

# BurstRadar

Practical Real-time Microburst Monitoring for  
Datacenter Networks

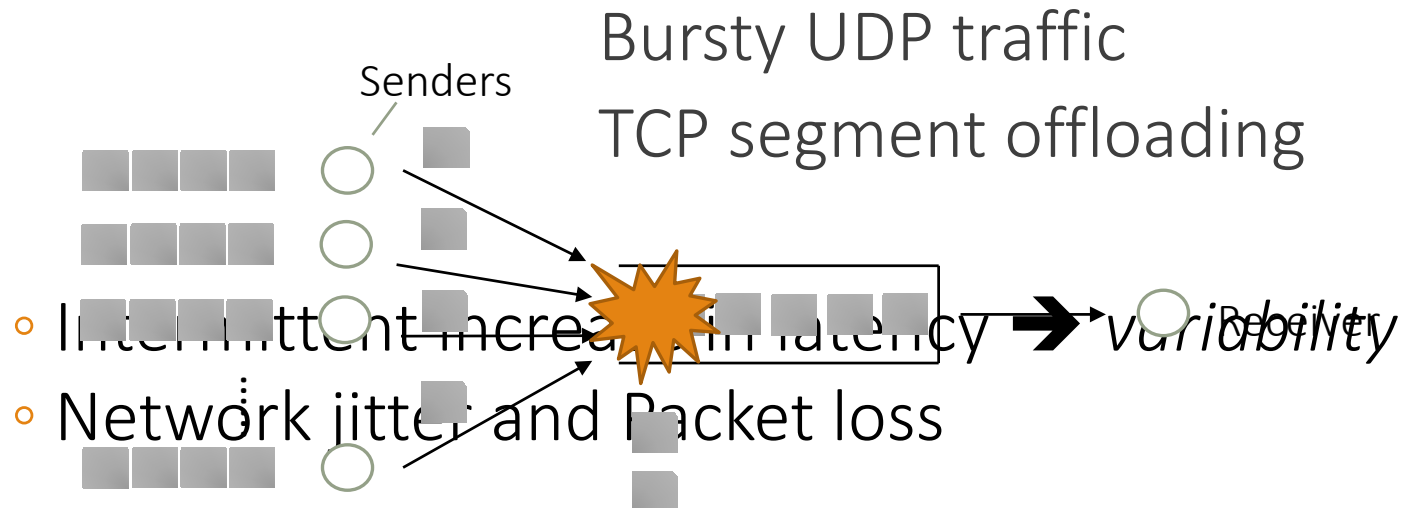
Raj Joshi<sup>1</sup>, Ting Qu<sup>2</sup>, Mun Choon Chan<sup>1</sup>, Ben Leong<sup>1</sup>, Boon Thau Loo<sup>3</sup>



# Microbursts ( $\mu$ bursts)

Events of **intermittent congestion** lasting 10's or 100's of  $\mu$ s

- Common Causes: TCP Incast

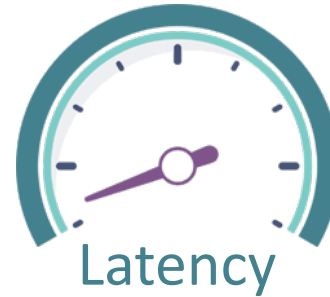


# Modern Datacenter Networks

---



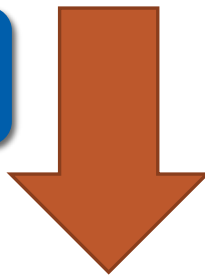
> 10 Gbps



10's of  $\mu$ s

Small amounts of queueing (**microbursts**):

Performance



# Modern Datacenter Networks

---



Detect the occurrence of  $\mu$ bursts  
&  
identify the contributing flows!

# Detecting & characterizing $\mu$ bursts is hard

---

## Measurement study from FB's datacenter

- Last for less than 200  $\mu$ s
- Occur unpredictably

## Traditional sampling-based techniques

- Cannot even detect microbursts



## Commercial Solutions

- Can detect the occurrence of microbursts
- Provide no information about the cause



# New Advancements

---

Programmable dataplanes  
and dataplane telemetry



## In-band Telemetry (INT)

- Adds queuing telemetry info into packets & exports it to monitoring servers from the last-hop switches

# Challenges: Effective & real-time monitoring

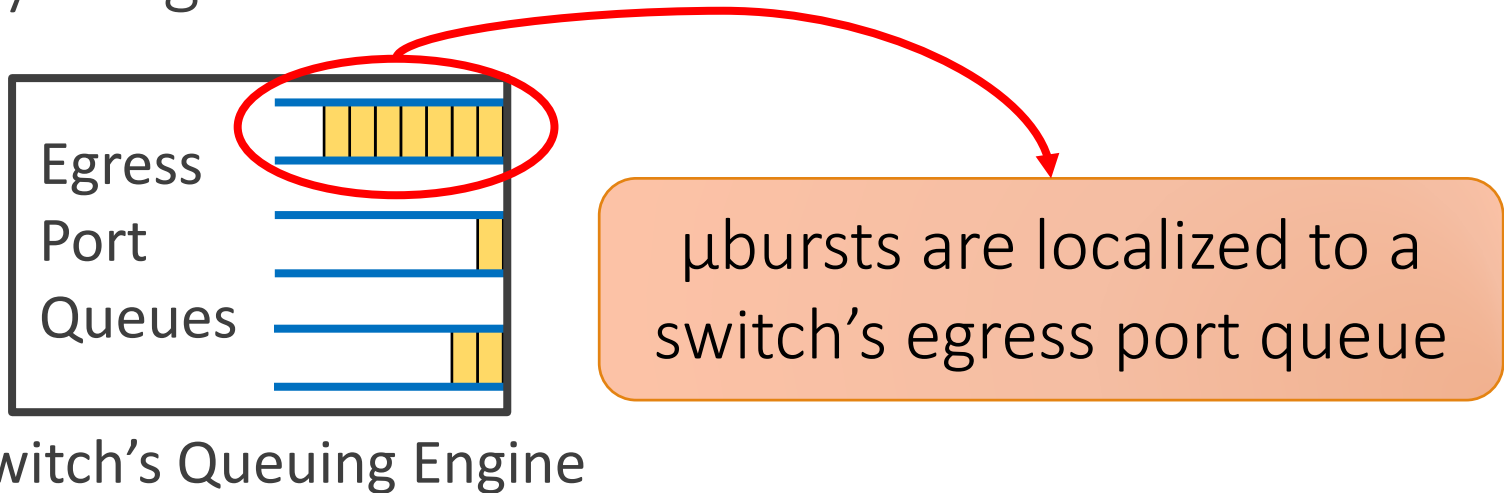
## Using INT to detect $\mu$ bursts is wasteful

- Need to capture and export/process telemetry data for **all packets**
  - Since  $\mu$ bursts are unpredictable
- Expensive computation and delay
  - Correlate monitoring data from different points in the network

# Solution: Out-band

---

Key Insight:



Key Idea:

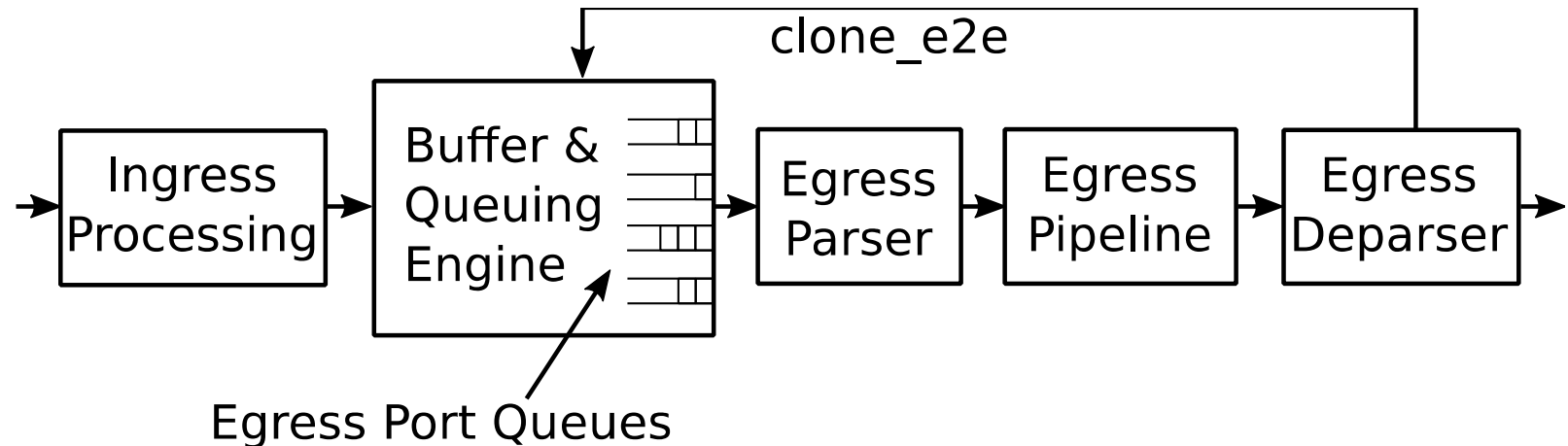
- We can detect the microburst directly on the switch where it happens



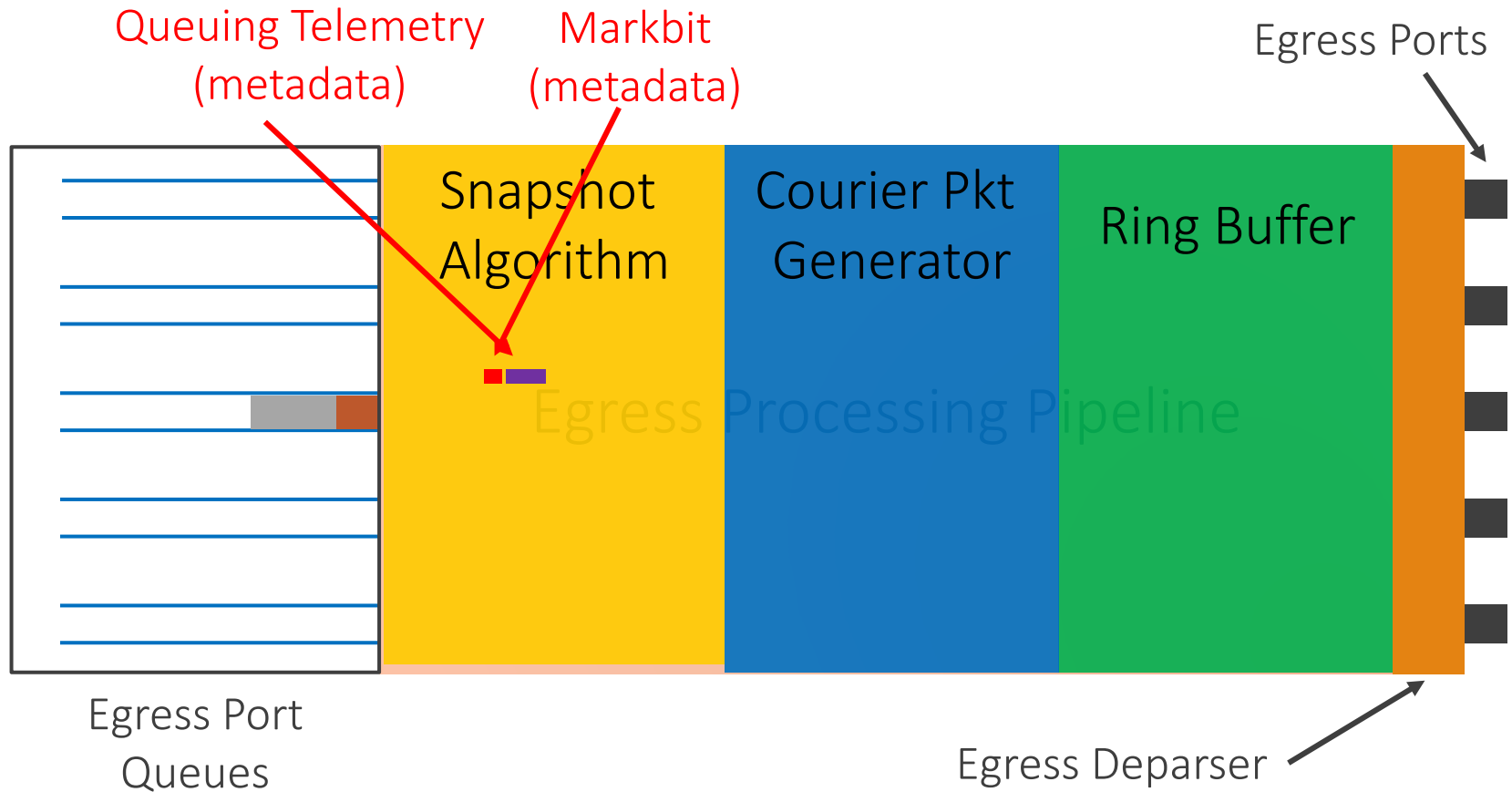
# Solution: egress pipeline

---

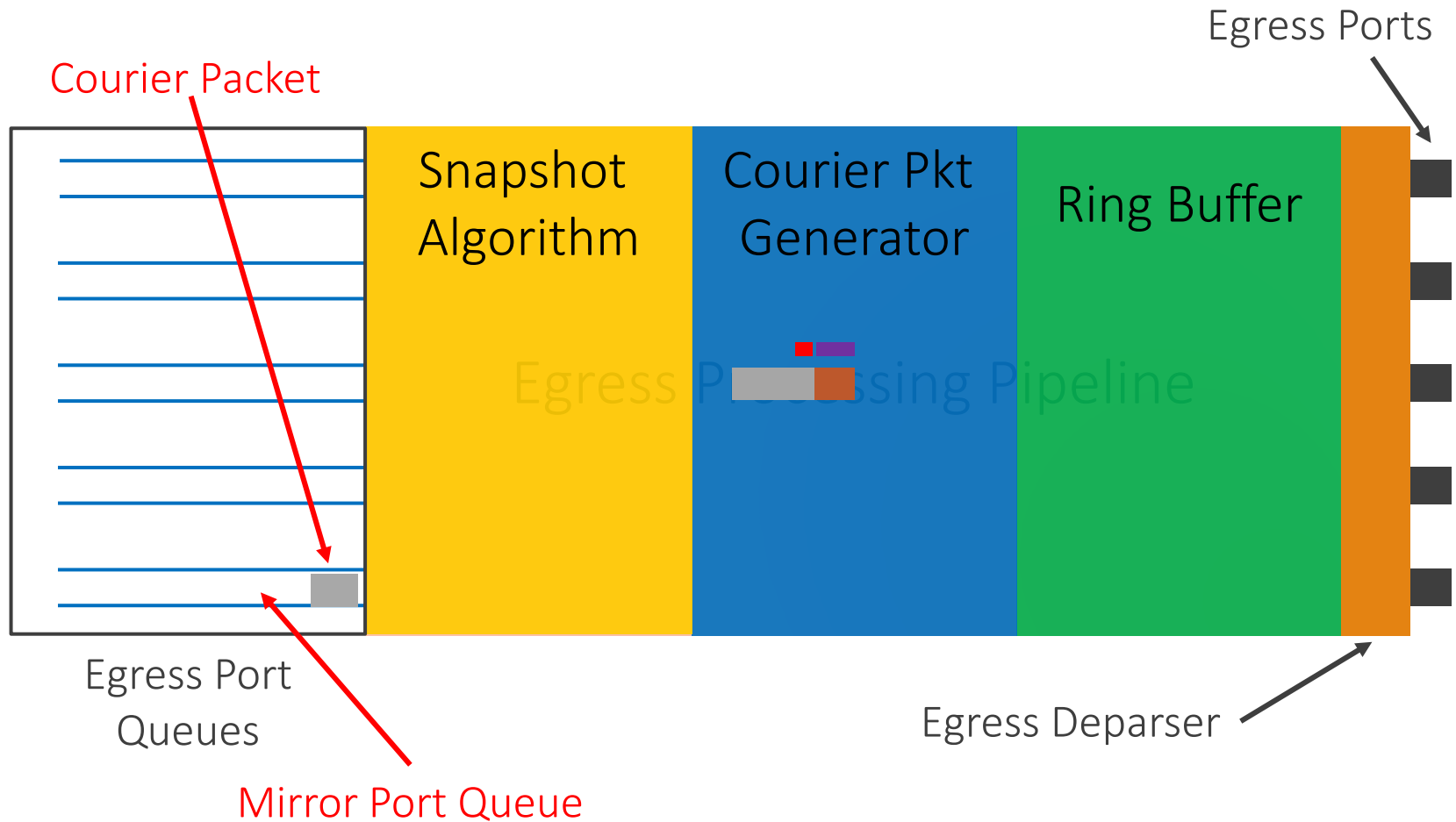
- Switching ASIC's "Buffer and Queuing Engine" (BQE) does not provide any support to peek into the contents of any queue



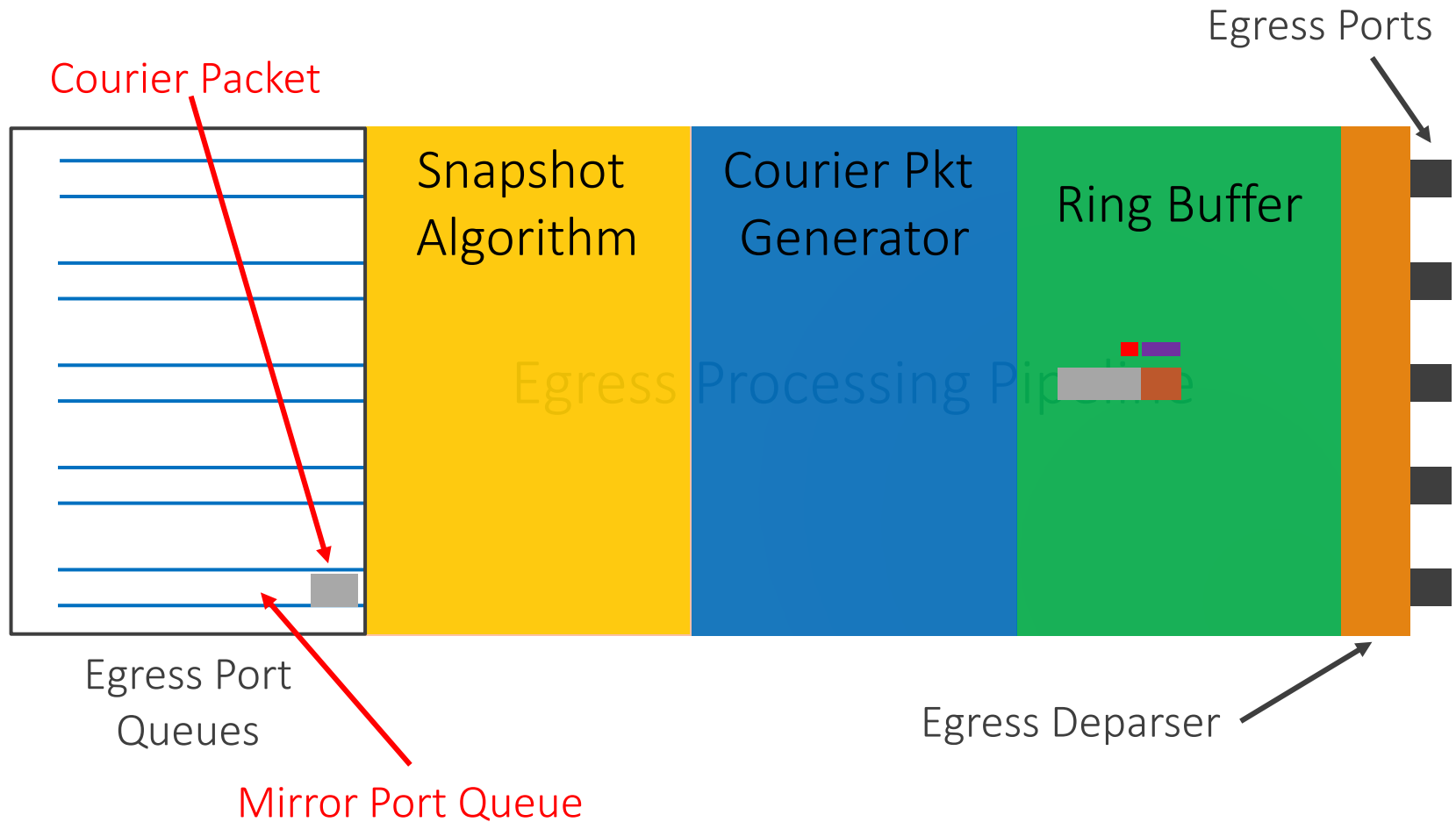
# BurstRadar Overview



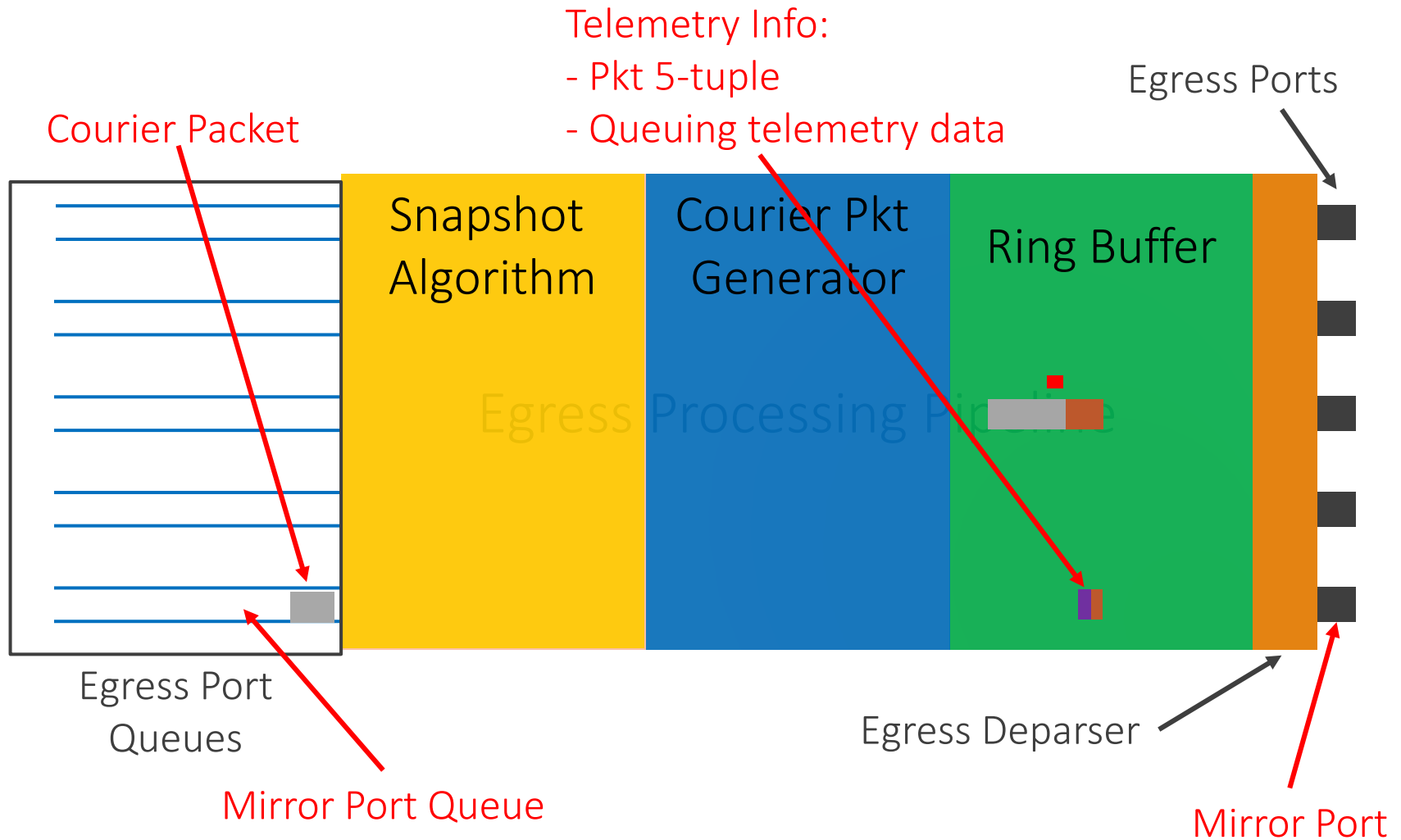
# BurstRadar Overview



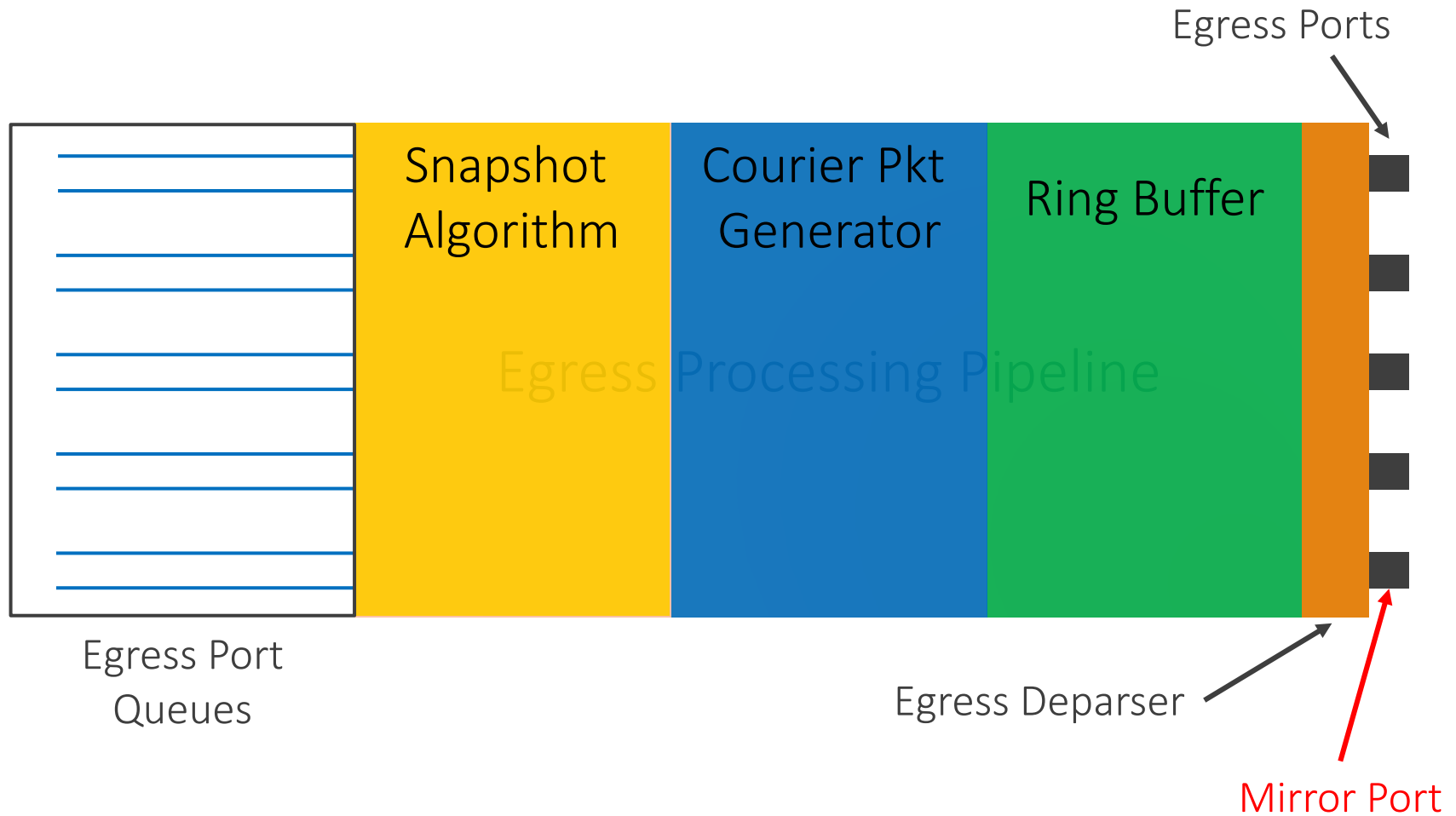
# BurstRadar Overview



# BurstRadar Overview



# BurstRadar Overview



Snapshot Algorithm



Courier Pkt Generator



Ring Buffer



“Snapshot” the telemetry info of *only* the packets involved in  $\mu$ bursts

Telemetry info: 5-tuple (packet header)

ingress/egress timestamps

enqueue/dequeue queue depths

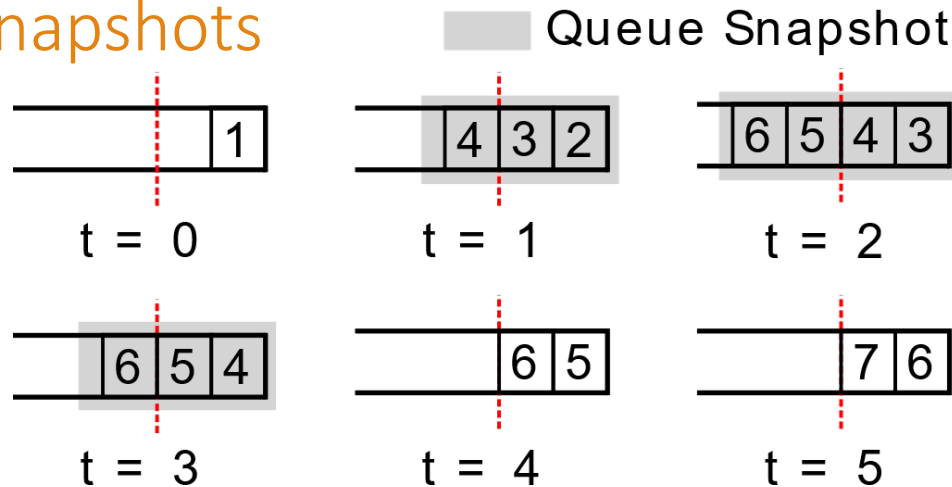
} (metadata)



## Latency Increase Threshold

- Eg. RTT =  $50\mu\text{s}$ , threshold = 30%, i.e., maximum delay =  $15\mu\text{s}$

## Queue Snapshots



## Snapshot Algorithm

- Each packet reports `deqQdepth`
- if `deqQdepth > threshold`, then mark pkt → snapshot
  - Track remaining bytes in the queue
- elif tracked bytes still remaining then mark pkt → snapshot

# Snapshot Algorithm



# Courier Pkt Generator



# Ring Buffer



“Courier” Packets transport the telemetry info via the switch’s mirror port (out-of-band)

All the data stays *together* → Avoids the expensive correlation on the monitoring servers

Each marked packet → generate new courier packet

clone egress to egress, `clone_e2e`

- Copy of the exiting marked packet
- Place it in the egress queue of the mirror port

# Snapshot Algorithm



# Courier Pkt Generator



# Ring Buffer



“Ring Buffer” temporarily stores the telemetry info of marked packets until they can be copied into the courier packets.

# Evaluation

---

# Evaluation Setup

---

## Hardware Testbed



BurstRadar Prototype



Send/Receive  $\mu$ burst Traffic

- About 550 lines of p4 code

## Generated $\mu$ burst Traffic Traces

- $\mu$ bursts data for “web” and “cache” traffic [IMC '17]

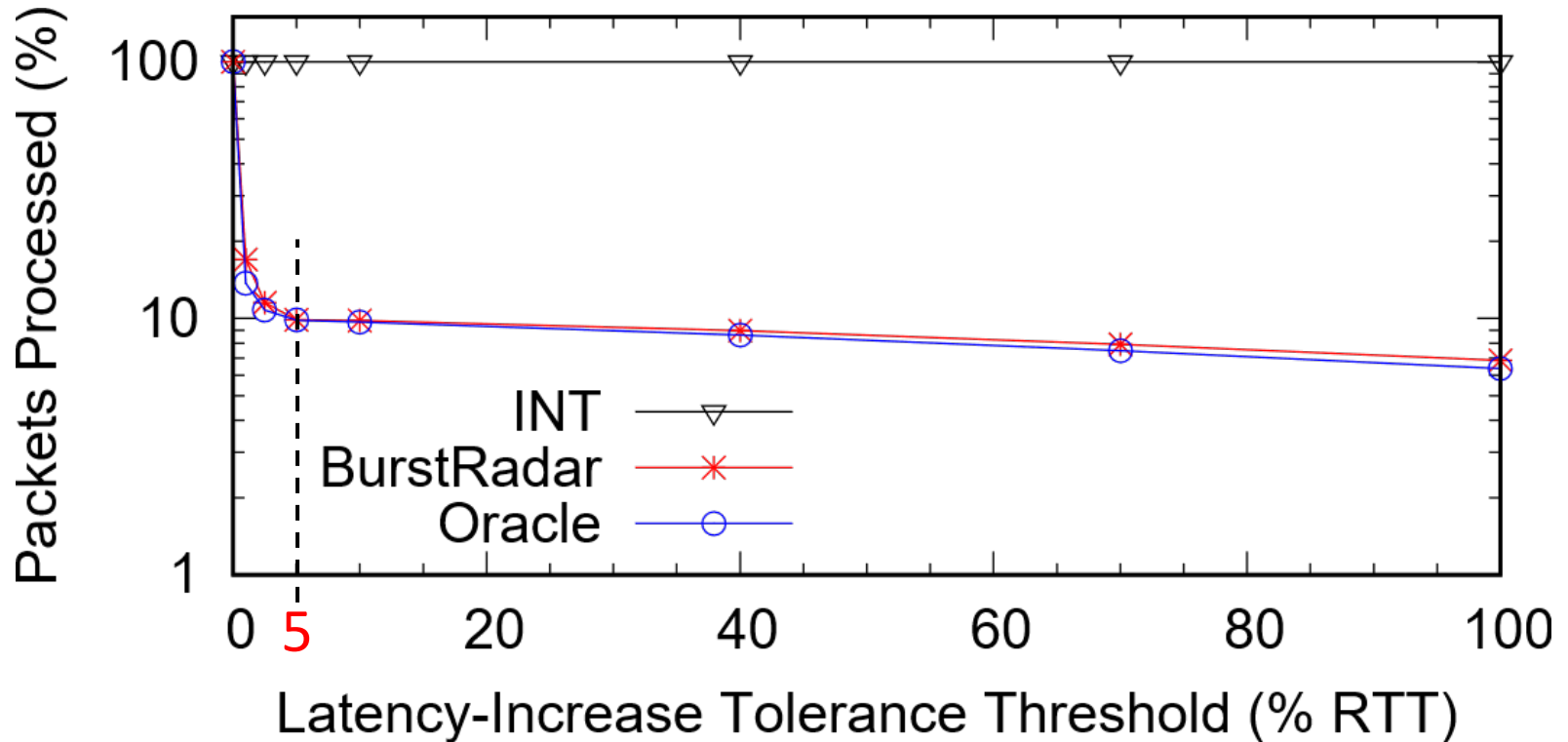


## Compare BurstRadar against

- In-band Telemetry (INT)  $\rightarrow$  dataplane-based solution
- “Oracle” Algorithm  $\rightarrow$  ground truth (exact pkts in  $\mu$ bursts)



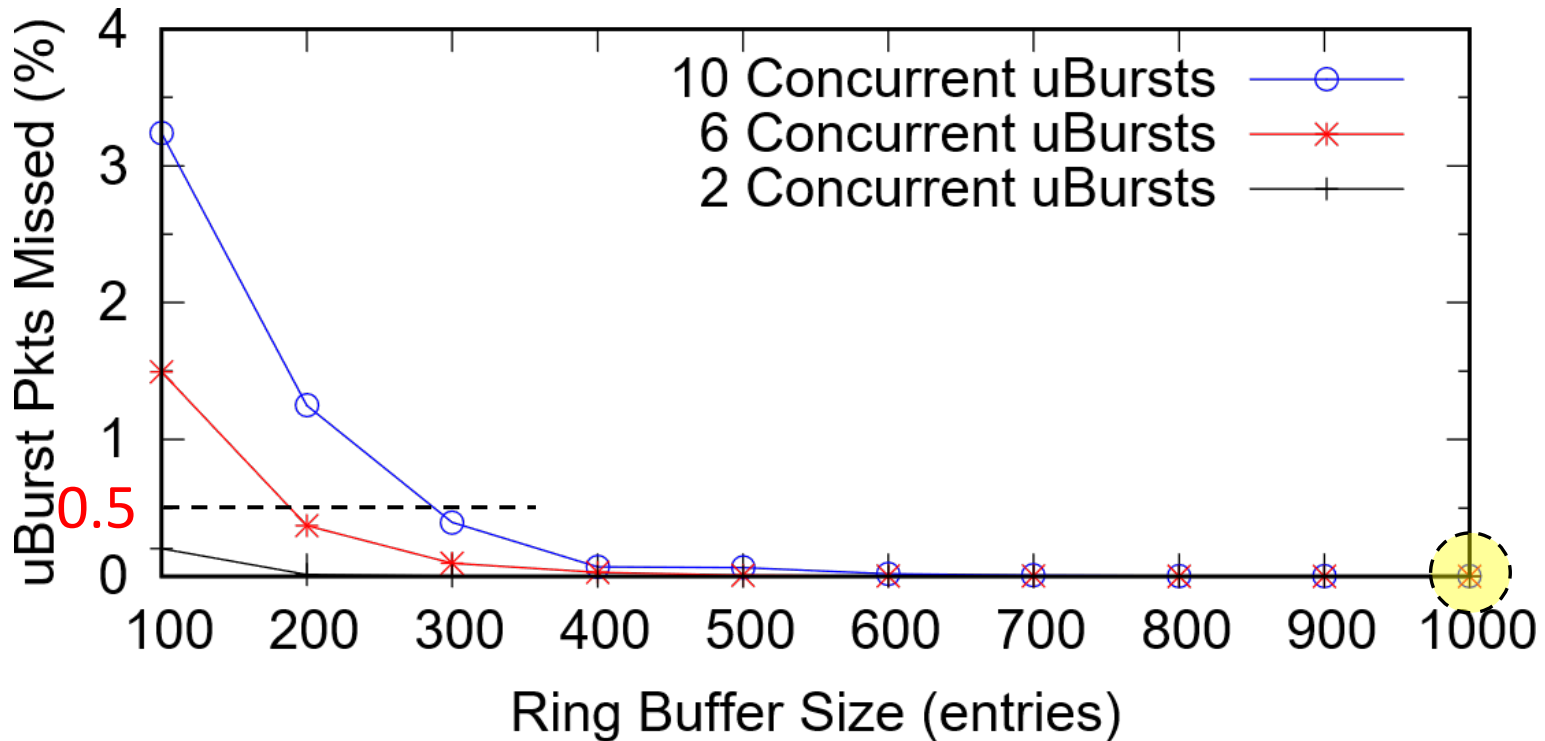
# Efficiency



5% RTT → 10 times less packets compared to INT

Note: 5% RTT  $\approx$  1.25 $\mu$ s of queuing @10Gbps in our testbed

# Handling Concurrent $\mu$ bursts



300 entries (8.7KB SRAM)  $\rightarrow$  10 concurrent  $\mu$ bursts ( $< 0.5\%$ )

**Note:** 1000 entries (29KB SRAM) fully handle 10 concurrent  $\mu$ bursts

# Dataplane Resource Utilization

Tofino Resource Utilization (Ring Buffer = 1000 entries)

Resource	switch.p4*	BurstRadar	Combined
Match Crossbar	50.13%	3.39%	53.52%
Hash Bits	32.35%	4.83%	37.18%
SRAM	29.79%	4.06%	33.85%
TCAM	28.47%	0.69%	29.16%
VLIW Actions	34.64%	4.69%	39.33%
Stateful ALUs	15.63%	12.5%	28.13%

\* resource utilization of a fully-featured datacenter ToR switch

Very little resources → combined with switch.p4

# Conclusion

---

- Microburst monitoring is important
  - High impact on performance
- BurstRadar can detect and identify Microbursts effectively and continuously
  - Capture and report the telemetry information of only the packets involved in microbursts
- BurstRadar demonstrates that modern programmable ASICs have made it practical to detect and characterize microbursts at multi-gigabit line rates in high-speed datacenter networks.

Vinaka  
 감사합니다  
 Dank Je  
 Blagodaram  
 Ngiyabonga  
 Dziukuje  
 Juspaxar  
 நன்றி  
 Ua Tsaug Rau Koj  
 Dėkuji  
 Suksama  
 Misaoira  
 Rahmat  
 Matur Nuwun  
 Bedankt  
 Dakujem  
 谢谢  
 Xbala  
 Welalin  
 Danke  
 Merci  
 Salamats  
 Maake  
 Kam Sah Hammida  
 Mauruuru  
 Grazas  
 Nirringrazzjak  
 Hvala  
 Di Ou Mesi  
 धन्यवाद  
 Asante  
 Shukria  
 Dhanyavadagalu  
 Manana  
 Dankon  
 Arigato  
 Gracias  
 cam ơn bạn  
 Kia Ora  
 Kop Khun Khap  
 Paldies  
 Obrigado  
 Djiere Dieuf  
 Tack  
 Taiku  
 Grazie  
 Mochchakkeram  
 Tingki  
 Gratias Tibi  
 ありがとう  
 Eskerrik Askos  
 Najis Tuke  
 Diolch i Chi  
 Terima Kasih  
 Matondo  
 Biyan  
 You  
 Thank