# Decoupling Optimizations and Algorithms in Network Functions

Omid Alipourfard and Minlan Yu
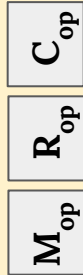
4.4 Mpps

$M_{op}$  $R_{op}$  $C_{op}$

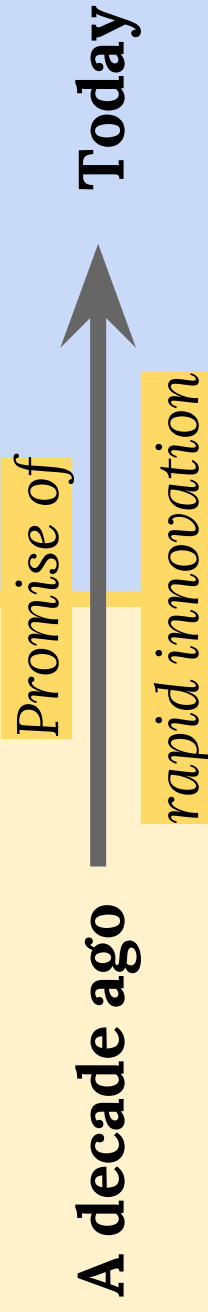Semantically equivalent NF pipeline

$\cong$

13.2 Mpps

$P_{op}$ → $P_{op}$  $M_{op}$  $R_{op}$  $C_{op}$

$M_{op}$ → $M_{op}$  $R_{op}$  $C_{op}$ → $M_{op}$

2

# Network functions are popular

**A decade ago** → **Today**

*Promise of rapid innovation*

## Hardware boxes

- Faster
- Slow to update
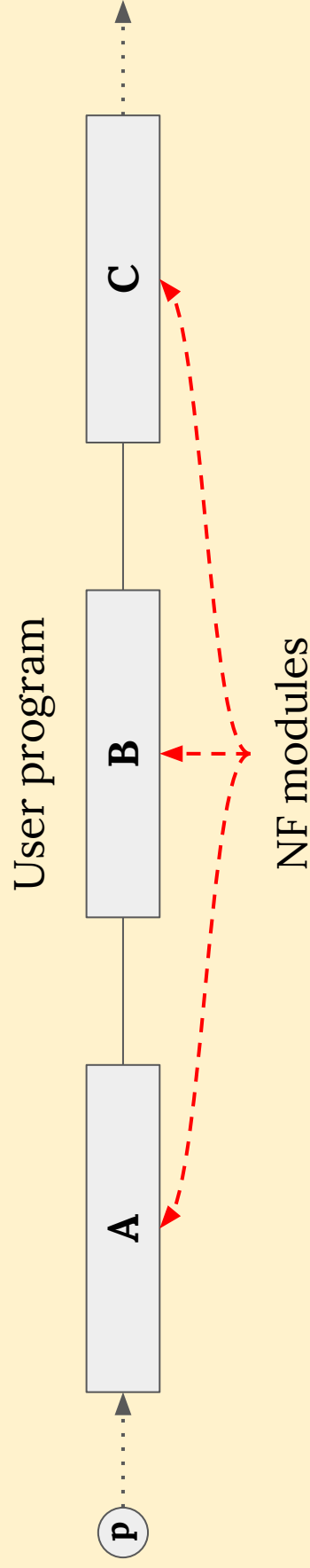- More expensive

## Software boxes

- Slower
- Fast update cycles
- Cheaper

3

# Performance is critical

Innovation and cost benefits only come with **good performance**:

- Lower tail and average **latency**
  - RPC like applications
  - Every μs counts
- Higher **throughput**
  - WAN, ISP, and storage like applications
  - Every additional bytes/cycle counts

# Many optimizations for packet processing

User program

NF modules

# Many optimizations for packet processing

**Batching packets**



Better throughput (amortize static cost)
Higher latency (wait for the batch to finish)

6

# Many optimizations for packet processing



**Prefetching**

P → A → B → C →

*Prefetch data*

P → A → B → C →

**Better throughput (depending on cache availability)**
**Less latency variability**

# Many optimizations for packet processing

Data structure tuning, e.g., layout, size, algorithm

P → A → B → C ⤑

P → A → B → C' ⤑

**Better throughput**
**Lower latency variability (Cache locality)**

# Many optimizations for packet processing

User program

Many other optimizations, e.g., fastpath, reorganizing the pipeline, end-to-end optimizations.

P → A — B — C ⇢

Lower latency variability (Cache locality)
Better throughput

*Whereas optimizations are well known, applying optimizations requires many **trials and errors.***

# Applying optimizations takes a huge effort, cont.

**User program**

| Checksum | — | Routing | — | Measurement |
|----------|---|---------|---|-------------|

(p)

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

10 mil unique flows
Top 1% has 10% volume

# Applying optimizations takes a huge effort, cont.

User program

| Checksum | — | Routing | — | Measurement |

**p**

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

10 mil unique flows
Top 1% has 10% volume

With DPDK & compiler optimizations on:

**4.4 Mpps**

# Applying optimizations takes a huge effort, cont.



**User program**

**Checksum** — **Routing** — **Measurement**

**Platform:**
Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**
10 mil unique flows
Top 1% has 10% volume

*Batching* $_{8\ pkts}$

*Batching* $_{128\ pkts}$

*Reordering*

*Prefetch inst.*

$\mathbf{P}_{op}$  $\mathbf{P}_{op}$  $\mathbf{M}_{op}$  $\mathbf{R}_{op}$  $\mathbf{C}_{op}$  $\mathbf{M}_{op}$  $\mathbf{M}_{op}$  $\mathbf{R}_{op}$  $\mathbf{C}_{op}$  $\mathbf{M}_{op}$

# Applying optimizations takes a huge effort, cont.



**User program**

Checksum — Routing — Measurement

**13.2 Mpps**

P$_{op}$ → P$_{op}$ M$_{op}$ R$_{op}$ C$_{op}$ M$_{op}$ → M$_{op}$ R$_{op}$ C$_{op}$ → M$_{op}$

*Batching* $_{8\ pkts}$

*Batching* $_{128\ pkts}$

*Prefetch inst.*

**Platform:**
Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**
10 mil unique flows
Top 1% has 10% volume

14

# Optimizations depend on the **workload**

User program

| Checksum | — | Routing | — | Measurement |

**p**

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

**Consider a different workload:**

- 100 k unique flows
- Top 1% has 99% traffic volume

# Optimizations depend on the **workload**

User program

| Checksum | Routing | Measurement |

**p**

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

**100 k unique flows**
**Top 1% has 99% volume**
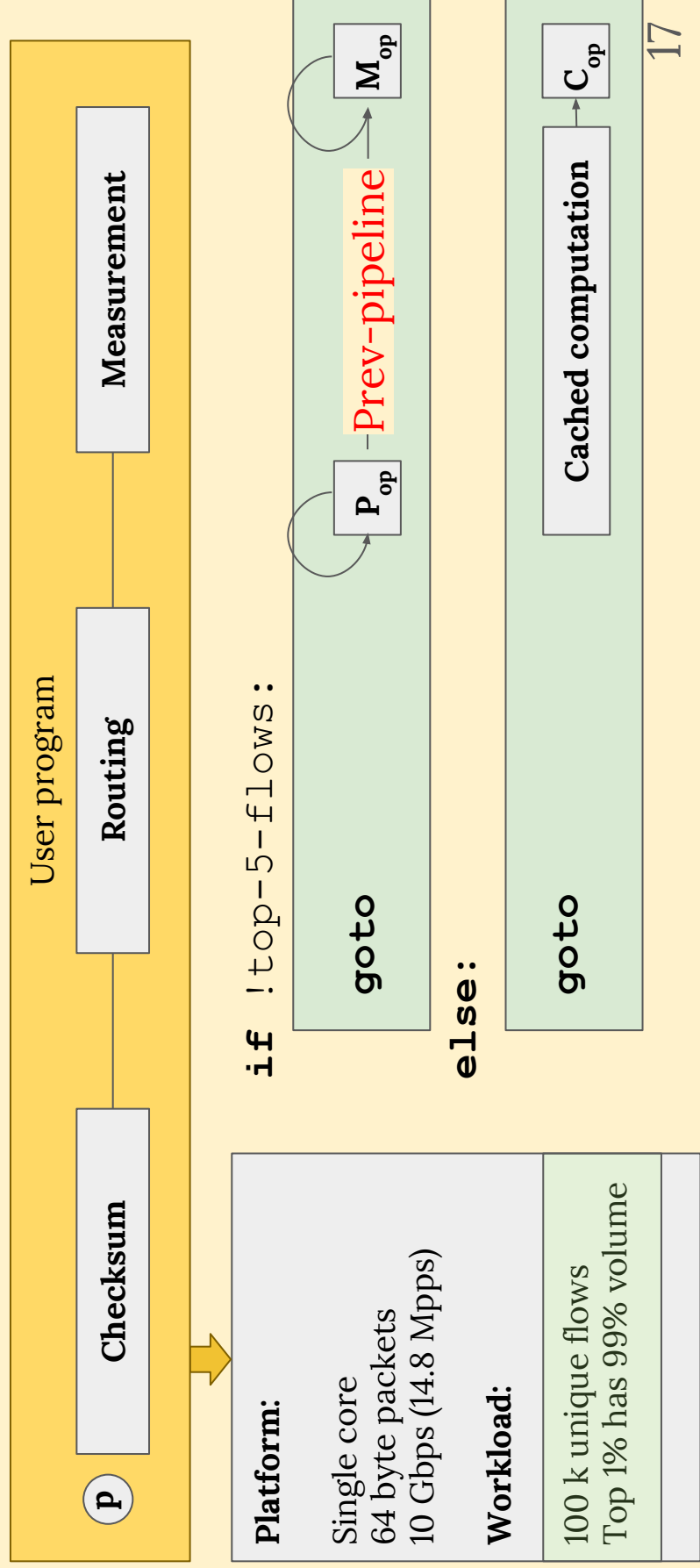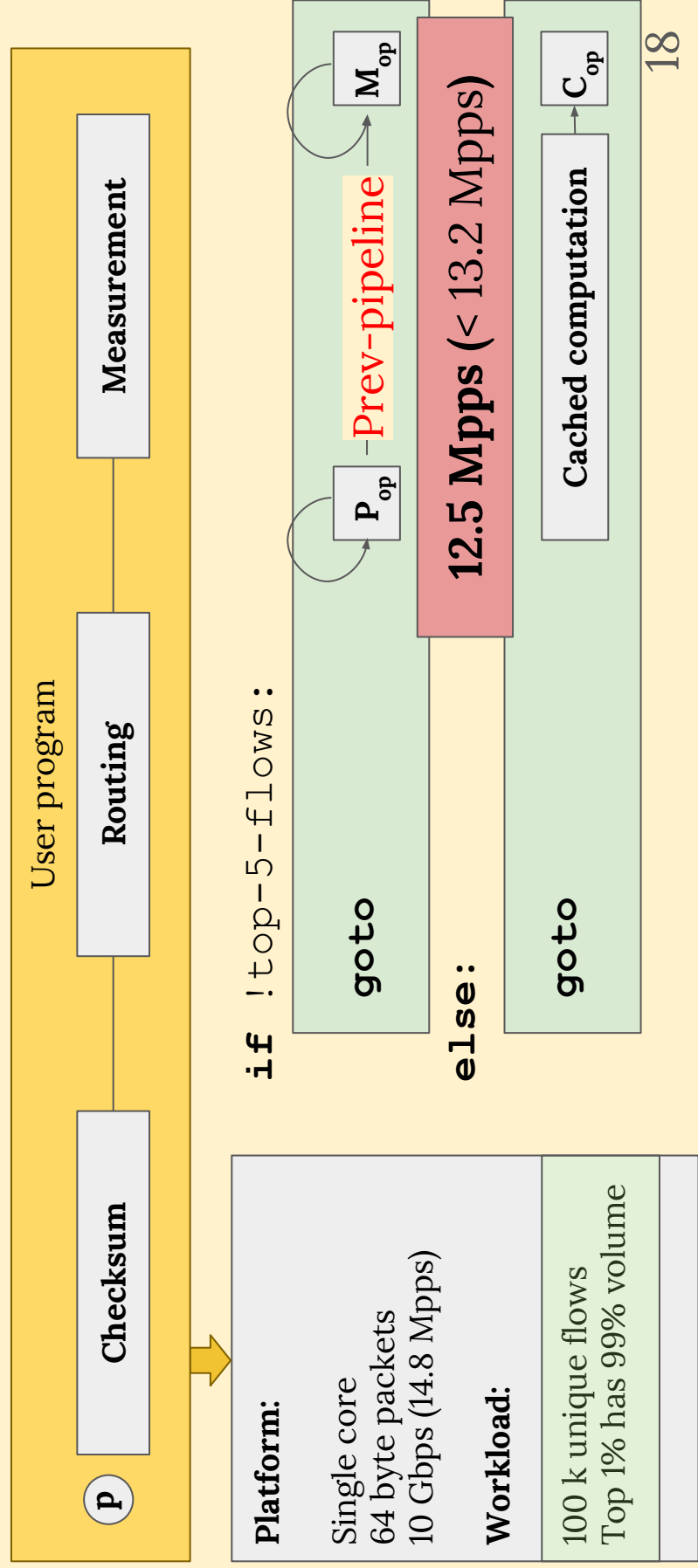
**Observation:**
Top–5 flows have 70% of traffic.

**Idea:**
Build a fastpath (cache the computation) of top–5 flows.

16

# Optimizations depend on the **workload**



User program

| p — Checksum — Routing — Measurement |

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

100 k unique flows
Top 1% has 99% volume

```
if !top-5-flows:
    goto
        P_op — Prev-pipeline — M_op
else:
    goto
        Cached computation — C_op
```

17

# Opt. depend on the **workload** and are (*very*) **HARD!**

## User program

**Checksum** — **Routing** — **Measurement**

**p**

```
if !top-5-flows:

    goto        $P_{op}$ — Prev-pipeline — $M_{op}$

else:          12.5 Mpps (<13.2 Mpps)

    goto        Cached computation → $C_{op}$
```

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

100 k unique flows
Top 1% has 99% volume

18

# Optimizations depend on the **workload**



User program

| p | Checksum | Routing | Measurement |

Platform:

Single core
64 byte packets
10 Gbps (14.8 Mpps)

Workload:

100 k unique flows
Top 1% has 99% volume

$P_{op}$ → $P_{op}$ $F_{op}$ $S_{op}$ $C_{op}$

*Slowpath extraction*

*Fastpath extraction*

*Cached–computation*

$CC_{op}$

$S_{op}$ $M_{op}$ $R_{op}$ $M_{op}$

19

# Optimizations depend on the **workload**



User program

| Checksum | Routing | Measurement |

**14.8 Mpps**

*Fastpath extraction*

*Cached-computation*

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

100 k unique flows
Top 1% has 99% volume

20

# Very different pipeline!

# Optimizations depend on the **platform**

User program

| Checksum | Routing | Measurement |

**Platform:**

Single core & NIC (checksum offloading)
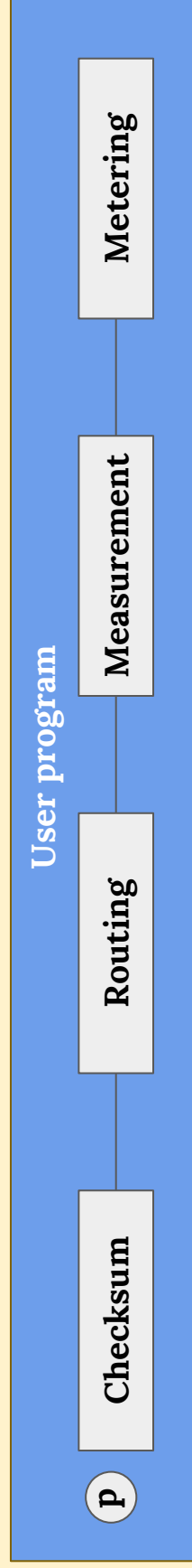64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

100 k unique flows
Top 1% has 99% volume

- Offload checksum to the NIC.

- Reoptimize measurement + routing

# Optimizations depend on the **NF-chain**

**User program**

p — Checksum — Routing — Measurement — Metering

**Platform:**

Single core
64 byte packets
10 Gbps (14.8 Mpps)

**Workload:**

100 k unique flows
Top 1% has 99% volume

Redo all the optimizations for the NF chain.

**Impossible to optimize for every NF-chain, platform, and workload.**

*What are the **fundamental challenges** to automatically optimize NFs for commodity servers?*

# x86 is not designed for packet processing.

- No packet pipelining (only instructions!)
  - Mapping requires knowledge of limited resources
- Non-determinism:
  - Variable memory access latency
  - Shared resources with other application

Optimizing **compiler** goals are different.

- Compiler goal: minimize **completion time** or **code size**

  – NF goals: minimize **latency** or maximize **throughput**

- **Packet optimizations** could change semantics

  – Reorder packets (keep each TCP conn. still in-order)

Optimizations **impact** each other.

- **Trial and error:**
  - Large batches help / Prefetching help
  - Large batches with prefetching pollutes the cache.
- **Proactively** optimizing the code is impossible
  - Workload/Platform/NF Chain

**Solution:** *Decouple **algorithms** and **optimizations** in network function design*

*Domain Specific Language*

**Solution:** *Decouple **algorithms** and **optimizations** in network function design*

*Optimizing Runtime*

# Domain Specific Language

- Express algorithm on a single packet

## Domain Specific Language

- Express algorithm on a single packet

  - Make packets **first class type**

# Domain Specific Language

- Express algorithm on a single packet
- Include the abstractions available in today's hardware

# Domain Specific Language

- Express algorithm on a single packet

- Include the abstractions available in today's hardware

  - Packet processing **keywords**

| LPMTable | Hash | TCPChecksum |
|----------|------|-------------|

# Domain Specific Language

- Express algorithm on a single packet

- Include the abstractions available in today's hardware

- Include "hints" to guide optimization choices

# Domain Specific Language

- Express algorithm on a single packet

- Include the abstractions available in today's hardware

- Include "hints" to guide optimization choices

  - Optimization **keywords**: pure, commutative, ...

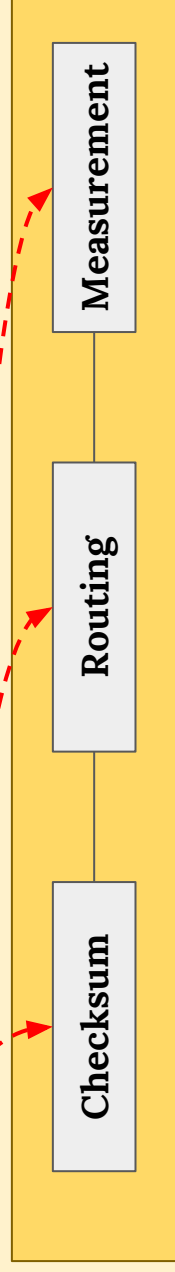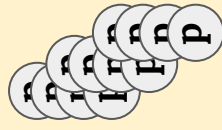| Extern | Pure | Commutative |
|--------|------|-------------|

# From **language** to **machine code**

Can we systematically make efficient code?

- Workload variations
- Available Platforms
- NF chain

**No single best optimization strategy!**

# Profile guided optimization

- Profile code while executing



- Profile traffic characteristics

# Template based optimization

- Abstract syntax tree transformation
- Templates with holes
    - Use well-known opt. templates: *batch, prefetch, ...*
    - Preserve the packet processing semantic

## Summary

- End-to-end NF optimization has meaningful gains

- Figuring out the right set of optimizations is difficult

  - But NF optimizations are well-known

- By decoupling algorithms and optimizations, we can automatically optimize NF functions.

**❝** In a given paradigm, ... programs become complicated for technical reasons that have no direct relationship to ... problem ... being solved. This is a sign that there is a new concept waiting to be discovered. **❞**

→*Peter Van Roy*

## *Thanks!*