



Determination of Throughput Guarantees for Processor-based SmartNICs

Johannes Krude, Jan R uth,
Daniel Schemmel, Felix Rath,
Johannes-Heorh Folbort, Klaus Wehrle

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



CoNEXT '21, December 7–10, 2021

⚠ packet drops and congestion testing with traffic traces gives no guarantee

Match-Action-based Programmable Switches

- + If a program compiles, it runs at ~1 pkt/cycle**
- Difficult to program** (subset of P4, no loops, few sequential operations)

⚠ packet drops and congestion testing with traffic traces gives no guarantee

Match-Action-based Programmable Switches

- + If a program compiles, it runs at ~ 1 pkt/cycle
- Difficult to program (subset of P4, no loops, few sequential operations)

Processor-based SmartNICs

- + Are easier and more freely programmable (C, BPF/XDP)
- Performance varies and is not obvious [our Netsoft2019 paper]

⚠ packet drops and congestion testing with traffic traces gives no guarantee

Match-Action-based Programmable Switches

- + If a program compiles, it runs at ~1 pkt/cycle
- Difficult to program (subset of P4, no loops, few sequential operations)

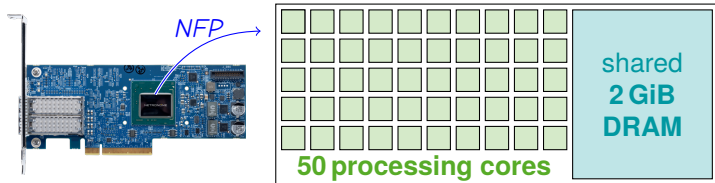
Processor-based SmartNICs

- + Are easier and more freely programmable (C, BPF/XDP)
- ~~Performance varies and is not obvious~~ [our Netsoft2019 paper]

Our Contribution: A tool to calculate the throughput of a program
...while developing a program. ...as part of regression tests.

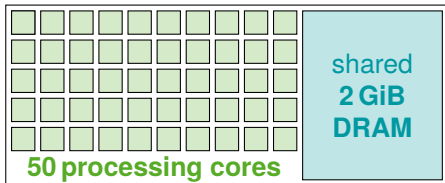
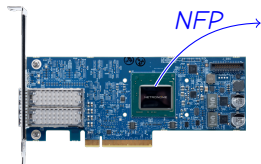
- **Netronome Agilio CX 2x40 GbE**

- ▶ Netronome Flow Processor (NFP)
- ▶ BPF/XDP programs compiled to NFP bytecode



- **Netronome Agilio CX 2x40 GbE**

- ▶ Netronome Flow Processor (NFP)
- ▶ BPF/XDP programs compiled to NFP bytecode



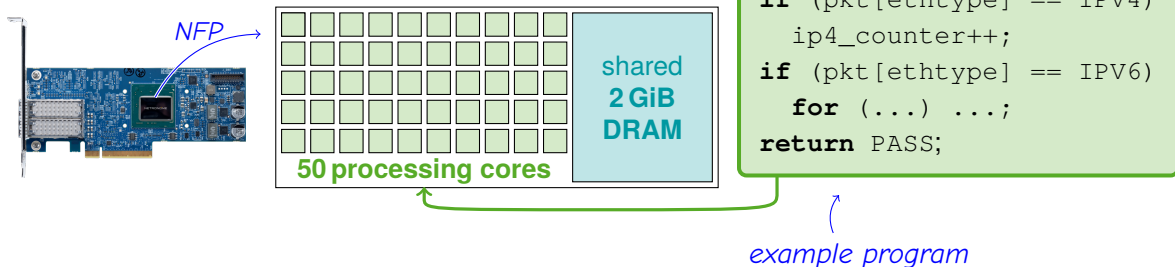
```
if (pkt.size < 100)
    return DROP;
if (pkt[ethertype] == IPV4)
    ip4_counter++;
if (pkt[ethertype] == IPV6)
    for (...) ...;
return PASS;
```

example program

Throughput Analysis

- **Netronome Agilio CX 2x40 GbE**

- ▶ Netronome Flow Processor (NFP)
- ▶ BPF/XDP programs compiled to NFP bytecode



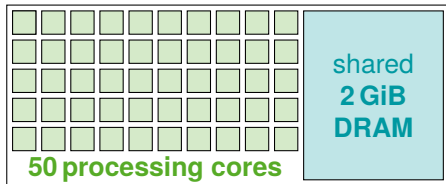
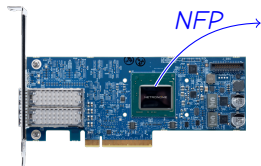
- **Challenges**

- ▶ The throughput depends on the executed instructions

Throughput Analysis

- **Netronome Agilio CX 2x40 GbE**

- ▶ Netronome Flow Processor (NFP)
- ▶ BPF/XDP programs compiled to NFP bytecode



```
if (pkt.size < 100)
    return DROP;
if (pkt[ethertype] == IPV4)
    ip4_counter++;
if (pkt[ethertype] == IPV6)
    for (...) ...;
return PASS;
```

example program

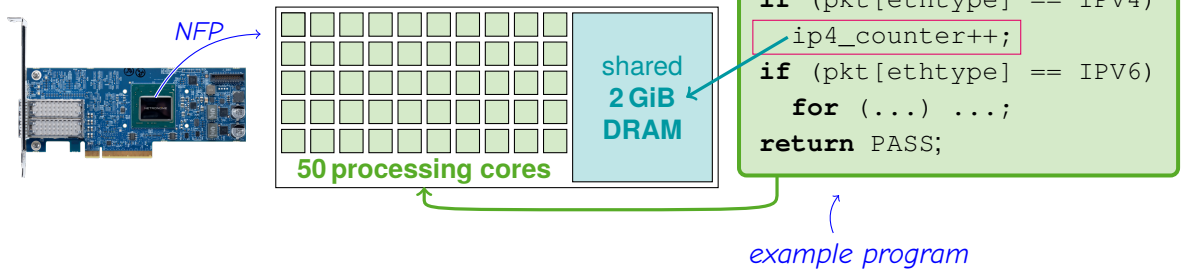
- **Challenges**

- ▶ The throughput depends on the executed instructions

Throughput Analysis

- **Netronome Agilio CX 2x40 GbE**

- ▶ Netronome Flow Processor (NFP)
- ▶ BPF/XDP programs compiled to NFP bytecode



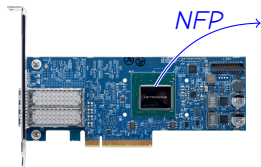
- **Challenges**

- ▶ The throughput depends on the executed instructions
- ▶ DRAM access can be a bottleneck

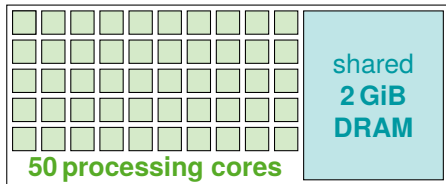
Throughput Analysis

- **Netronome Agilio CX 2x40 GbE**

- ▶ Netronome Flow Processor (NFP)
- ▶ BPF/XDP programs compiled to NFP bytecode



NFP →



```
if (pkt.size < 100)
    return DROP;
if (pkt[ethertype] == IPV4)
    ip4_counter++;
if (pkt[ethertype] == IPV6)
    for (...) ...;
return PASS;
```

example program

- **Challenges**

- ▶ The throughput depends on the executed instructions
- ▶ DRAM access can be a bottleneck
- ▶ Packet sizes influence throughput

Throughput depends on executed instructions

- **Different paths cause different throughput**
 - ⇒ The slowest path establishes a lower bound

```
if (pkt.size < 100)
  return DROP;
if (pkt[ethertype] == IPV4)
  ip4_counter++;
if (pkt[ethertype] == IPV6)
  for (...) ...;
return PASS;
```

Throughput depends on executed instructions

- **Different paths cause different throughput**
 - ⇒ The slowest path establishes a lower bound
- **Not every path can be triggered by a packet**
 - ▶ Perhaps a tighter lower bound is possible
 - ⇒ Check each path for contradictions

```
if (pkt.size < 100)
    return DROP;
if (pkt[ethertype] == IPV4)
    ip4_counter++;
if (pkt[ethertype] == IPV6)
    for (...) ...;
return PASS;
```

Throughput depends on executed instructions

- **Different paths cause different throughput**
 - ⇒ The slowest path establishes a lower bound
- **Not every path can be triggered by a packet**
 - ▶ Perhaps a tighter lower bound is possible
 - ⇒ Check each path for contradictions
- **Too many paths (path explosion)**
 - ▶ Up to 2^n paths for n `ifs`
 - ⇒ Only analyze the slowest paths

```
if (pkt.size < 100)
    return DROP;
if (pkt[ethertype] == IPV4)
    ip4_counter++;
if (pkt[ethertype] == IPV6)
    for (...) ...;
return PASS;
```

Throughput depends on executed instructions

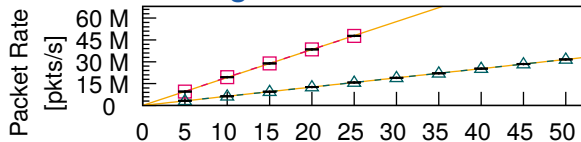
- **Different paths cause different throughput**
 - ⇒ The slowest path establishes a lower bound
- **Not every path can be triggered by a packet**
 - ▶ Perhaps a tighter lower bound is possible
 - ⇒ Check each path for contradictions
- **Too many paths (path explosion)**
 - ▶ Up to 2^n paths for n `ifs`
 - ⇒ Only analyze the slowest paths

```
if (pkt.size < 100)
    return DROP;
if (pkt[ethertype] == IPV4)
    ip4_counter++;
if (pkt[ethertype] == IPV6)
    for (...) ...;
return PASS;
```

How do individual instructions influence the throughput?

Processing Cost

- **Linear scaling over 50 cores**

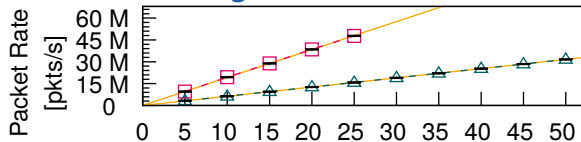


- **Most NFP instructions take 1 cycle**

- ▶ Branches take 2-3 cycles when taken
- ▶ Deterministic DRAM access times when not overloaded

Processing Cost

- **Linear scaling over 50 cores**



- **Most NFP instructions take 1 cycle**
 - ▶ Branches take 2-3 cycles when taken
 - ▶ Deterministic DRAM access times when not overloaded

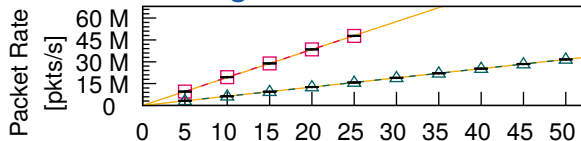
DRAM Cost

- **Too many memory instructions overload the DRAM**
- **DRAM throughput varies up to $\times 4$**
 - ⇒ Lower bound

| | | |
|-------------|--------------------------|--------------------------|
| Atomic Inc: | $\geq 248\text{M ops/s}$ | $cost = \frac{50}{rate}$ |
| Read32: | $\geq 197\text{M ops/s}$ | |

Processing Cost

- **Linear scaling over 50 cores**



- **Most NFP instructions take 1 cycle**
 - ▶ Branches take 2-3 cycles when taken
 - ▶ Deterministic DRAM access times when not overloaded

A program path is either **processing or **DRAM** bottlenecked**

- Per instruction cost tuple: (**processing cost**; **DRAM cost**)
- Per path bottleneck: maximum over **processing** & **DRAM** cost

DRAM Cost

- **Too many memory instructions overload the DRAM**
- **DRAM throughput varies up to $\times 4$**
 - ⇒ Lower bound

Atomic Inc: $\geq 248\text{M ops/s}$
Read32: $\geq 197\text{M ops/s}$ $cost = \frac{50}{rate}$

- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)

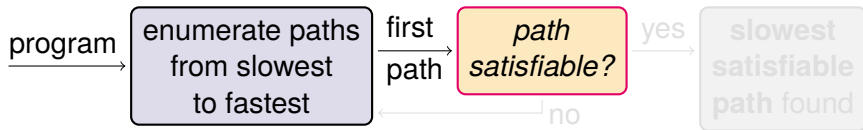
- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



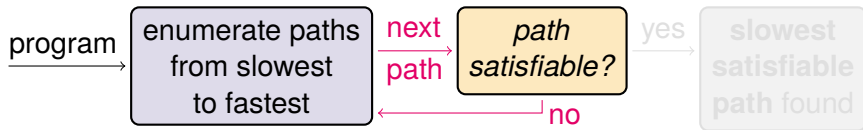
- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



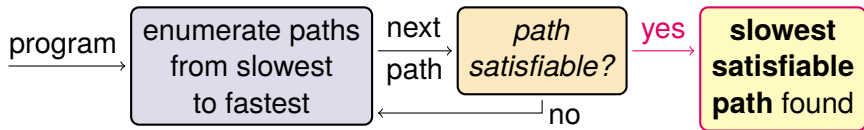
- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



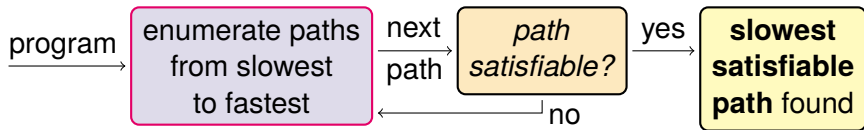
- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



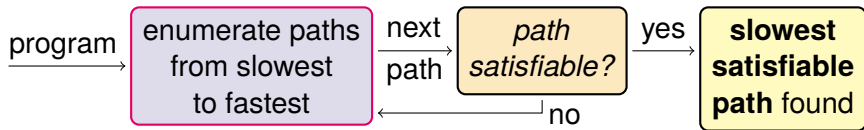
- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



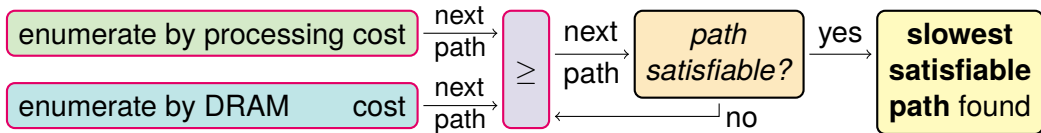
- **Two separate cost functions (processing cost; DRAM cost)**

Analysing only the slowest paths

- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)

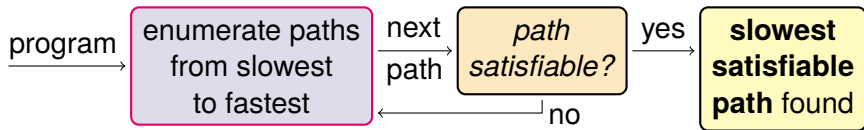


- **Two separate cost functions (processing cost; DRAM cost)**

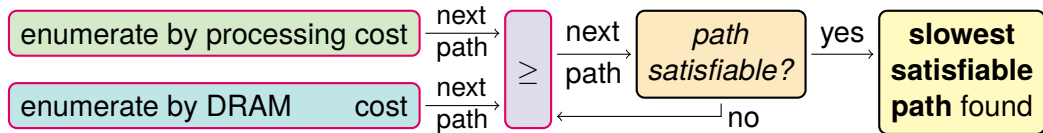


Analysing only the slowest paths

- **Can not enumerate all paths:** Incremental longest path algorithm [Kundu, 1994]
- **Check paths for contradictions:** Use SMT solver (similar to symbolic execution)



- **Two separate cost functions (processing cost; DRAM cost)**



Determines **packet-rate** guarantees

- **Bit-Rate throughput influenced by**
 - ▶ The time it takes to process a packet
 - ▶ The size of the packet

Programs which process fixed sized headers

- **Processing mostly independent from packet size**
 - Packet-rate guarantees

Programs which loop over multiple headers and payload

- **Processing depends on packet sizes**
 - Bit-rate guarantees

- **Static program analysis to identify packet size classes**

Example Program

```
if (pkt.size < 100)
    return DROP;
if (...) ...;
if (...) ...;
return PASS;
```

- **Static program analysis to identify packet size classes**

Example Program

```
if (pkt.size < 100)
  return DROP;
if (...) ...;
if (...) ...;
return PASS;
```

- **Static program analysis to identify packet size classes**

| Example Program | pkt.size: 60-99 | pkt.size \geq 100 |
|---|-----------------|---------------------|
| <pre>if (pkt.size < 100) return DROP; if (...) ...; if (...) ...; return PASS;</pre> | | |

- **Static program analysis to identify packet size classes**

| Example Program | pkt.size: 60-99 | pkt.size \geq 100 |
|---|--|---------------------|
| <pre>if (pkt.size < 100) return DROP; if (...) ...; if (...) ...; return PASS;</pre> | <pre>if (pkt.size < 100) return DROP;</pre> | |

- Static program analysis to identify packet size classes

| Example Program | pkt.size: 60-99 | pkt.size \geq 100 |
|---|--|--|
| <pre>if (pkt.size < 100) return DROP; if (...) ...; if (...) ...; return PASS;</pre> | <pre>if (pkt.size < 100) return DROP;</pre> | <pre>if (pkt.size < 100) if (...) ...; if (...) ...; return PASS;</pre> |

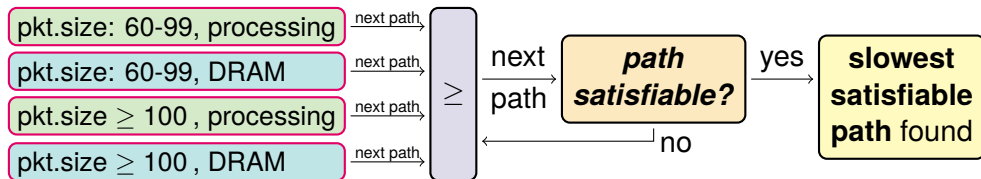
- Static program analysis to identify packet size classes

| Example Program | pkt.size: 60-99 | pkt.size \geq 100 |
|---|--|--|
| <pre>if (pkt.size < 100) return DROP; if (...) ...; if (...) ...; return PASS;</pre> | <pre>if (pkt.size < 100) return DROP;</pre> | <pre>if (pkt.size < 100) if (...) ...; if (...) ...; return PASS;</pre> |

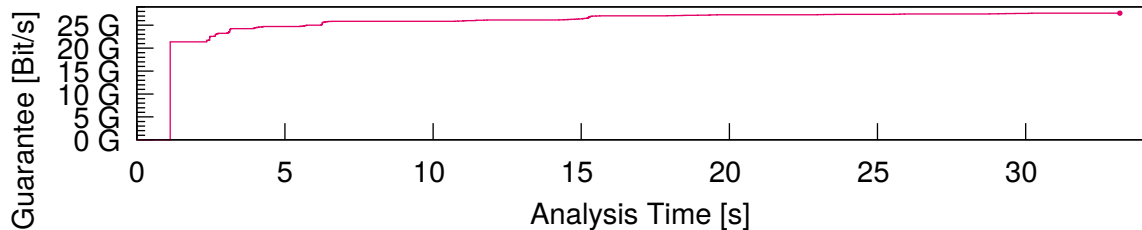
- Static program analysis to identify packet size classes

| Example Program | pkt.size: 60-99 | pkt.size \geq 100 |
|---|--|--|
| <pre>if (pkt.size < 100) return DROP; if (...) ...; if (...) ...; return PASS;</pre> | <pre>if (pkt.size < 100) return DROP;</pre> | <pre>if (pkt.size < 100) if (...) ...; if (...) ...; return PASS;</pre> |

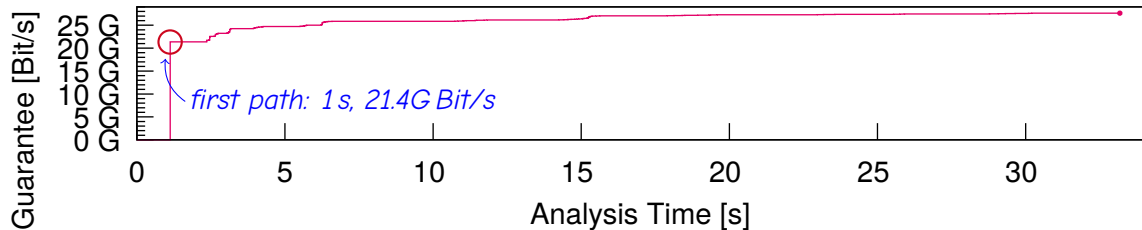
- Enumerate paths ordered by bit-rate



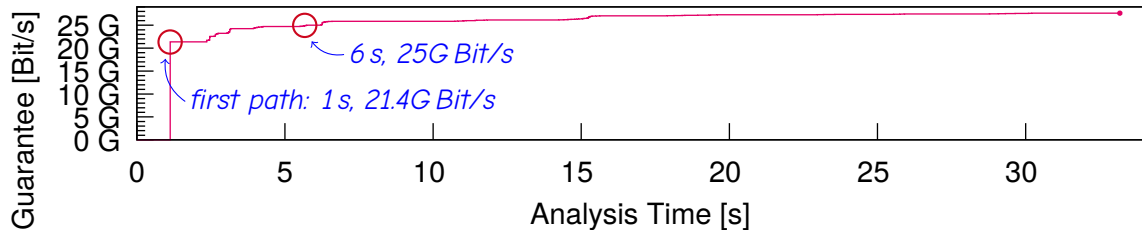
Example Program: QUIC Load Balancer with IPv6 Option parsing



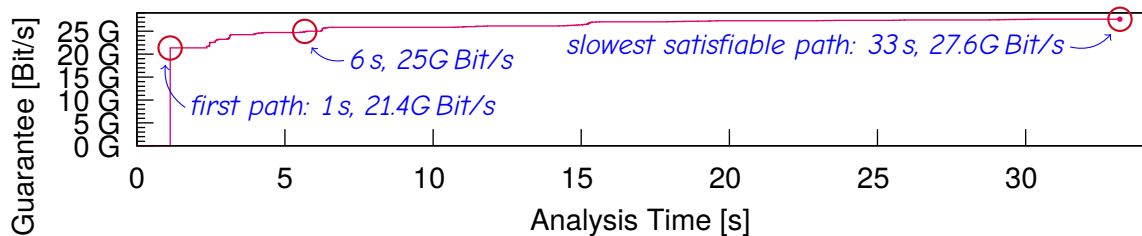
Example Program: QUIC Load Balancer with IPv6 Option parsing



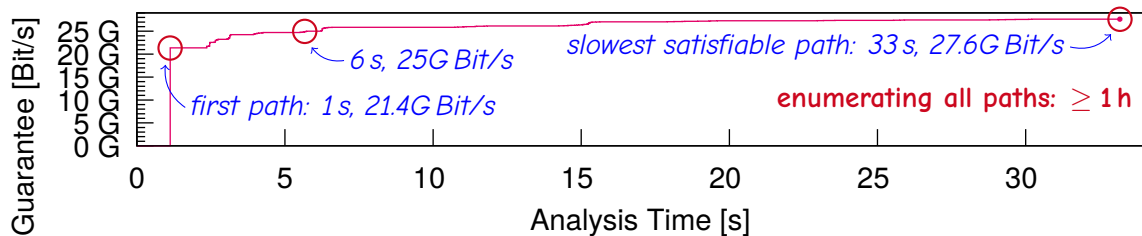
Example Program: QUIC Load Balancer with IPv6 Option parsing



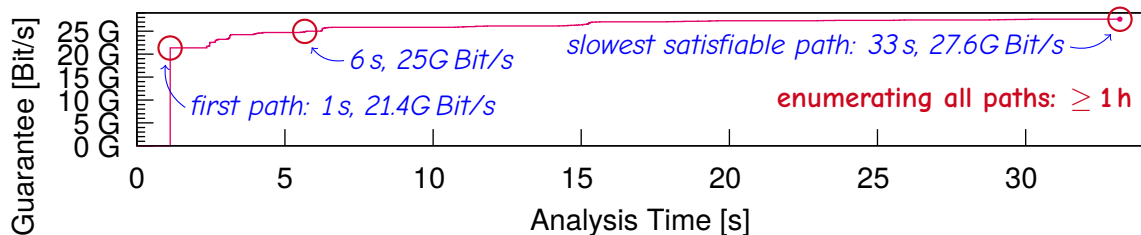
Example Program: QUIC Load Balancer with IPv6 Option parsing



Example Program: QUIC Load Balancer with IPv6 Option parsing



Example Program: QUIC Load Balancer with IPv6 Option parsing



- Analyzed 9 real XDP/BPF programs on a desktop PC (Core i7, 16 GiB)
 - ▶ Up to 102 s analysis time (except maliciously crafted example)
 - ▶ SMT checking improved throughput guarantees by up to 44%

Is **fast** enough to be part of the regular development cycle

- Measured throughput of **21k paths** through all analyzed programs
 - ▶ Minimally sized example packets produced by the SMT solver

- Measured throughput of **21k paths** through all analyzed programs
 - ▶ Minimally sized example packets produced by the SMT solver

| <i>Program</i> | <i>Slowest Path</i> |
|-------------------------|-----------------------|
| switch.p4 (parser) | ✓ |
| Cloudflare DoS | ✓ |
| QUIC LB (IPv4) | ✓ |
| QUIC LB (IPv6) | ✓ |
| RTP a→ μ -law | ✓ |
| RTP a→ μ -law (opt) | ✓ |
| DNS Cache | 7th |
| Count-Min (5) | ✓ |
| Count-Min (20) | ✓ |

Evaluation: Estimation Accuracy

- Measured throughput of **21k paths** through all analyzed programs
 - Minimally sized example packets produced by the SMT solver

| <i>Program</i> | <i>Slowest Path</i> | <i>Estimated</i> | <i>Measured</i> | <i>Accuracy</i> | |
|--------------------|-----------------------|------------------|-----------------|-----------------|------------------------------|
| switch.p4 (parser) | ✓ | 24.7G Bit/s | 25.8G Bit/s | +4.1% | <i>not a lower bound</i> |
| Cloudflare DoS | ✓ | 35.2G Bit/s | 34.7G Bit/s | -1.0% | ← |
| QUIC LB (IPv4) | ✓ | 22.8G Bit/s | 23.5G Bit/s | +2.9% | |
| QUIC LB (IPv6) | ✓ | 27.6G Bit/s | 28.5G Bit/s | +2.9% | |
| RTP a→μ-law | ✓ | 2.97G Bit/s | 2.97G Bit/s | ✓ | <i>worst underestimation</i> |
| RTP a→μ-law (opt) | ✓ | 4.70G Bit/s | 4.70G Bit/s | ✓ | |
| DNS Cache | 7th | 9.0G Bit/s | 10.4G Bit/s | +13.1% | ← |
| Count-Min (5) | ✓ | 21.6G Bit/s | 22.2G Bit/s | +2.4% | |
| Count-Min (20) | ✓ | 6.0G Bit/s | 6.0G Bit/s | ✓ | |

Evaluation: Estimation Accuracy

- Measured throughput of **21k paths** through all analyzed programs
 - Minimally sized example packets produced by the SMT solver

| Program | Slowest Path | Estimated | Measured | Accuracy | |
|--------------------|-----------------------|-------------|-------------|---------------|------------------------------|
| switch.p4 (parser) | ✓ | 24.7G Bit/s | 25.8G Bit/s | +4.1% | <i>not a lower bound</i> |
| Cloudflare DoS | ✓ | 35.2G Bit/s | 34.7G Bit/s | -1.0% | ← |
| QUIC LB (IPv4) | ✓ | 22.8G Bit/s | 23.5G Bit/s | +2.9% | |
| QUIC LB (IPv6) | ✓ | 27.6G Bit/s | 28.5G Bit/s | +2.9% | |
| RTP a→μ-law | ✓ | 2.97G Bit/s | 2.97G Bit/s | ✓ | <i>worst underestimation</i> |
| RTP a→μ-law (opt) | ✓ | 4.70G Bit/s | 4.70G Bit/s | ✓ | |
| DNS Cache | 7th | 9.0G Bit/s | 10.4G Bit/s | +13.1% | ← |
| Count-Min (5) | ✓ | 21.6G Bit/s | 22.2G Bit/s | +2.4% | |
| Count-Min (20) | ✓ | 6.0G Bit/s | 6.0G Bit/s | ✓ | |

Accurate throughput guarantees

Throughput Guarantees for Processor-based SmartNICs

- Which **packet-rate** or **bit-rate** is achievable by a given program?
- Similar determinism as match-action pipelines and FPGAs

- **Our approach is fast**

- ▶ Because, we combine incremental ordered enumeration with SMT checks

- **Our approach is accurate**

- ▶ Because, we model the SmartNIC performance characteristics



all the details 

Full Implementation & All Measurement Data

<https://github.com/johannes-krude/nfp-pred-artifacts>

