ethernet and IPv4 → Clone = Layer4

- Secure String → Authorization

- MC/BC Option → enables Multi-/Broadcast
  - In combination with IP and Netmask for MC/BC

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

# Goals and future prospects

- Generation of flexible test-traffic >= 100 GBit/s

- Suited for analyzing and learning network protocols

  - Header structures and implementation

  - Basic algorithms

  - See what a network device has to do to forward a packet

- Front-End Application for INT, Scapy

htw.
Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences

# Thank you for your attention!

# Contact

- Marcel Beausencourt
  - [Marcel.beausencourt@student.htw-berlin.de](mailto:Marcel.beausencourt@student.htw-berlin.de)
- Max Julius Bode
  - [Max.bode@student.htw-berlin.de](mailto:Max.bode@student.htw-berlin.de)

Prof. Dr. Thomas Scheffler
  - [Thomas.scheffler@htw-berlin.de](mailto:Thomas.scheffler@htw-berlin.de)
- Git Repo incl. Bachelor thesis for getting an easier start into P4
  - https://github.com/Selltowitz/p4

Hochschule für Technik
und Wirtschaft Berlin
University of Applied Sciences