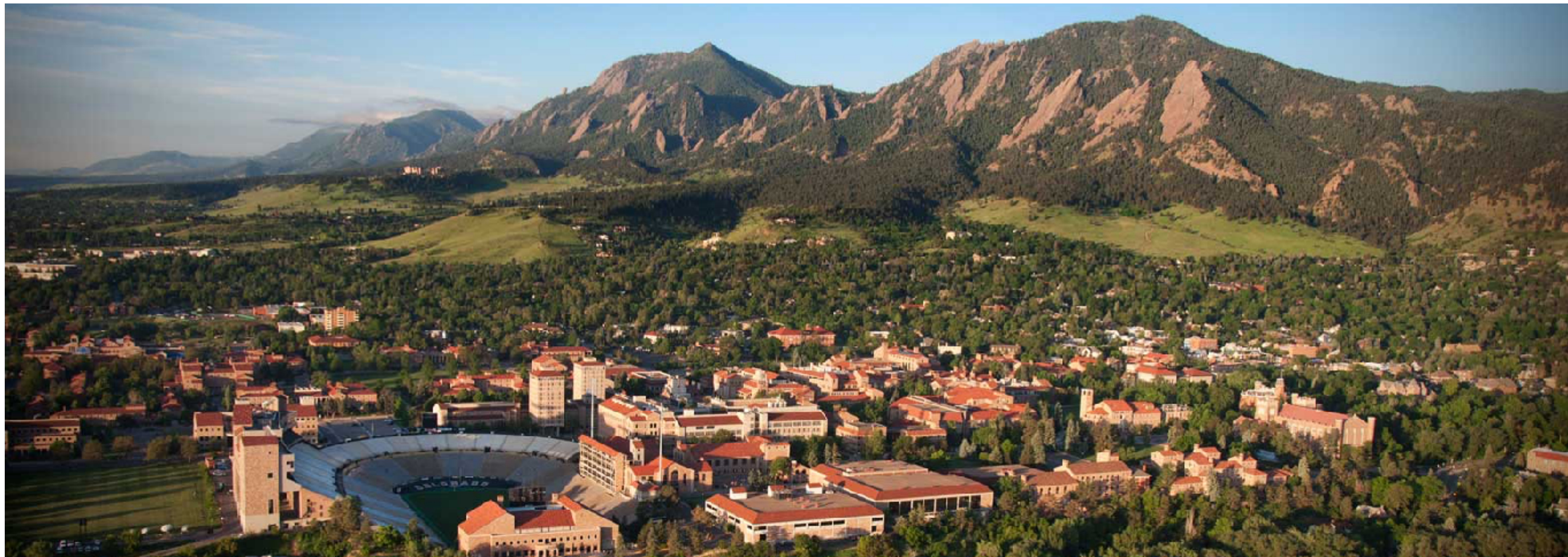# Packet-Level Analytics in Software without Compromises

HotCloud '18, July 9th, 2018, Boston, MA
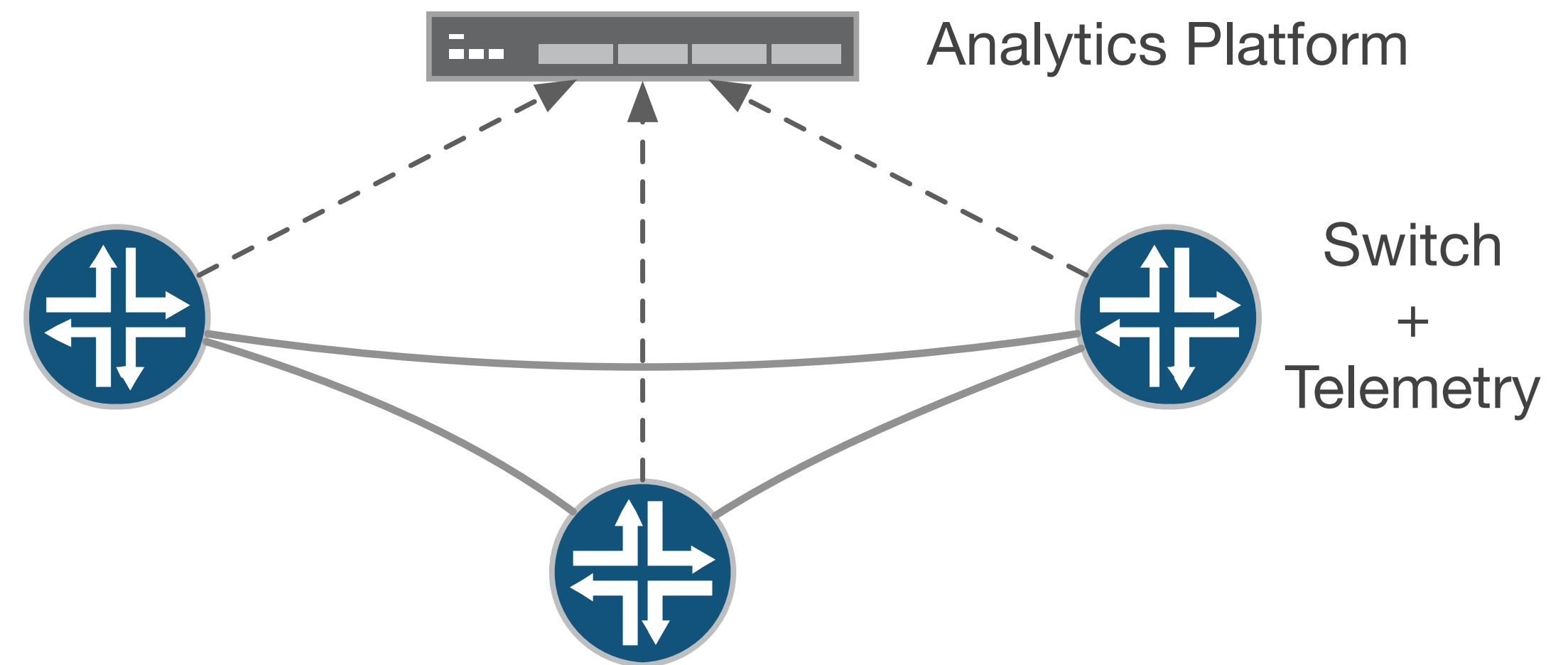
Oliver Michel
John Sonchack
Eric Keller
Jonathan M. Smith
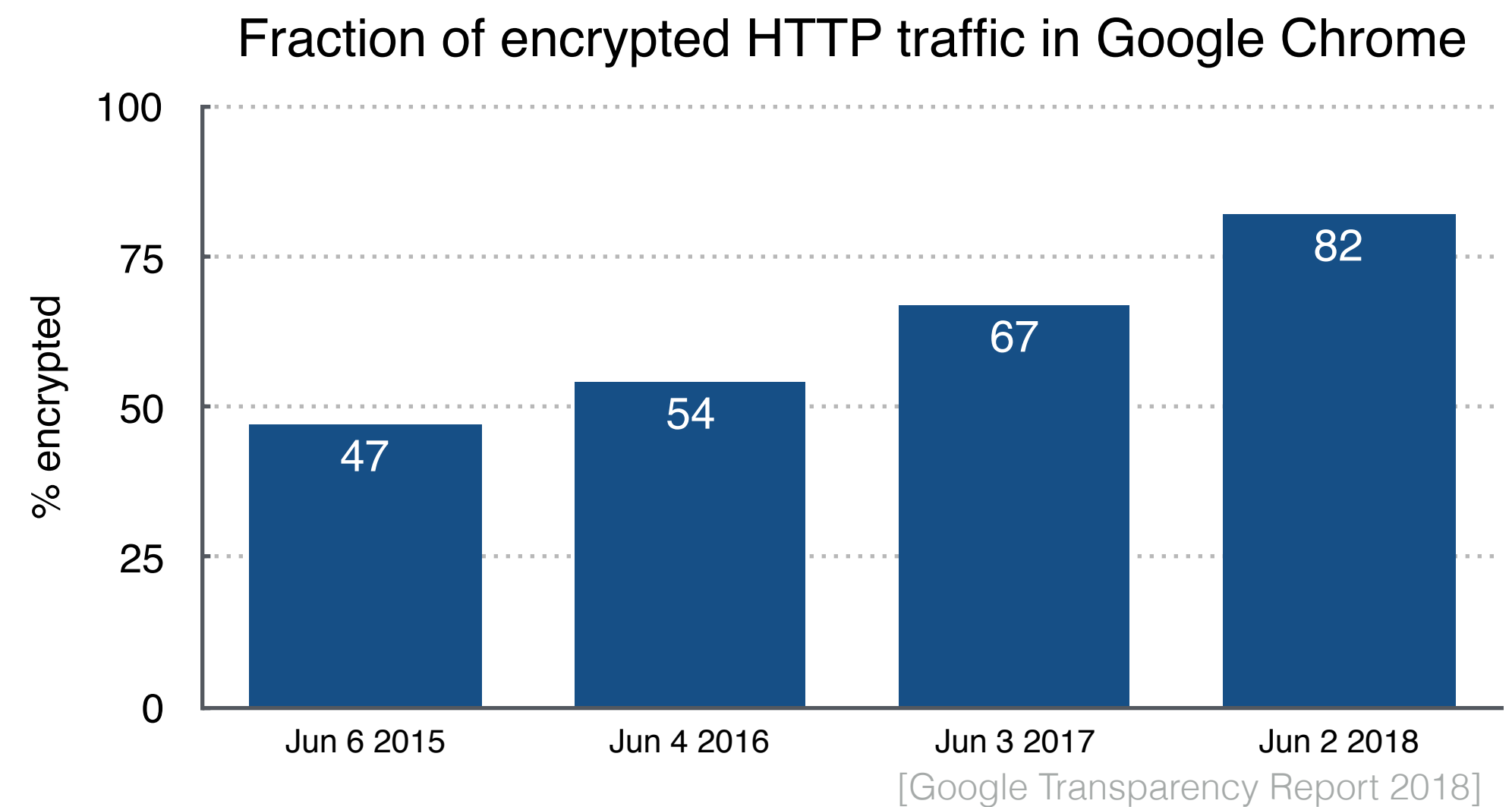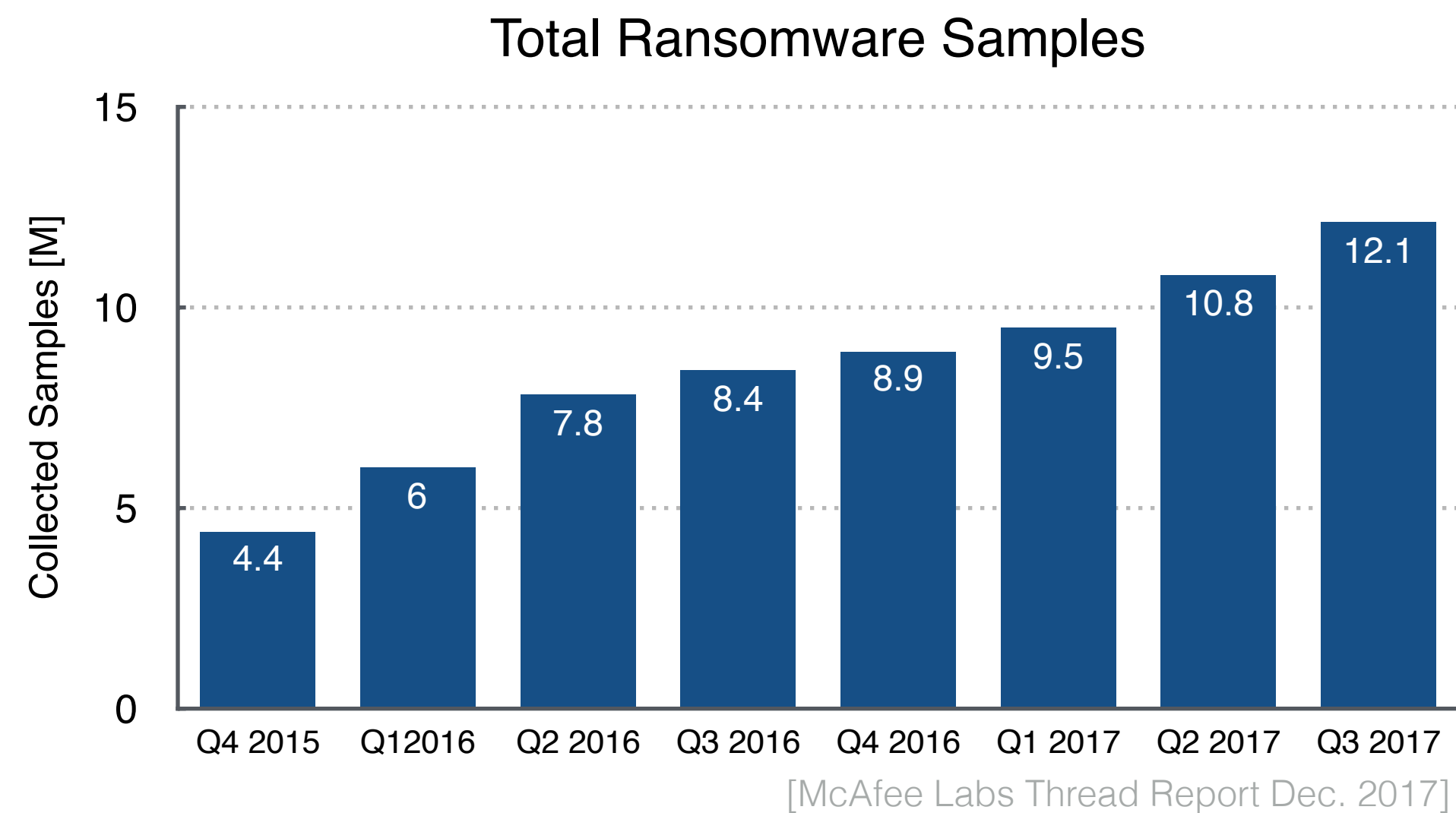
University of Colorado Boulder

Penn

# Network monitoring is important

- Security issues

- Performance issues

- Equipment failure

- Misconfiguration

Analytics Platform

Switch
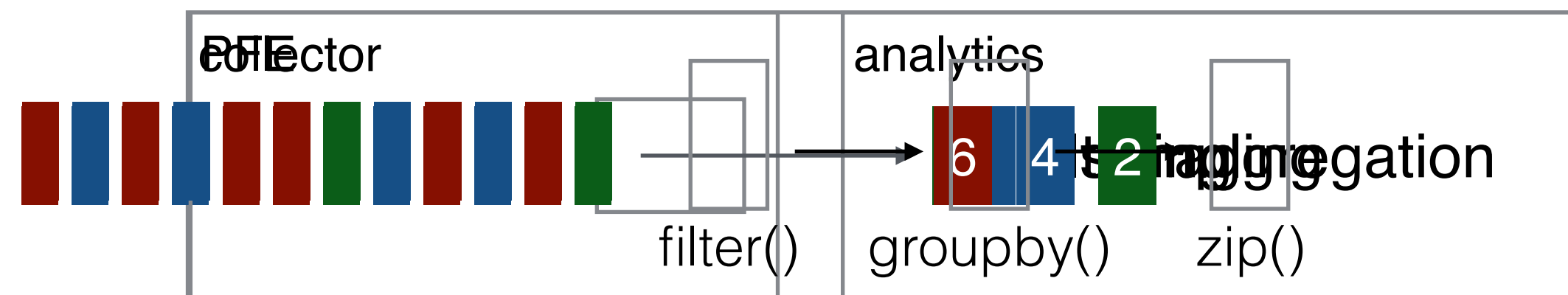+
Telemetry

# Challenging environment

- more traffic

- more threats

- encrypted traffic

**Total Ransomware Samples**



[McAfee Labs Thread Report Dec. 2017]

**Fraction of encrypted HTTP traffic in Google Chrome**



[Google Transparency Report 2018]

# Existing systems make compromises
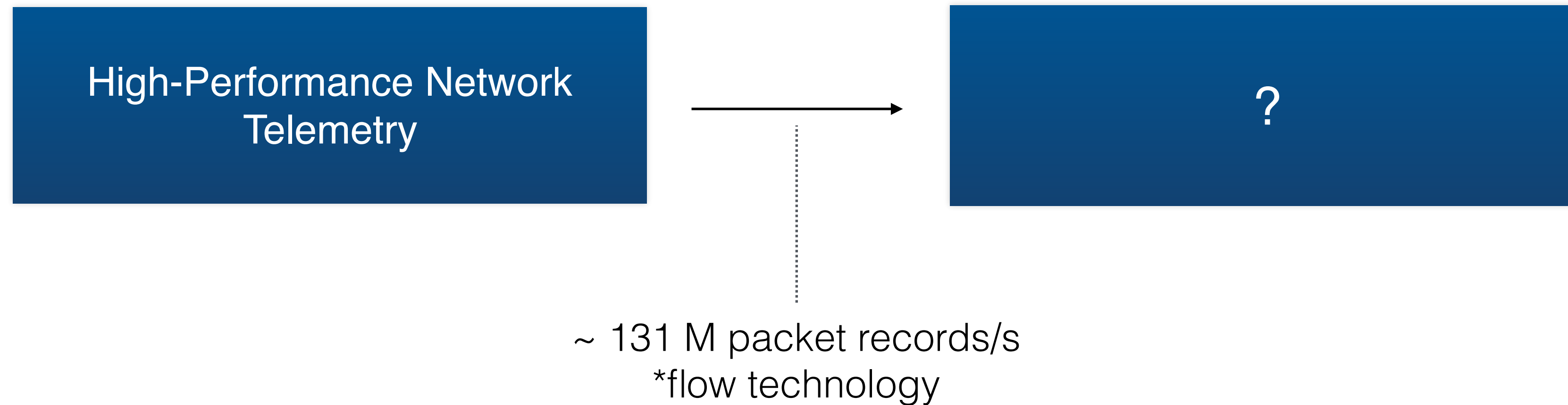


loss of information

loss of capability

# Programmable Forwarding Engines

- Programmable Forwarding Engines

  - Marple [SIGCOMM 2017]

  - *flow [ATC 2018]



High-Performance Network Telemetry → ?

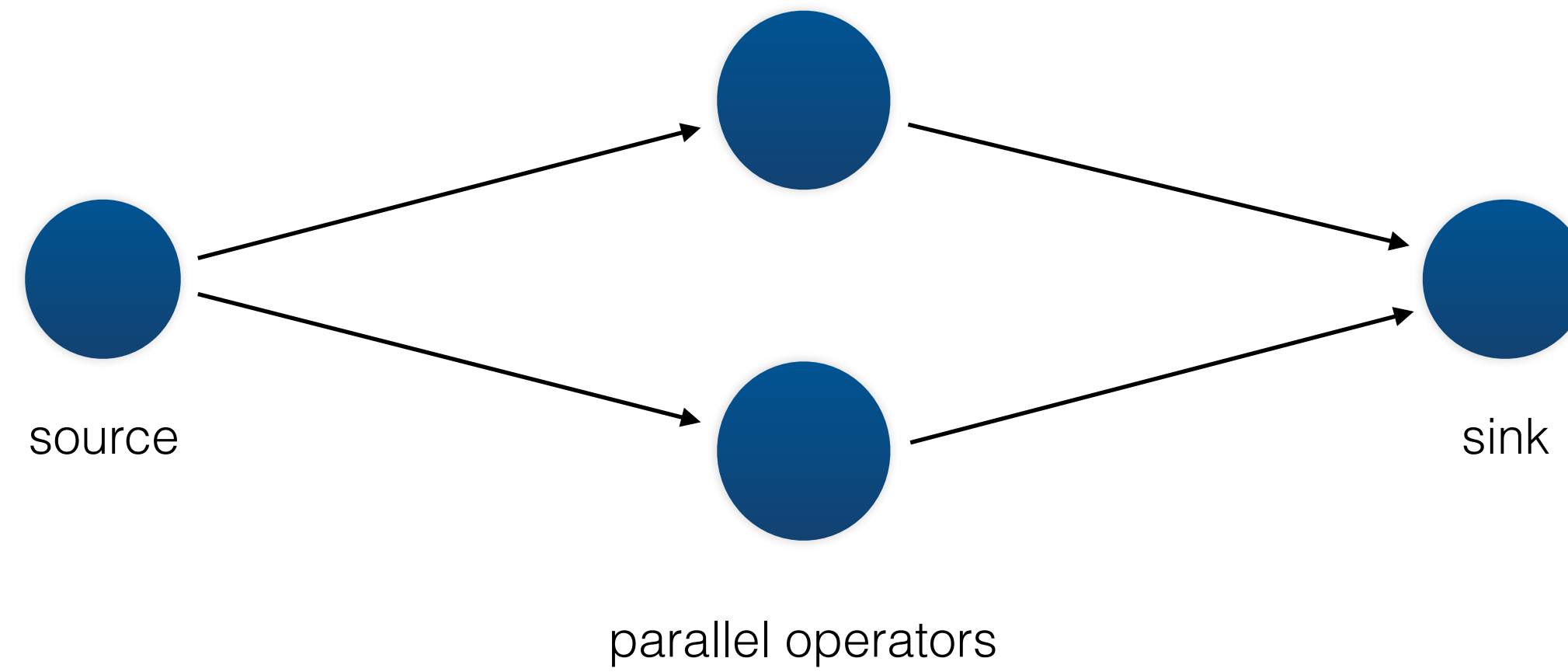~ 131 M packet records/s
*flow technology
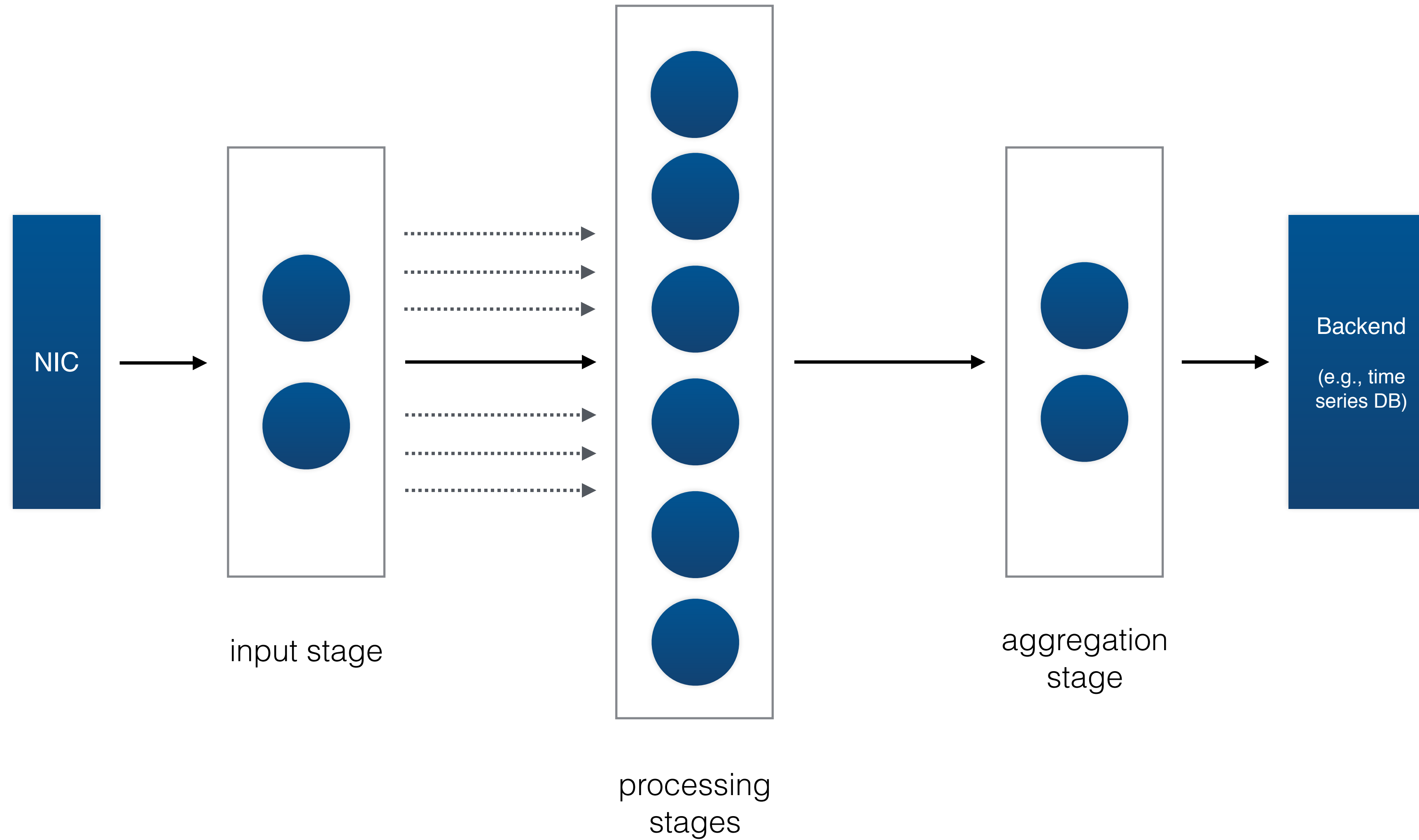
# The ideal network analytics system

*Is it possible to perform packet-level analytics on cloud-scale infrastructures without compromises?*

- per-packet records

-  x86 / general purpose programming language

- ~5M pps per core

# Leveraging parallel architectures



source

parallel operators

sink

# Leveraging parallel architectures



NIC

input stage

processing
stages

aggregation
stage

Backend

(e.g., time
series DB)

# Characteristics of packet record workloads

*Can we use properties of packet analytics workloads to our advantage?*
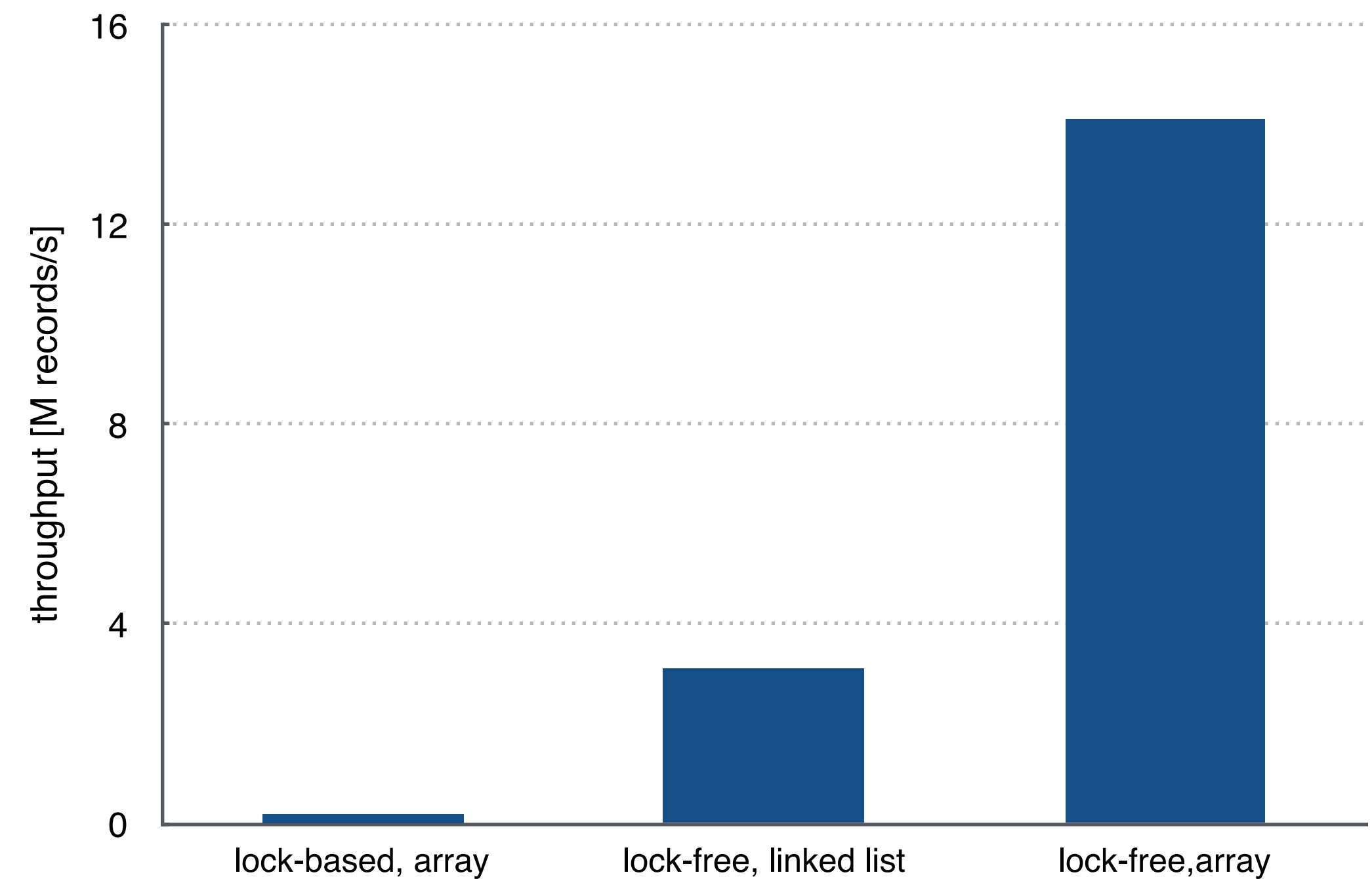
- Network attached input

- Partitionability/aggregation

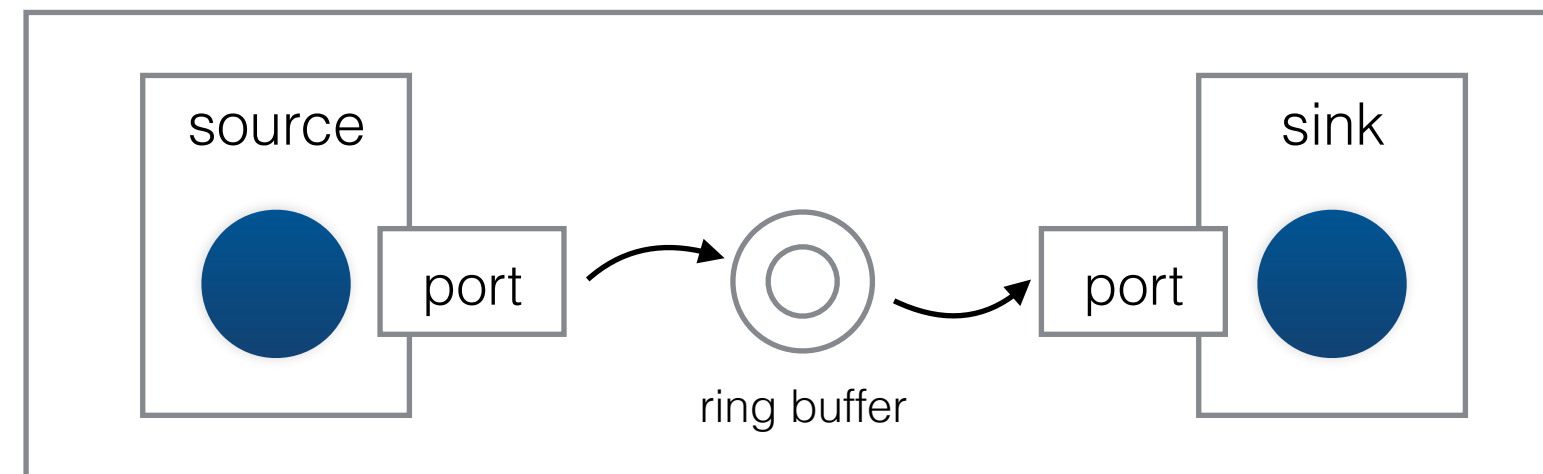- High rates, small, well-formed records

# Network attached input

# Many small records

- Array vs. linked list

- Lock-free design

- Wait-free design

- Zero-copy operations
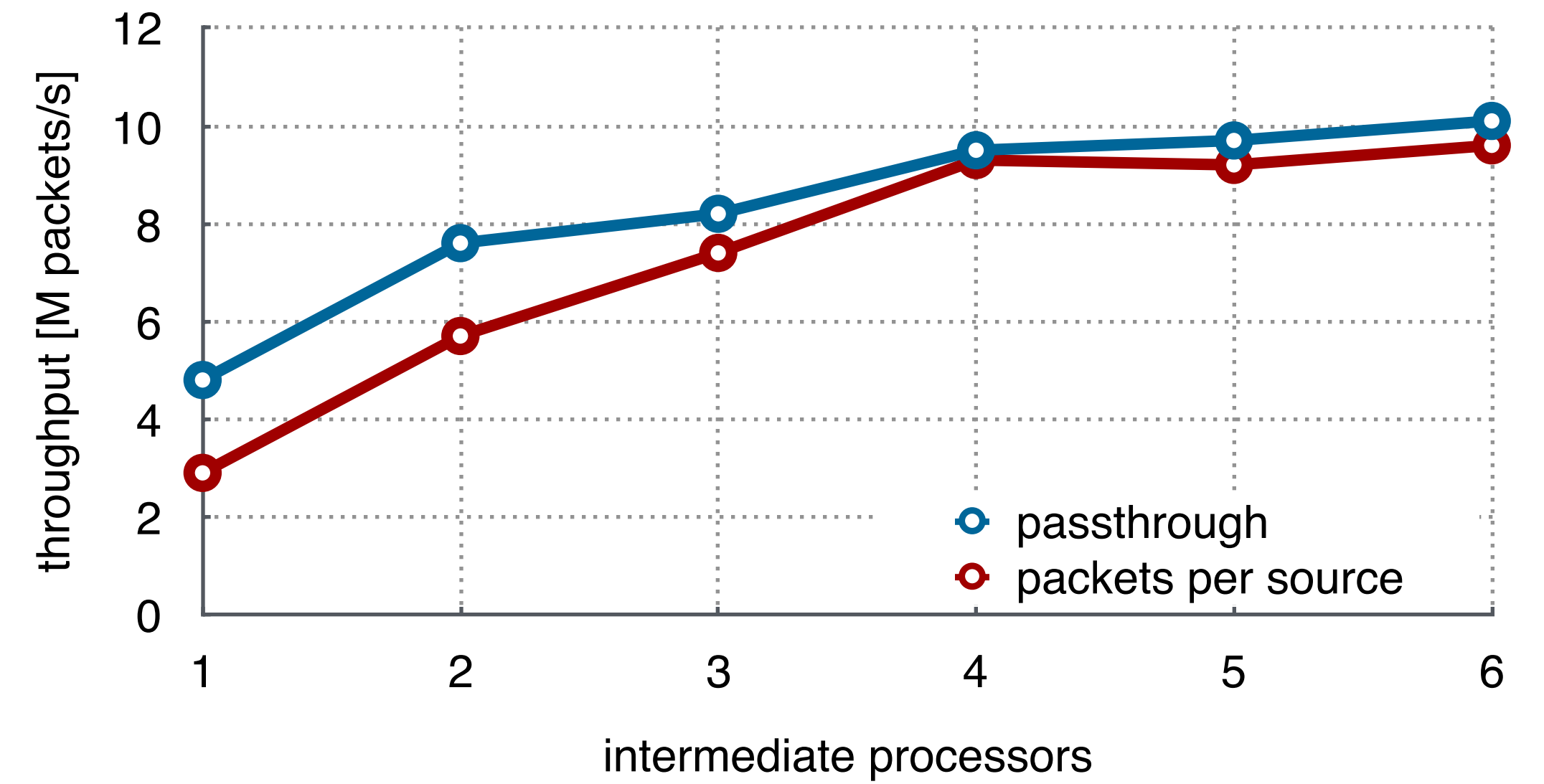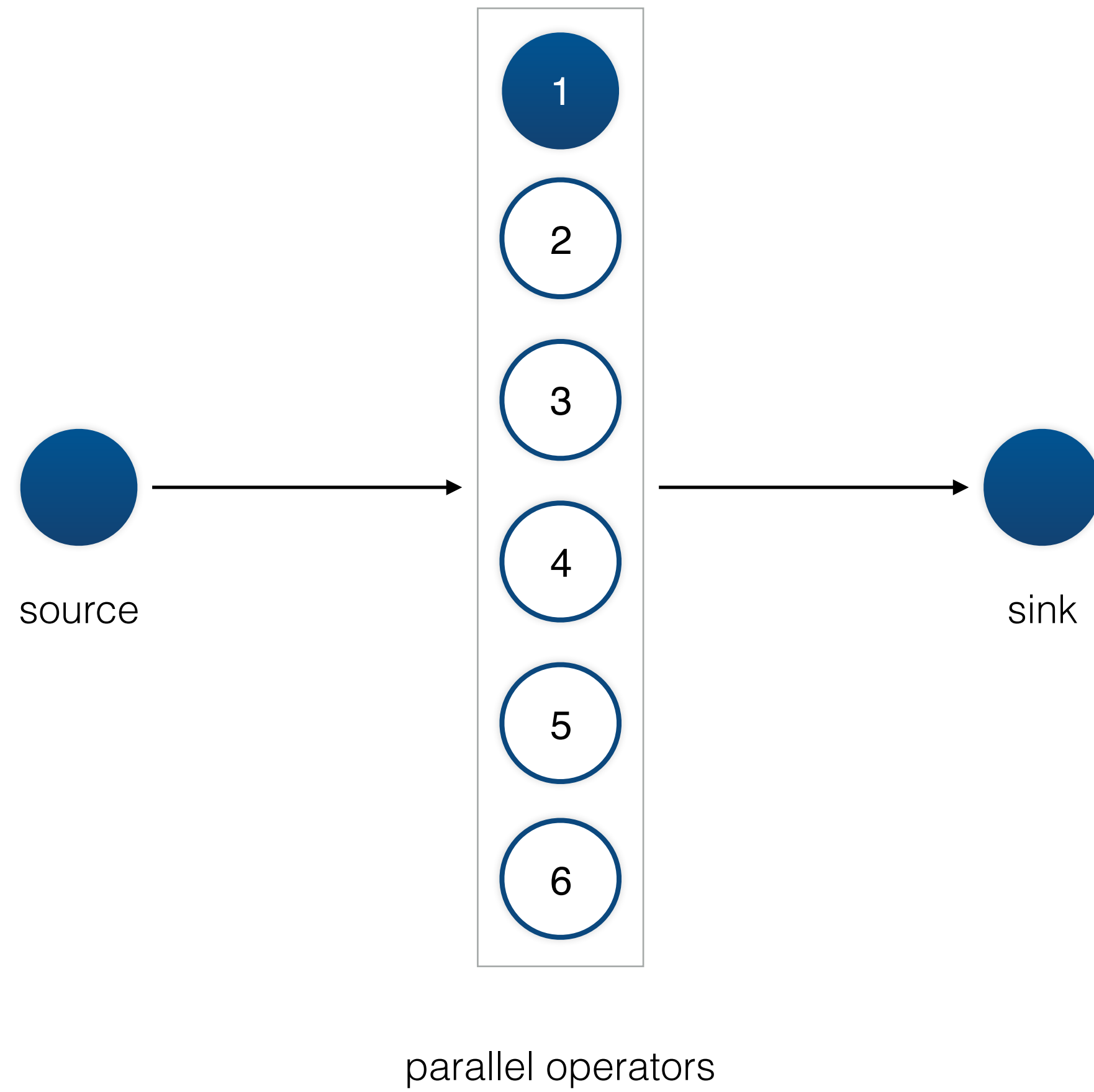
# Programming Abstraction



```
1   int main(int argc, char** argv)
2   {
3       jetstream::app app;
4       auto source = app.add_stage<source>(1, "enp6s0f0");
5       auto sink   = app.add_stage<sink>(1, std::cout);
6       app.connect<jetstream::pkt_t>(source, sink);
7       app();
8       return 0;
9   }
```
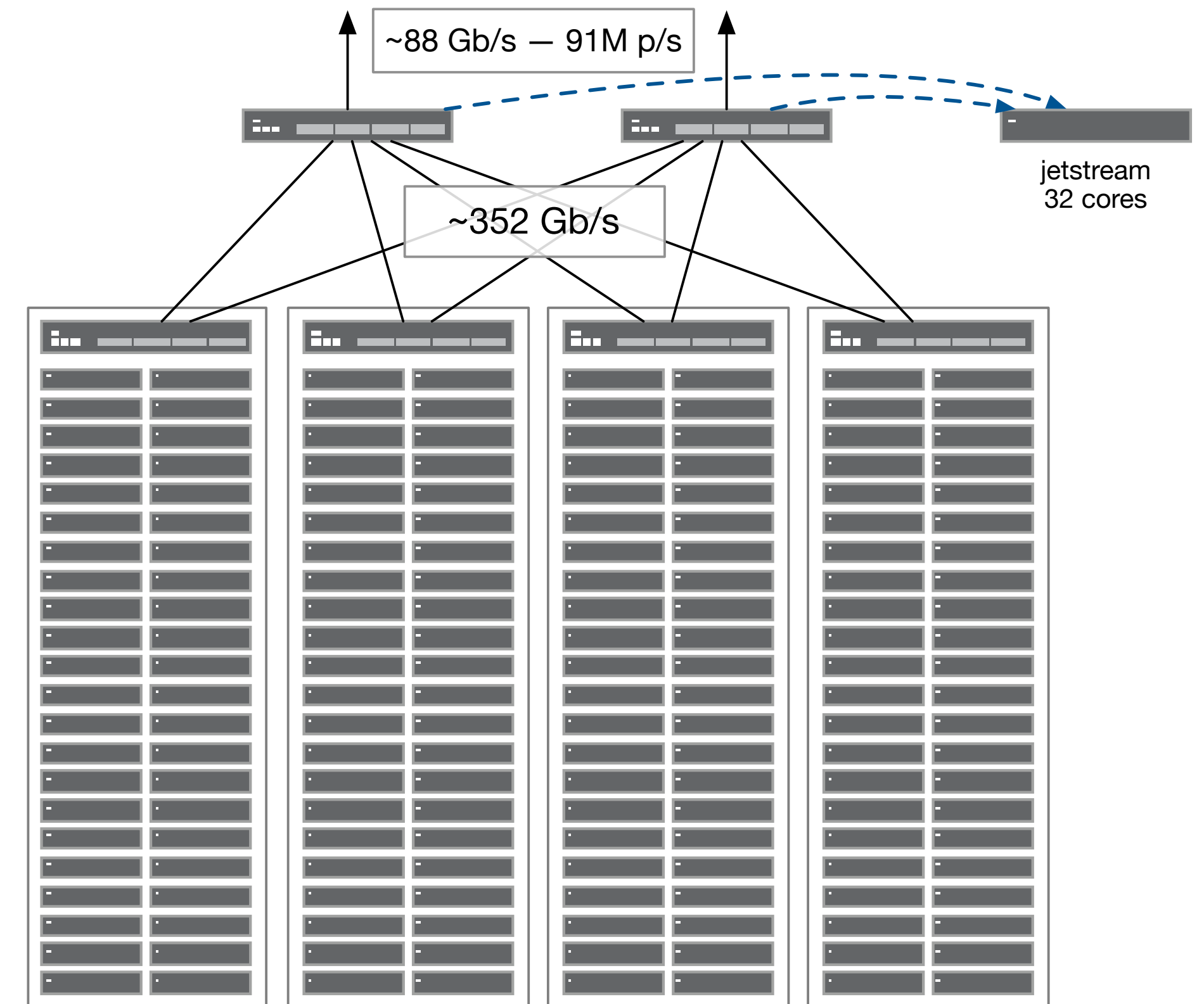
# Performance



source

parallel operators

sink

throughput [M packets/s] vs intermediate processors

- passthrough
- packets per source

# Performance



~88 Gb/s — 91M p/s

jetstream
32 cores

~352 Gb/s

- Facebook web cluster: ~ 91M egress pps

- ~32 cores for basic packet-level insight

- 176 web servers — 1 analytics server: ~0.5% of cluster capacity

[Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C. Snoeren. 2015. Inside the Social Network's (Datacenter) Network. SIGCOMM Comput. Commun. Rev. 45, 4 (August 2015), 123-137]

# Conclusion / Discussion

*Is it possible to perform packet-level analytics on cloud-scale infrastructures without compromises?*

jetstream $\longrightarrow$ high-performance, software network analytics platform

# Q&A / Discussion

Oliver Michel

oliver.michel@colorado.edu

http://nsr.colorado.edu/oliver

University of Colorado **Boulder**

The *right* approach for network monitoring and analytics?

packet-level ←→ flow-level

software ←→ hardware

What data do we *need* for monitoring/debugging?

PANEL OPENING SLIDE

# Packet-Level Analytics in Software without Compromises

Oliver Michel, John Sonchack, Eric Keller, Jonathan M. Smith
University of Colorado Boulder, University of Pennsylvania

BACKUP SLIDES

[Apache Flink]



[StreamBox Miao '18]

# Programming abstraction

Processor definition

```
1   class source : public jetstream::proc {
2     [...]
3   };
```

```
1   explicit source(const std::string& iface_name_) : proc() {
2     add_out_port<jetstream::pkt_t>(0);
3     [...]
4   }
```
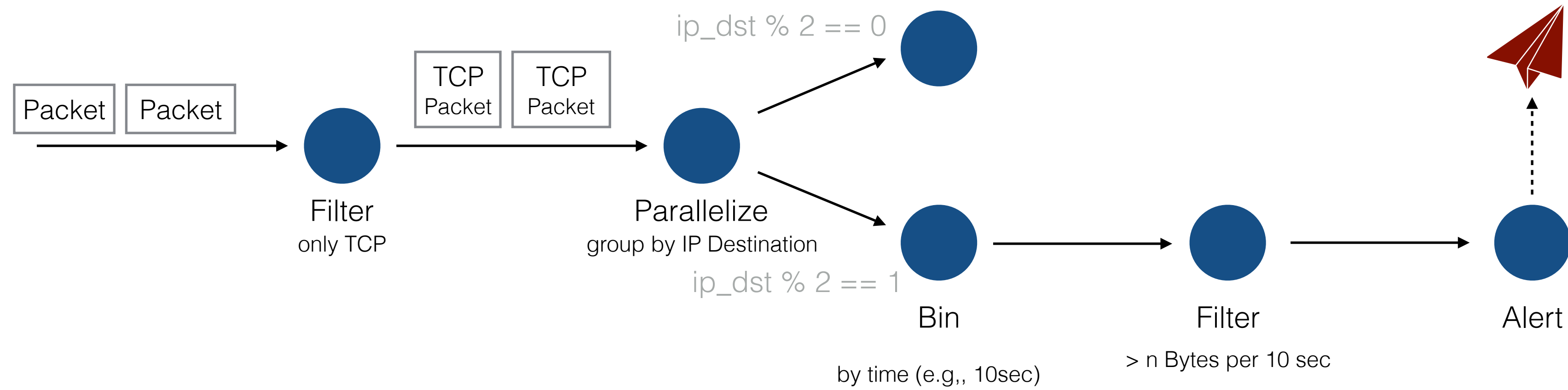
```
1   jetstream::signal operator()() override {
2       out_port<pkt_t>(0)->enqueue(read_from_nic(_pkt),
                                      jetstream::signal::continue);
3       return jetstream::signal::continue;
4   }
```

# Jetstream architecture

NUMA awareness



pipeline 1 → CPU socket 1

NIC

Backend

(e.g., time series DB)

pipeline 2 → CPU socket 2

# Stream Processing



Packet   Packet → **Filter** (only TCP) → TCP Packet   TCP Packet → **Parallelize** (group by IP Destination)

ip_dst % 2 == 0

ip_dst % 2 == 1

**Bin** by time (e.g,, 10sec) → **Filter** > n Bytes per 10 sec → **Alert**

# Reducing copy operations

Packet Buffer

Pointer
Passing

queue<pkt*>                    queue<pkt*>

# Reducing copy operations

```
1  packet p;
2  p.ip_proto = 6;
3  q.enqueue(p);
```

pointer directly
into queue

Pointer
Passing

queue<pkt>

```
1  auto p = q.enqueue();
2  p->ip_proto = 6;
```

# Technologies

- Programmable switches and PISA: Protocol Independent Switch Architecture

  - Reconfigurable match-action tables in hardware

  - multiple stages with TCAM/ALU pair, fixed processing time, guarantees line rate