

# Direct Telemetry Access

**Jonatan Langlet**

*Queen Mary University of  
London*

Ran Ben Basat  
*University College London*

Sivaramakrishnan  
Ramanathan  
*University of Southern  
California*

Gabriele Oliaro  
*Harvard University*

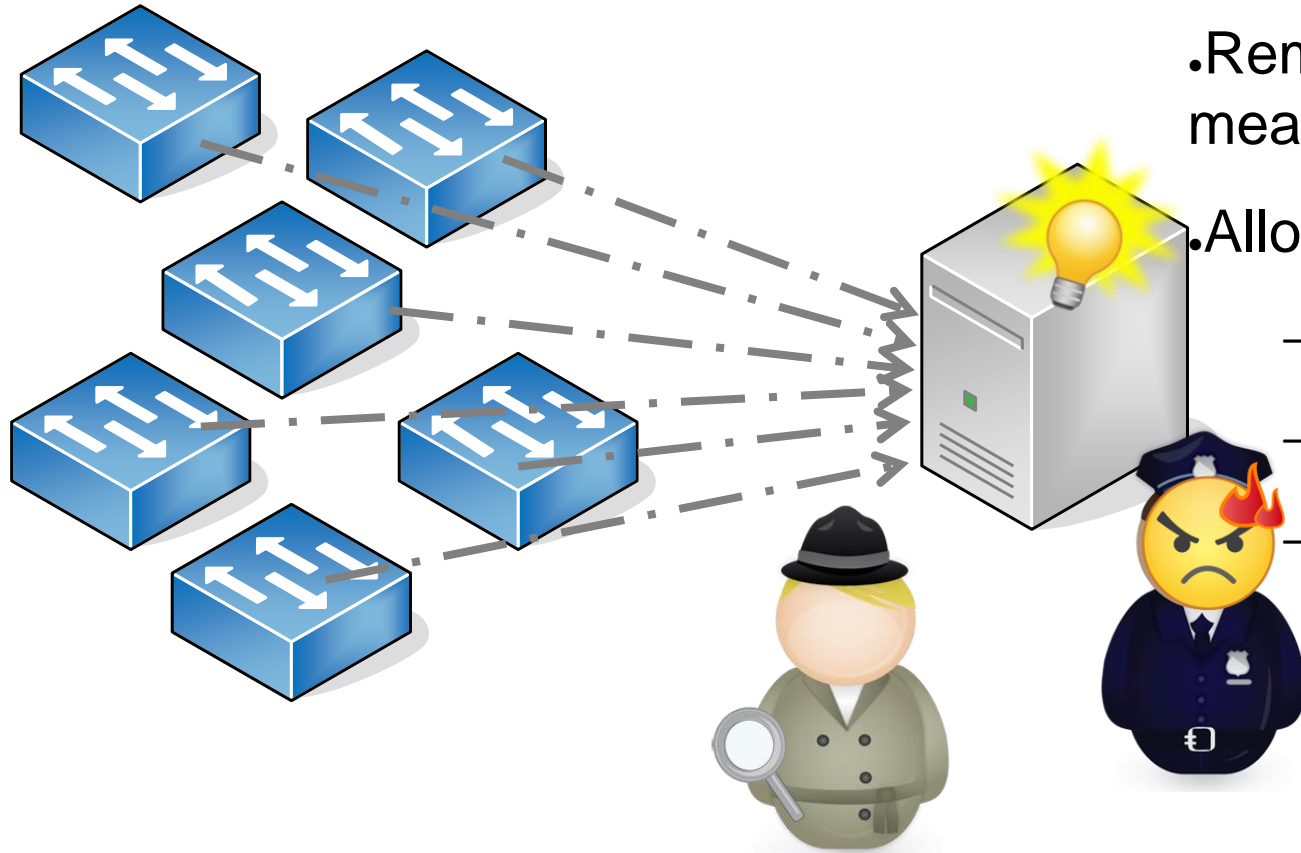
Michael Mitzenmacher  
*Harvard University*

Minlan Yu  
*Harvard University*

Gianni Antichi  
*Queen Mary University of  
London*



# What is network telemetry?



• Remotely reading network measurements

• Allows deep network insight

- Troubleshooting
- SLA compliance
- Network reactivity

# Modern network control loop

- .Flow paths
- .Queue depths
- .Latency spikes
- .Packet losses

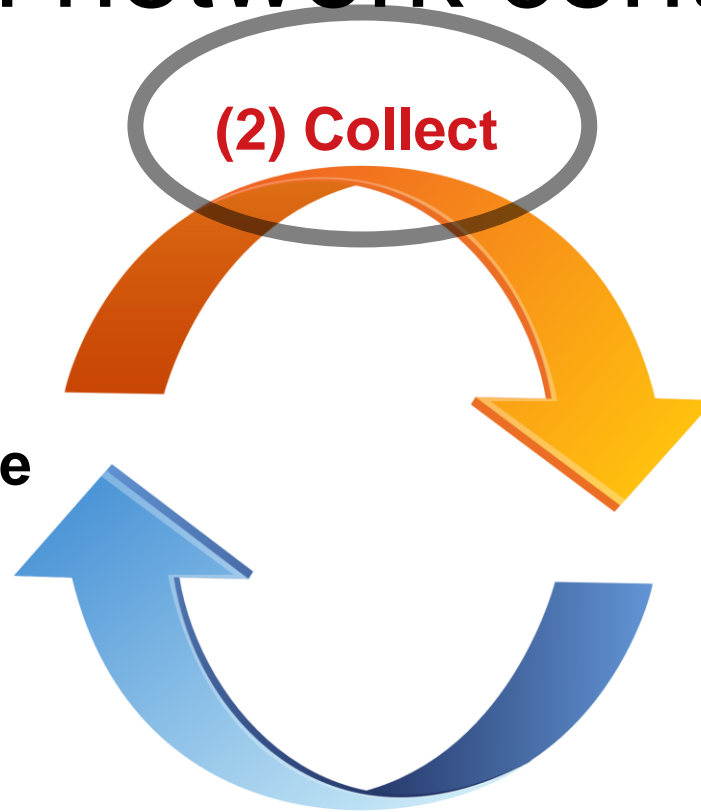
**(1) Measure**

**(2) Collect**

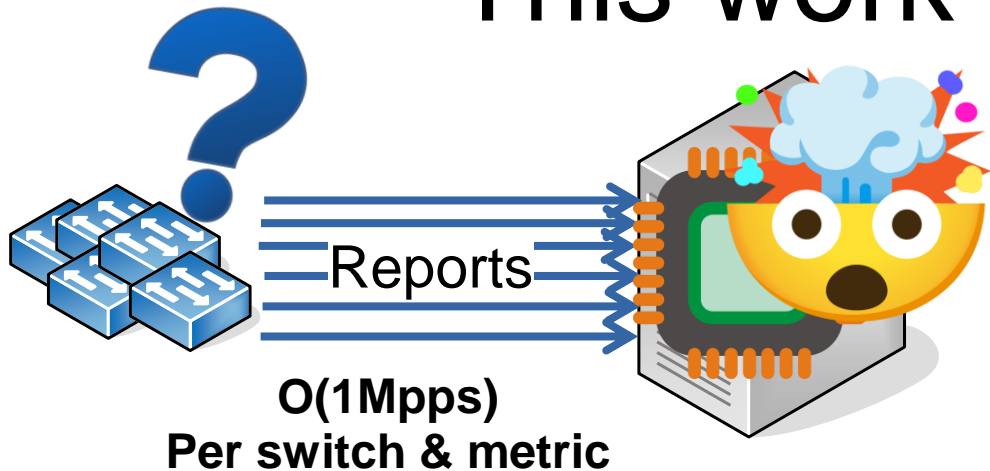
- 1)Continuously measure
- 2)Collecting for analysis
- 3)Adapting the network
- 4)Push changes

**(3) Adapt**

**(4) Reconfigure**

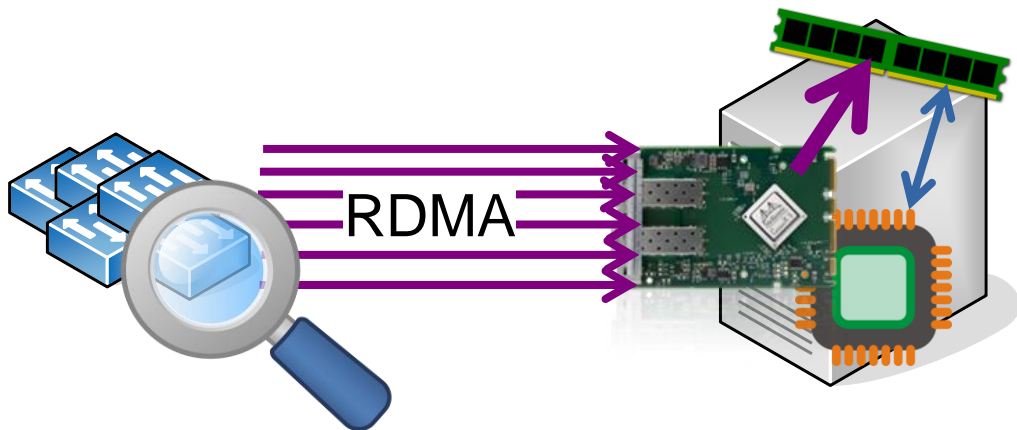


# This work in one slide



• Telemetry generates **tons** of data!

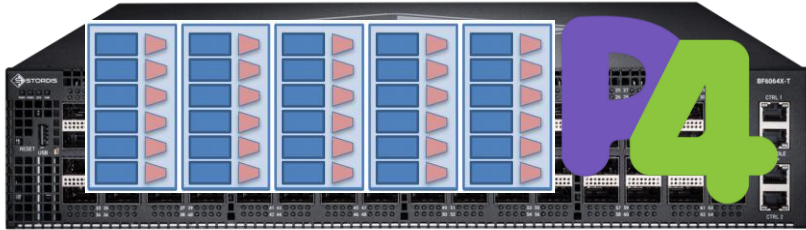
- Overloads collector CPU
- Reduced insight



• We redesign telemetry collection

- Bypassing the CPU
- Increased insight

# Why is this now a problem?



.Programmable switches enable incredible traffic insight  
 – Custom per-packet logic

.Tons of monitoring papers!

## In-band Network Telemetry via Programmable Dataplanes

Changhoon Kim\*, Anirudh Sivaraman\*\*, Nago Katta\*\*\*, Antonin Bas\*, Advait Dixit\*, Lawrence J Wobler\*  
 \*Barefoot Networks, \*\*Massachusetts Institute of Technology, \*\*\*University of California, San Diego  
 {chang, antonin, adixit, ljw}@barefootnetworks.com, anirudh@csail.mit.edu

## Flow Event Telemetry on Programmable Data Plane

Yu Zhou\*†, Chen Sun\*, Hongqiang Harry Liu\*, Rui Miao\*, Shi Bai\*, Bo Li\*, Zhilong Zheng\*†, Lingjun Zhu\*, Zhen Shen\*, Yongqing Xi\*, Pengcheng Zhang\*, Dennis Cai\*, Ming Zhang\*, Mingwei Xu†  
 \*Alibaba Group †JNSC and BMDI, ‡University of California, San Diego  
 nathan@cs.berkeley.edu, Yuliang Li, Harvard University, yuliangli@g.harvard.edu, Michael Mitzenmacher, Harvard University, michaelm@eecs.harvard.edu

## Sonata: Query-Driven Streaming Network Telemetry

Arpit Gupta, Princeton University, Rob Harrison, Princeton University, Marco Canini, KAUST, Nick Feamster, Princeton University, Jennif, Princeton

## Language-Directed Hardware Design for Network Performance Monitoring

PacketScope: Monitoring the Packet Lifecycle Inside a Switch  
 Ross Teixeira, Princeton University, rapt@cs.princeton.edu, Arpit Gupta, UC Santa Barbara, arpitgupta@cs.ucsb.edu, Rob Harrison, United States Military Academy, Vikram Nathan<sup>1</sup>, Prateesh Goyal<sup>1</sup>, umar Jeyakumar<sup>3</sup>, Changhoon Kim<sup>4</sup>

## One Sketch to Rule Them All: Rethinking Network Flow Monitoring with UnivMon

Zaoxing Liu<sup>†</sup>, Antonis Manousis\*, Gregory Vorsanger<sup>†</sup>, Vyas Sekar\*, Vladimir Braverman<sup>†</sup>  
<sup>†</sup> Johns Hopkins University · Carnegie Mellon University

# Telemetry is INTENSE

System	Metric	Generation Rate
INT 0.5%	Raw	~19 Mpps
Marple	TCP Out-of-Sequence	~6.7 Mpps
Marple	Flow counting	~4.3 Mpps
NetSeer	Packet losses	~1.0 Mpps

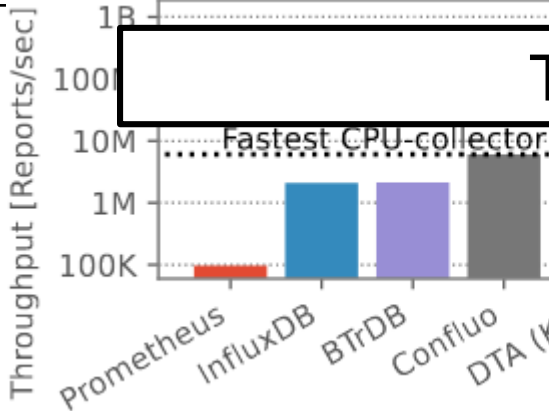
.Single monitoring query  
 -  $O(1M)$  per switch query

.Collection performance

Max  $O(10M)$

The collector CPU is bottlenecking!

Spoiler



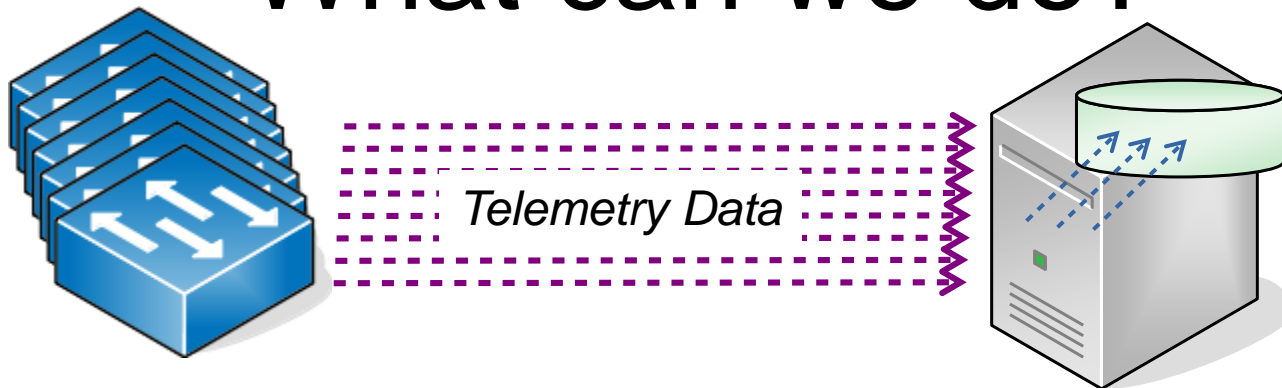
Narayana, Srinivas, et al., "Language-Directed Hardware Design for Network Performance Monitoring", SIGCOMM'17

Zhou, Yu, et al., "Flow Event Telemetry on Programmable Data Plane", SIGCOMM'20

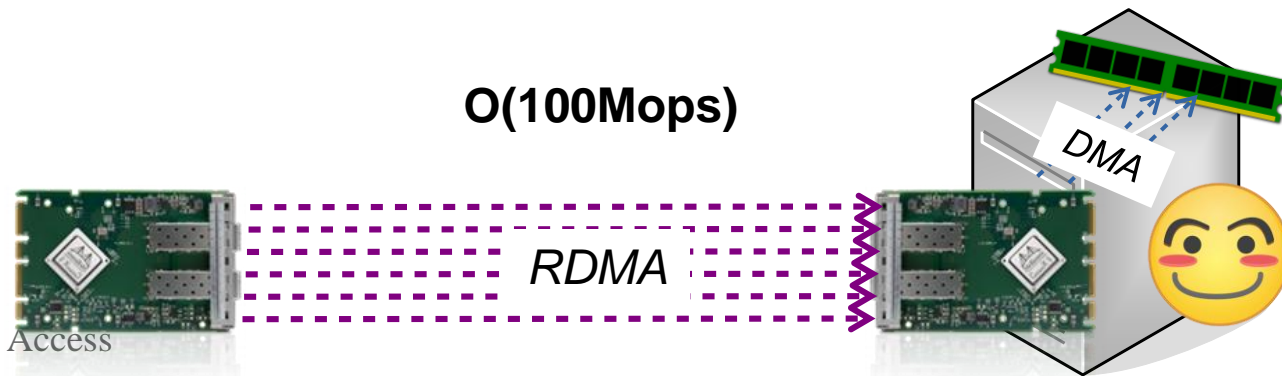
Khandelwal, Anurag, et al., "Confluo: Distributed Monitoring and Diagnosis Stack for High-speed Networks", NSDI'19

Fig: 16-core performance

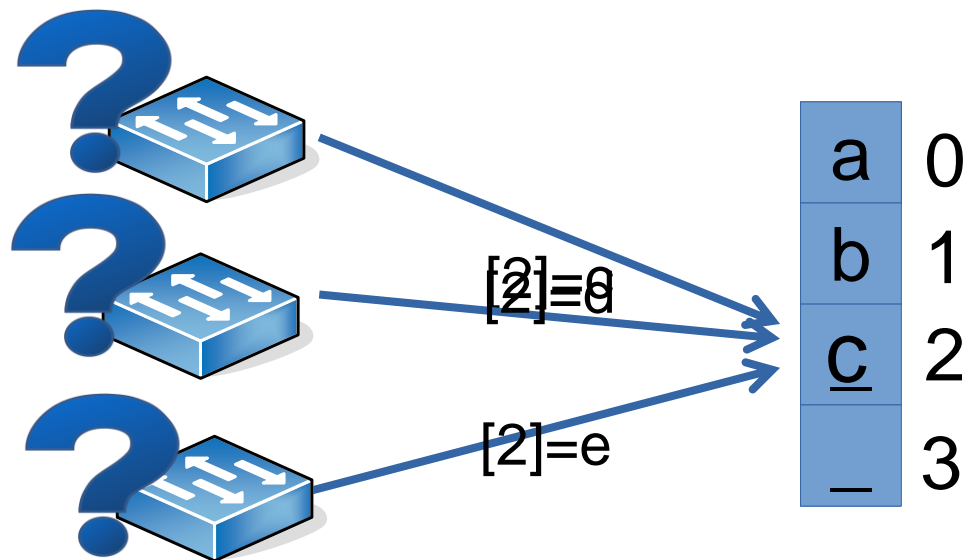
# What can we do?



This is what **RDMA** is designed for!



# But RDMA is limited...



• Two high-speed operations

– Read

– Write

• High-speed RDMA is limited

– Dumb memory operations

– Explicit addresses



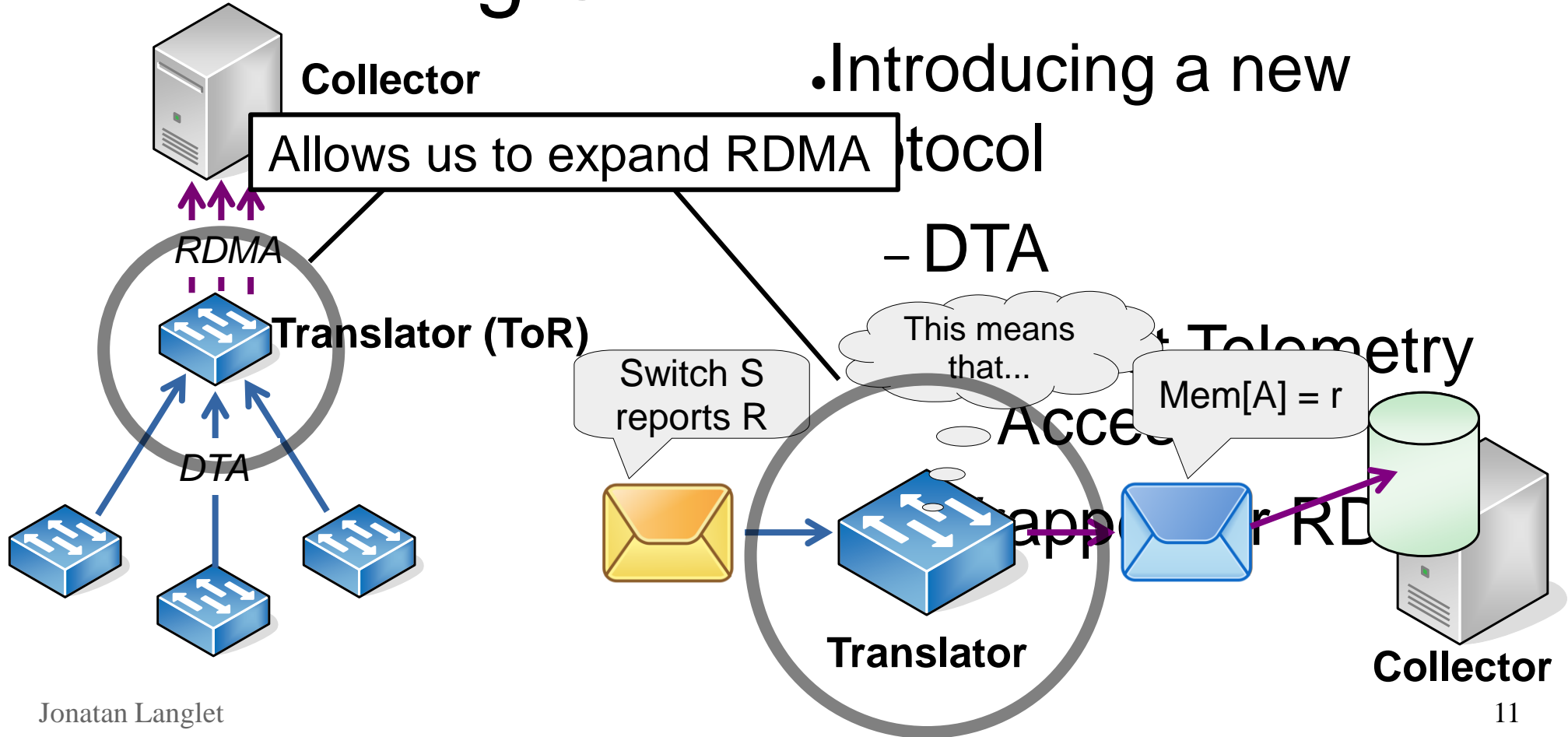
# Making Switch-RDMA work

Collector

Introducing a new

Allows us to expand RDMA protocol

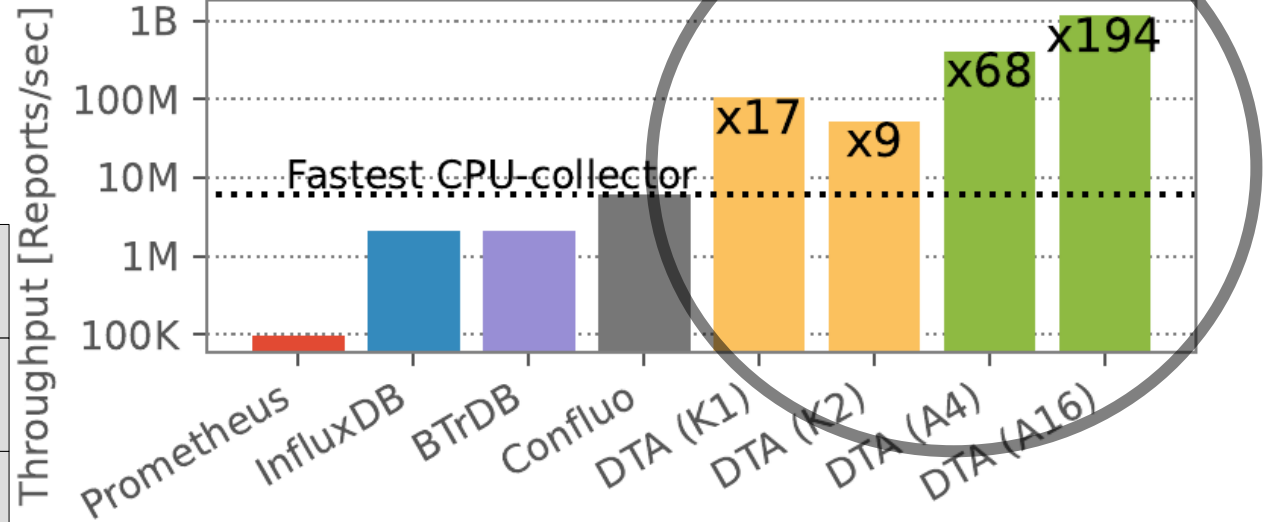
- DTA



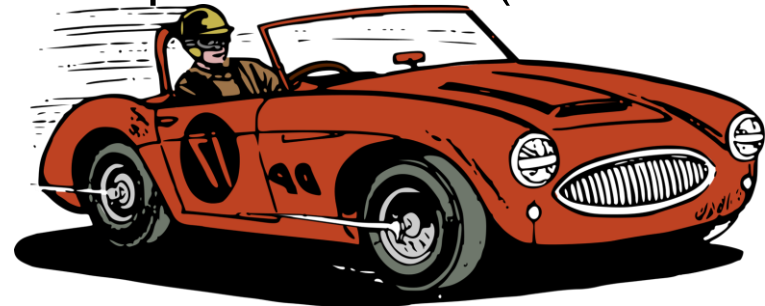


# Fast and Generic

<b>INT</b>	<b>HyperLogLog</b>
<b>Count Sketch</b>	<b>CM Sketch</b>
<b>Sonata</b> (SIGCOMM'18)	<b>Marple</b> (SIGCOMM'17)
<b>AROMA</b> (IFIP'20)	<b>TurboFlow</b> (EuroSys'18)
<b>NetSeer</b> (SIGCOMM'20)	<b>PacketScope</b> (SOSR'20)
	<b>PINT</b> (SIGCOMM'20)

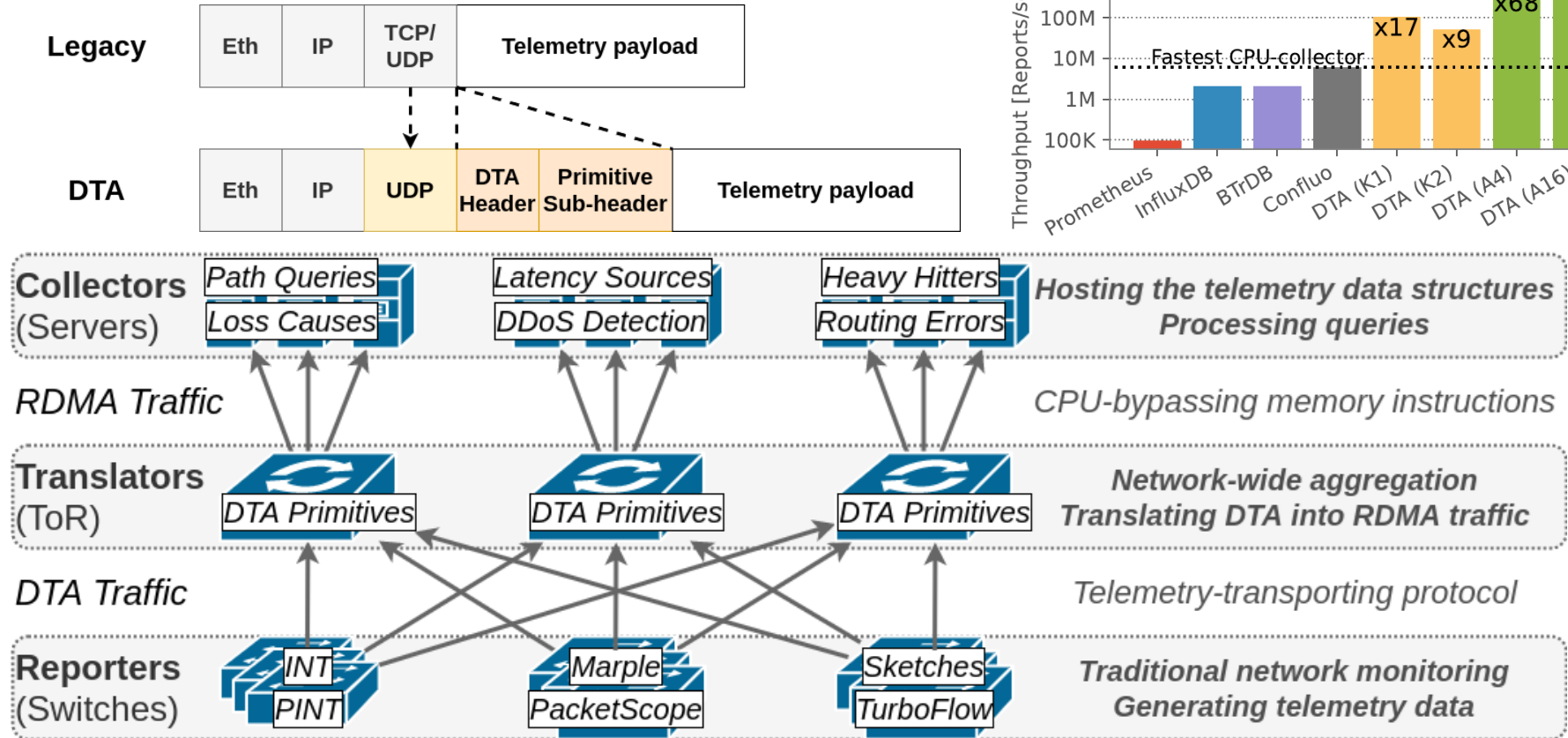


**Fig:** 16-core performance (vs 0-core DTA)

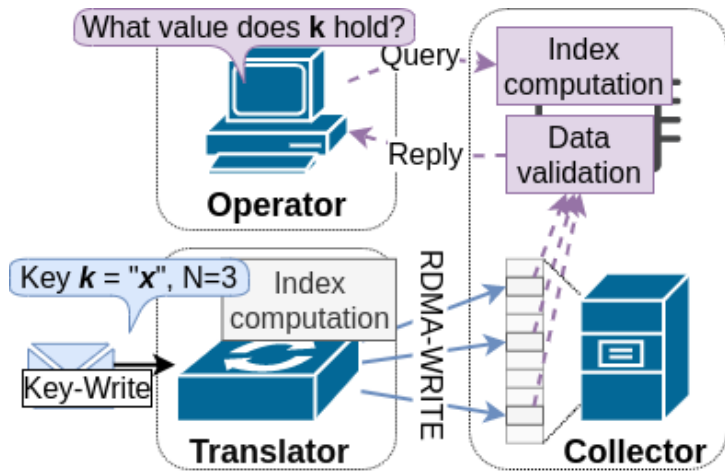


# Backup slides

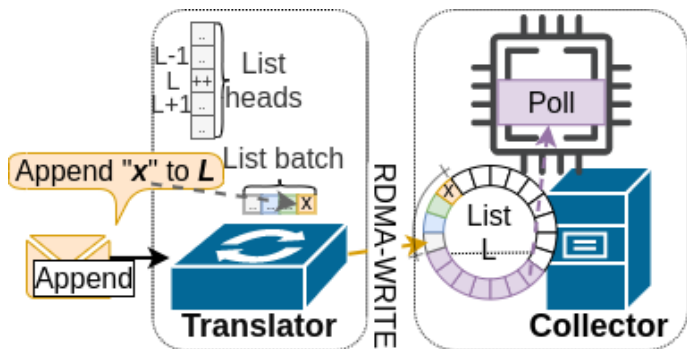
# Overview



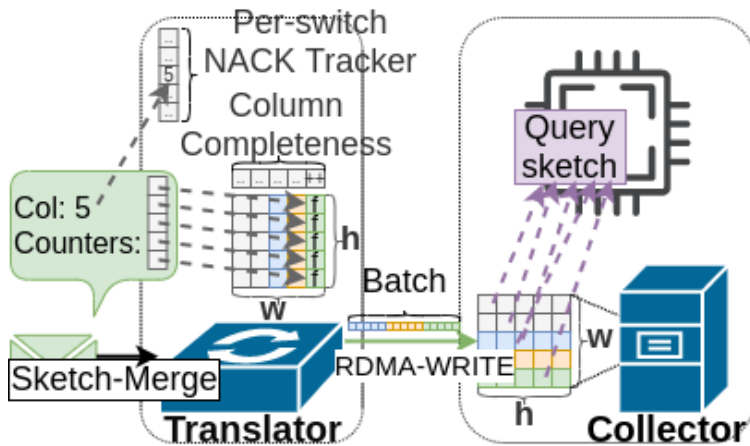
# Primitives



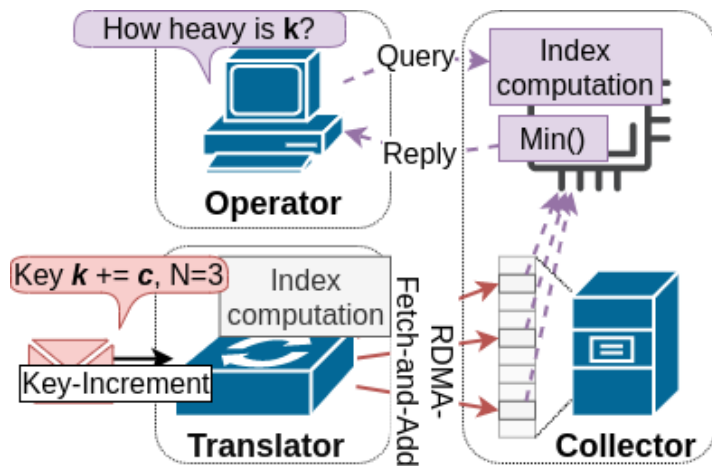
**Fig: Key-Write**



**Fig: Append**

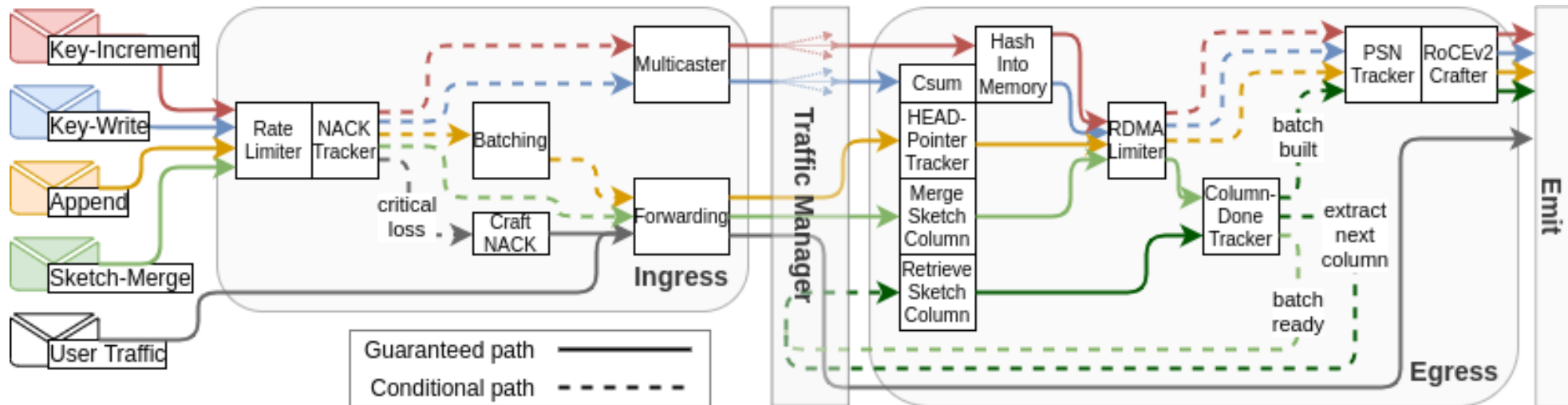


**Fig: Sketch-Merge**



**Fig: Key-Increment**

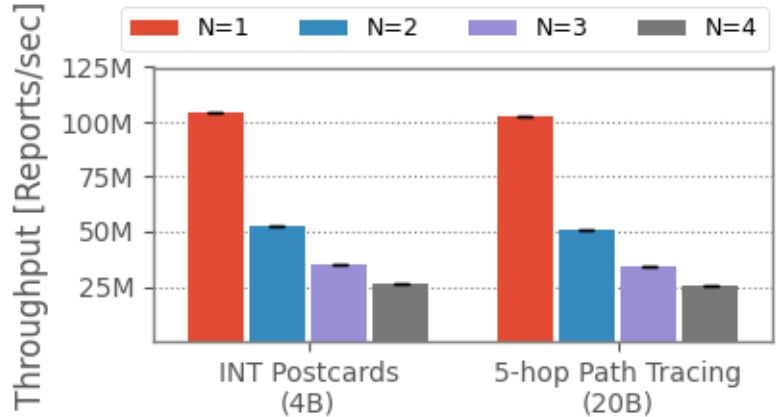
# ASIC Pipeline



Resource	Base footprint	Batching	Retransmission
SRAM	5.5%	+3.0%	+0.5%
Match Crossbar	5.9%	+9.0%	+0.2%
Table IDs	27.6%	+7.8%	+1.0%
Hash Dist Unit	13.9%	+20.8%	-
Ternary Bus	20.3%	+7.8%	+1.1%
Stateful ALU	10.4%	+31.3%	+2.1%

Table 3: Resource costs of the translator. Append batching creates batches of 16x4B data payloads, and retransmission supports tracking 65K reporter sequence numbers with 256 in-transit reports each.

# Querying



Telemetry Data

Fig: Key-Write speed

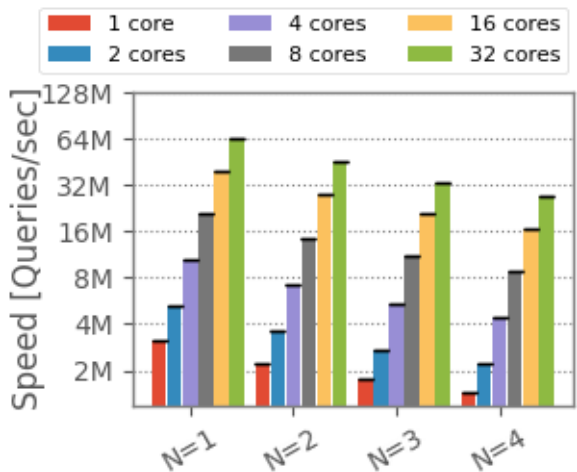


Fig: Key-Write Querying

Jonatan Langlet  
Direct Telemetry Access

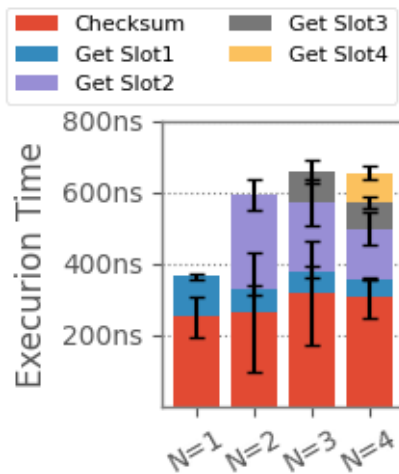
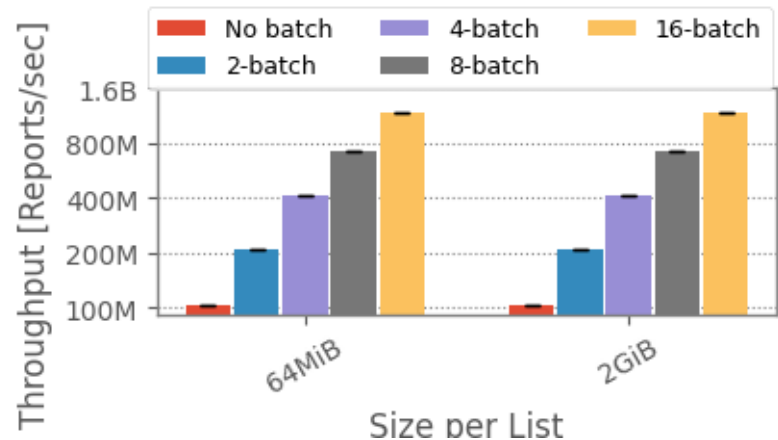


Fig: Key-Write breakdown



Size per List

Fig: Append speed

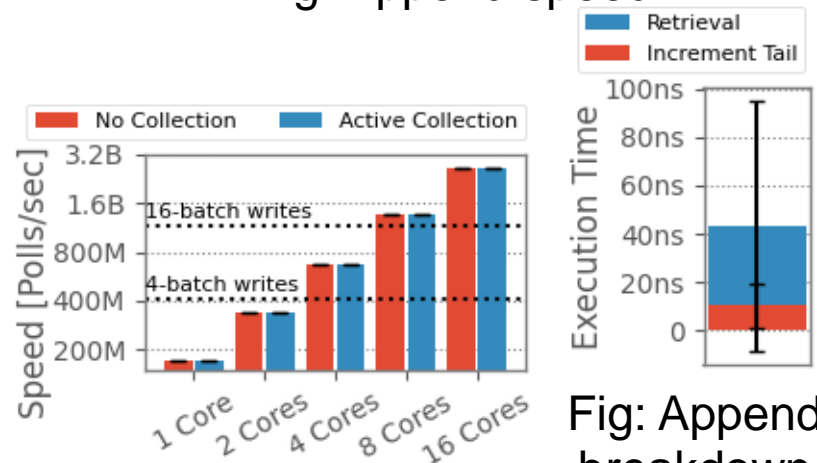
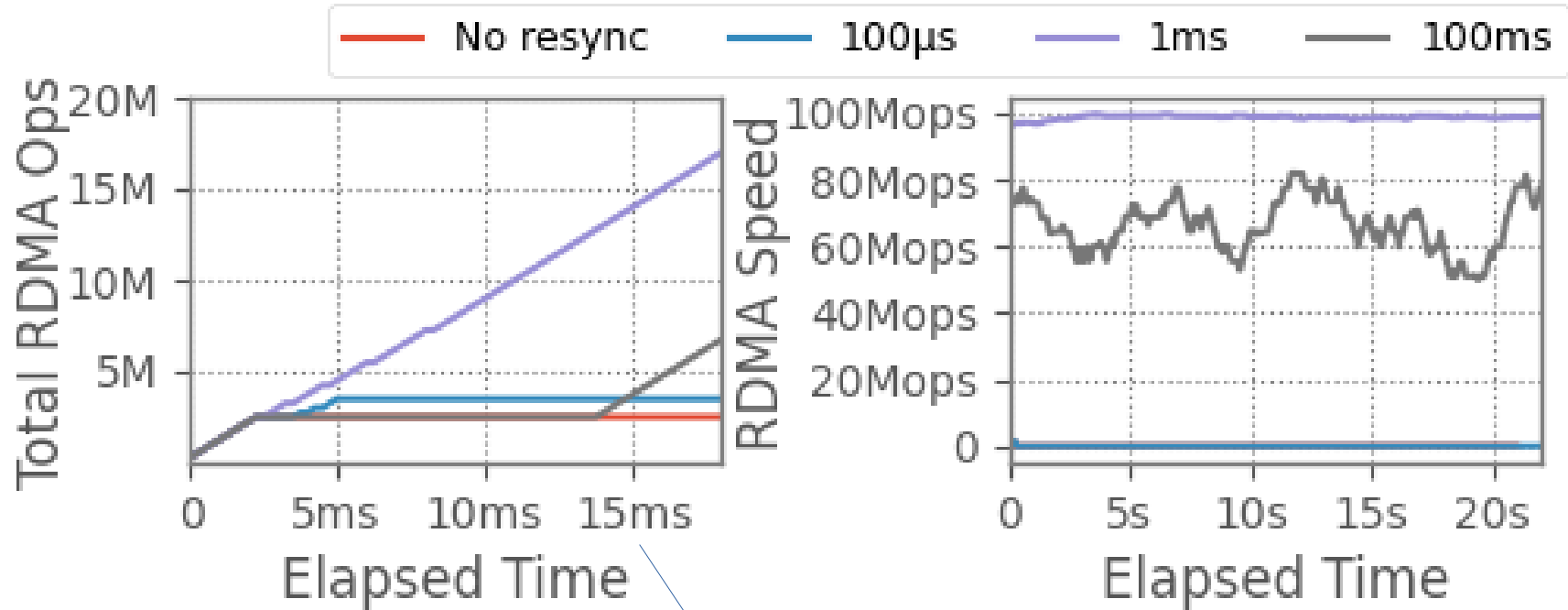


Fig: Append querying

Fig: Append breakdown

# Resync



x10 :)