Chair of Network Architectures and Services
School of Computation, Information and Technology
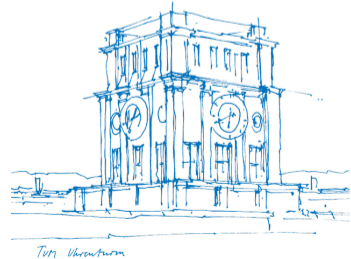Technical University of Munich

ΠΠΠ

# Honey for the Ice Bear – Dynamic eBPF in P4

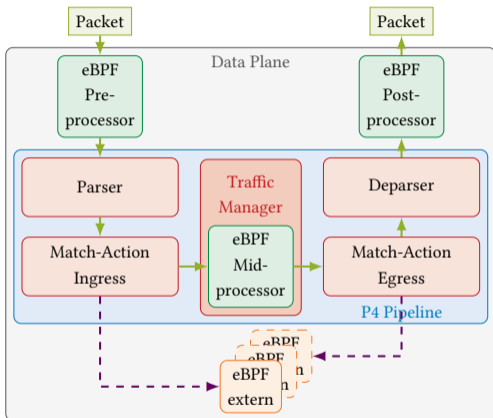**Manuel Simon,** Henning Stubbe, Sebastian Gallenmüller, Georg Carle

Sunday 4th August, 2024

Chair of Network Architectures and Services
School of Computation, Information and Technology
Technical University of Munich

TUM Uhrenturm

# Motivation

- Interrupt-free, dynamic updates increase network resilience
  - ⇒ application migration
  - ⇒ tenant-specific processing

- P4 and eBPF are well-established languages for programmable packet processing
  - ⇒ P4: restricted, simple language, optimized for high performance
  - ⇒ eBPF: JIT compiled, more high-level language features

- Both languages bring advantages for specific use-cases
  - ⇒ eBPF programs as well-defined API for P4 externs to *extend* functionality

# Dynmiac eBPF in P4

- Extension of P4 pipeline with updatable eBPF modules
  - Fixed position
  - Extern
- Allows runtime re-programmability
  - Exchange using pre-compiled **byte code**
  - JIT compiled to **machine code**
- Extends P4 functionality with well-defined API

# Dynamic modes

### Static

- Fixed, non-changeable functionality

### Pre-defined

- Pre-implemented, fixed set of functionality
- Defined before *initialization*, switchable during *runtime*

### Extensible

- New functionality is sent as source or byte code
- JIT compiled and bound during *runtime*

# Related Work

Reprogrammable P4:

- Das et al., ActiveRMT [1]: Instruction set in P4 allowing changegable functionality
- Xing et al., FlexCore [6]: Runtime partial reprogrammable switch architecture
- Feng et al., In-situ Programmability Data Plane [3]: Switch architecture and reconfigurable P4 (rP4) for runtime updates

  $\Rightarrow$ single-language P4 approaches

P4/eBPF:

- P4 to eBPF [4]: Translation of P4 program to eBPF [4]

- Implementation for software target *T4P4S* [5]
- eBPF execution using DPDK `rte_bpf` library
  - batched tx/rx eBPF callback execution for *fixed* position
  - non-batched execution for *flexible* externs
- User space eBPF execution
- Optional BLAKE3-based MACs ensuring authenticity of code updates

## Overhead of eBPF execution at different positions (Throughput)



Three programs for basic overhead:

- *dummy*: returns 0
- *filter*: filters for one UDP port and IP address
- *change*: changes a header field

Overhead of eBPF execution at different positions (modeled per-packet CPU cycles)



Three programs for basic overhead:

- *dummy*: returns 0
- *filter*: filters for one UDP port and IP address
- *change*: changes a header field

Cost model:

$$C = \frac{f_{CPU}}{r_{testbase}} - \frac{f_{cpu}}{r_{baseline}}$$

Median costs of dynamic updates—ten runs (100 Mbit/s)



⇒ Update of fixed-position functionality more expensive

⇒ Dynamic eBPF byte code installation at reasonable costs

⇒ Authentication possible

- eBPF offers fixed API for P4 externs
- eBPF hardware offloading solutions exist
- eBPF execution within P4 allows additional applications
- Functionality can be updated during runtime ( 200 µs)

Read the paper if you want more information about:

- Security considerations
- Discussion of different processor positions
- Detailed analysis of program change

# Bibliography

[1] R. Das and A. C. Snoeren.
Memory management in activermt: Towards runtime-programmable switches.
In *Proceedings of the ACM SIGCOMM 2023 Conference*, ACM SIGCOMM '23, page 10431059, New York, NY, USA, 2023. Association for Computing Machinery.

[2] P. Emmerich, S. Gallenmüller, D. Raumer, F. Wohlfart, and G. Carle.
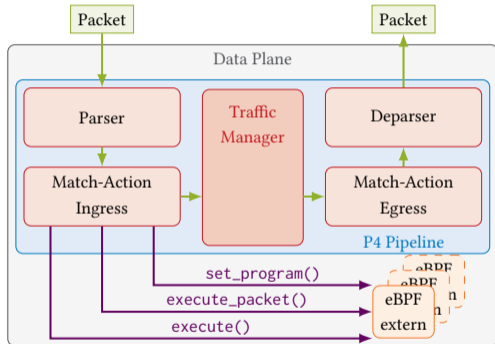Moongen: A scriptable high-speed packet generator.
In K. Cho, K. Fukuda, V. S. Pai, and N. Spring, editors, *Proceedings of the 2015 ACM Internet Measurement Conference, IMC 2015, Tokyo, Japan, October 28-30, 2015*, pages 275–287. ACM, 2015.

[3] Y. Feng, Z. Chen, H. Song, W. Xu, J. Li, Z. Zhang, T. Yun, Y. Wan, and B. Liu.
Enabling in-situ programmability in network data plane: From architecture to language.
In A. Phanishayee and V. Sekar, editors, *19th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2022, Renton, WA, USA, April 4-6, 2022*, pages 635–649. USENIX Association, 2022.

[4] p4lang.
GitHub p4c/backends/ebpf - eBPF Backend, 2024.
Last accessed: 2024-05-24.

[5] P. Vörös, D. Horpácsi, R. Kitlei, D. Leskó, M. Tejfel, and S. Laki.
T4P4S: a target-independent compiler for protocol-independent packet processors.
In *IEEE 19th International Conference on High Performance Switching and Routing, HPSR 2018, Bucharest, Romania, June 18-20, 2018*, pages 1–8. IEEE, 2018.

[6] J. Xing, K.-F. Hsu, M. Kadosh, A. Lo, Y. Piasetzky, A. Krishnamurthy, and A. Chen.
Runtime programmable switches.
In *19th USENIX Symposium on Networked Systems Design and Implementation (NSDI 22)*, pages 651–665, Renton, WA, Apr. 2022. USENIX Association.
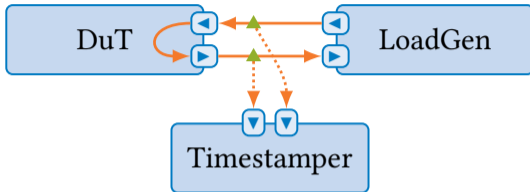
**Fixed position**

- *Pre*-, *Mid*-, or *Postprocessor*
- Processes every packet
- Access to whole packet
- Potentially easier implementation
- E.g., prefilter, preprocessing, hashing/crypto

Dynamic eBPF in P4



**Extern**

- Flexible position as P4 extern
- Conditional execution
- Return value usable
- Access to whole packet *or* restricted to selected header fields

### Setup



#### DuT

- Intel Xeon D-1518  2.2 GHz, 32 RAM
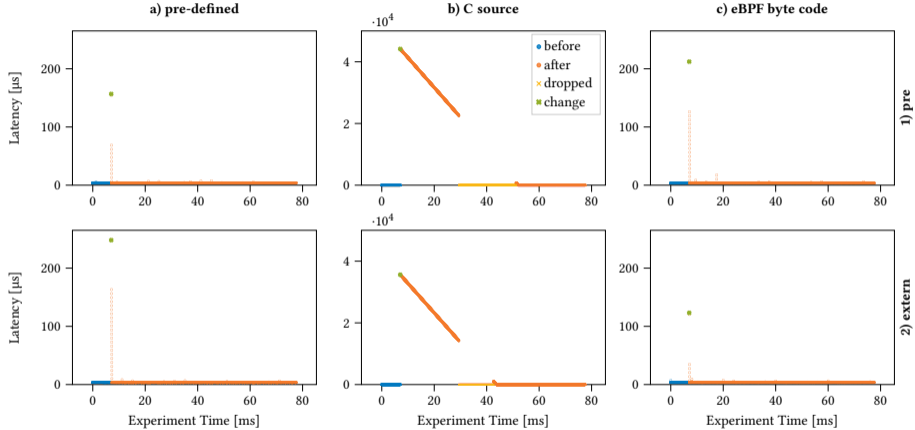- Latency optimized *T4P4S* → batch size of one

#### LoadGen

- MoonGen [2] is used to generate traffic
- Packet size 84 B

#### Timestamper

- Packet streams duplicated using optical splitter
- Timestamps each packet incoming packet
- Resolution: 12.5 ns

## Costs of dynamic updates—single run (100 Mbit/s)



⇒ eBPF byte code swapping during runtime possible without packet loss