

# Scaling Hardware Accelerated Network Monitoring to Concurrent and Dynamic Queries with **\*Flow**

**John Sonchack**, Oliver Michel,  
Adam J. Aviv, Eric Keller, Jonathan M. Smith



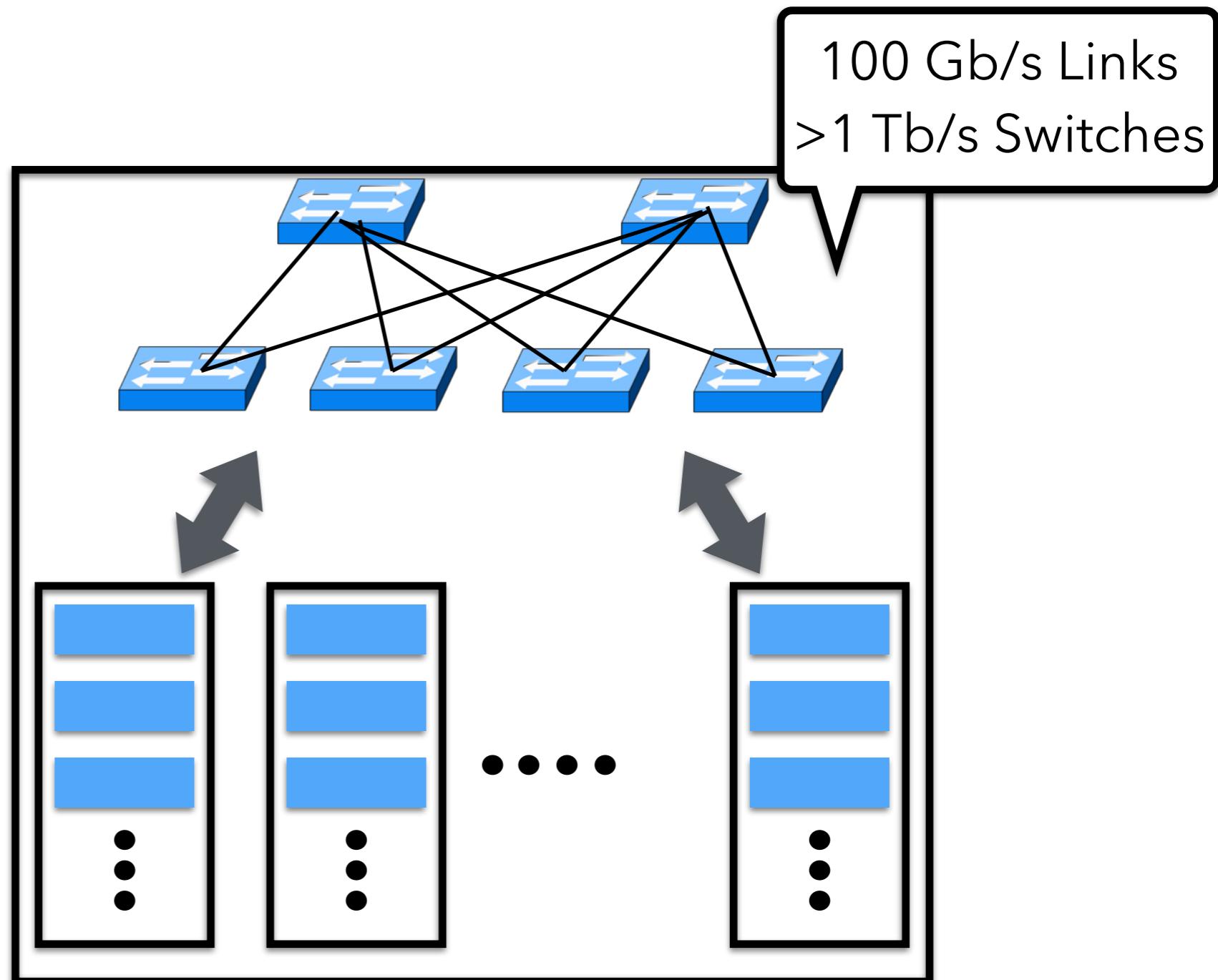
**Penn**  
UNIVERSITY of PENNSYLVANIA



University of Colorado  
Boulder



# Measuring High Speed Networks

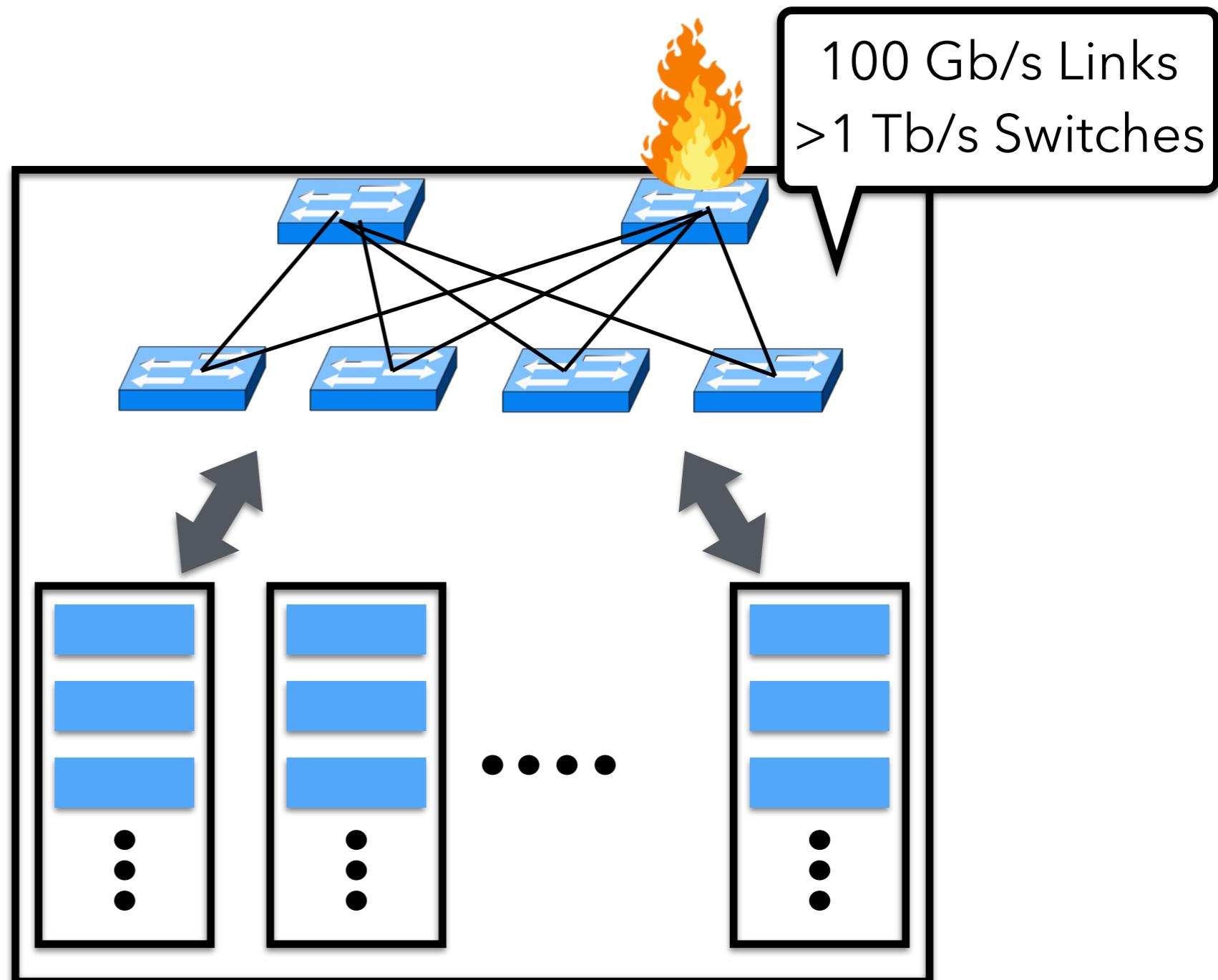


# Measuring High Speed Networks



*Who is causing microbursts?*

- Queue Lengths
- Drop Counts



# Measuring High Speed Networks



Debugging

*Who is causing microbursts?*

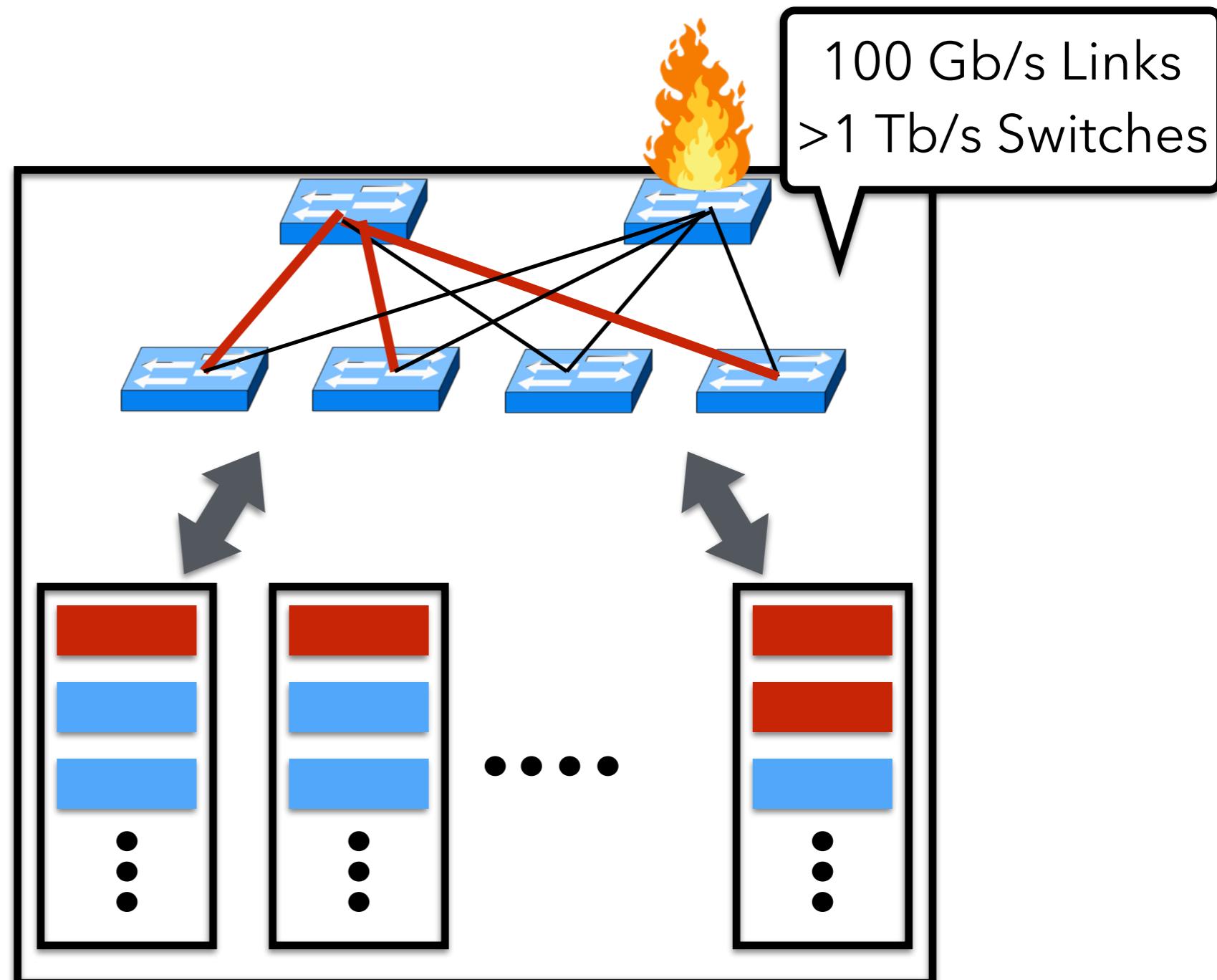
- Queue Lengths
- Drop Counts



Traffic Engineering

*Which flows are colliding?*

- Utilization



# Measuring High Speed Networks



Debugging

*Who is causing microbursts?*

- Queue Lengths
- Drop Counts



Traffic Engineering

*Which flows are colliding?*

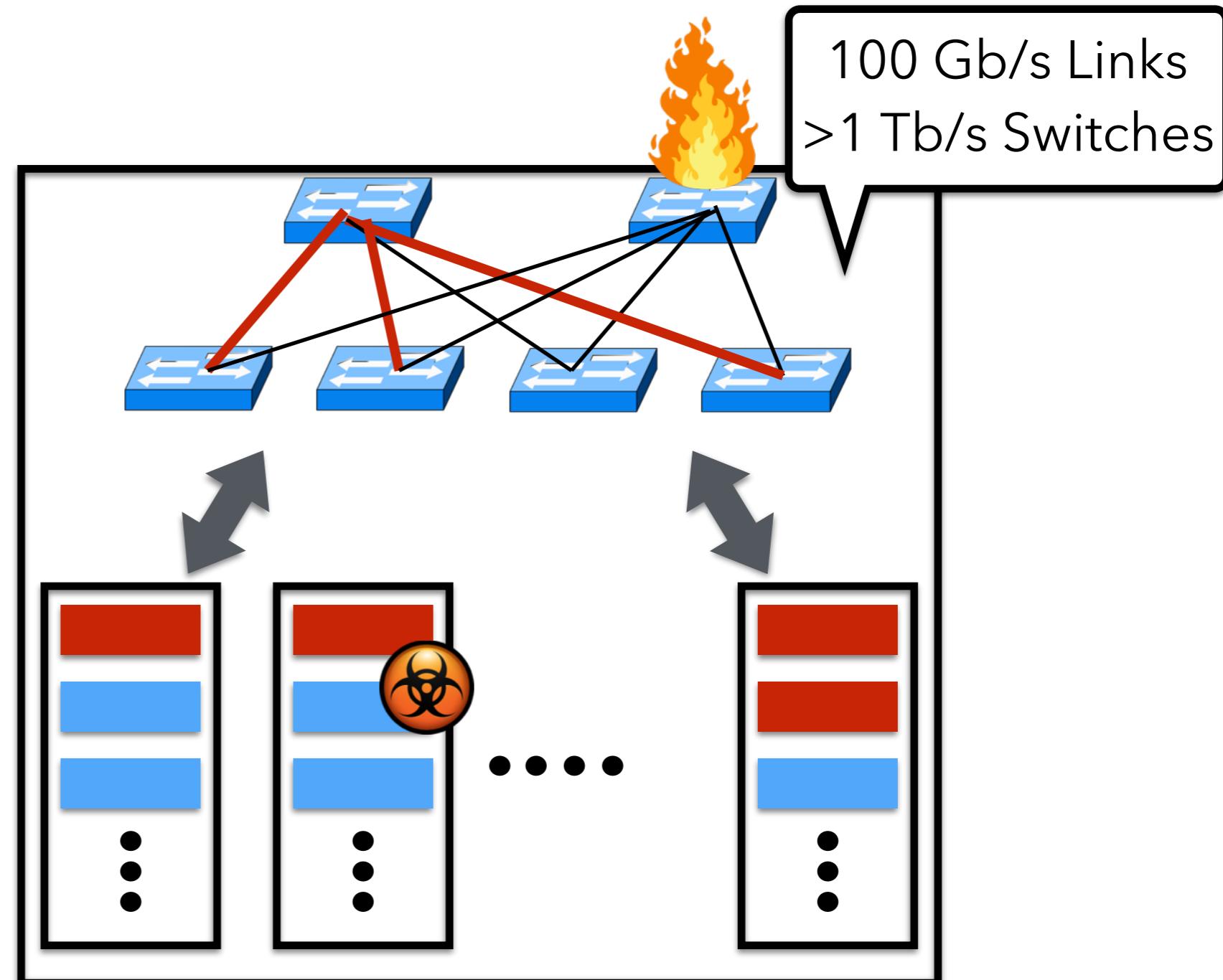
- Utilization



Security

*Are any hosts compromised?*

- Packet timing statistics



# Measurement Challenges

## Concurrent Applications      High Packet Rates



*Who is causing microbursts?*

- **Queue Lengths**
- **Drop Counts**



*Which flows are colliding?*

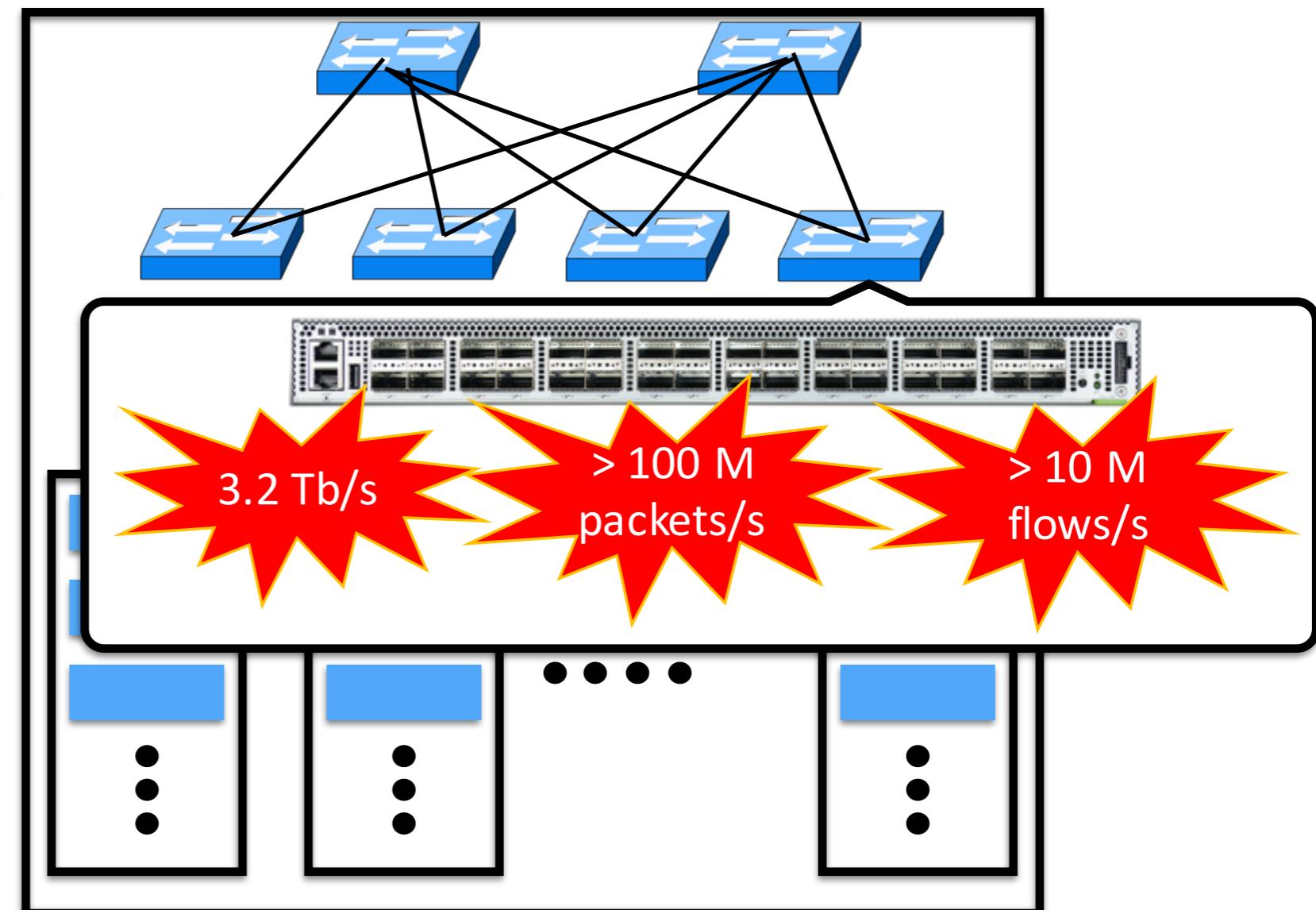
- **Utilization**

Traffic Engineering



*Are any hosts compromised?*

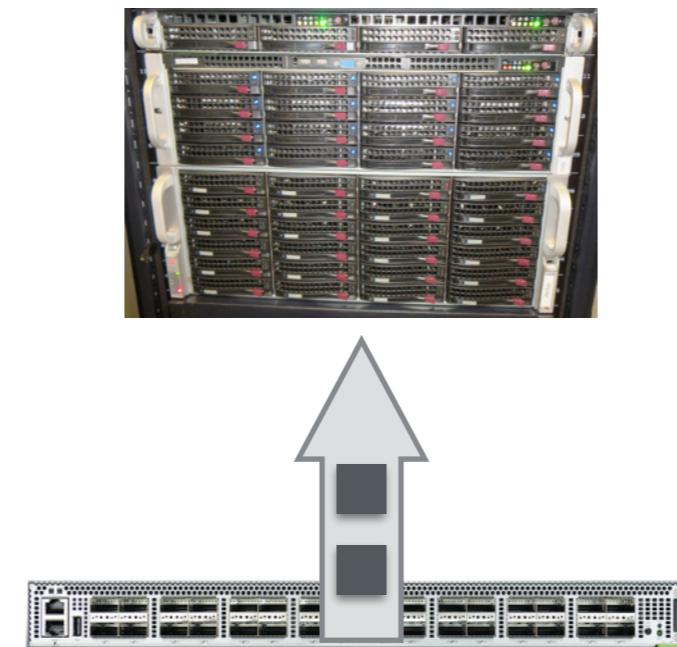
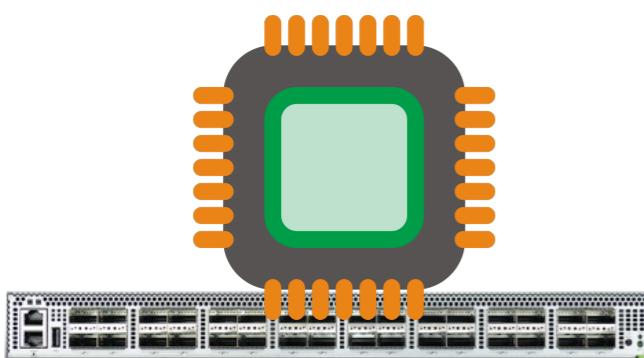
- **Packet timing statistics**



Security

# Prior Measurement Systems

Custom Hardware Processing  
(e.g., NetFlow)      Cluster

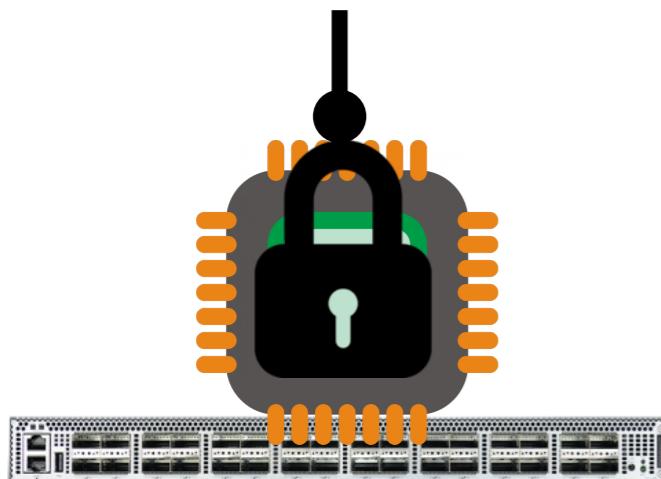


# Prior Measurement Systems

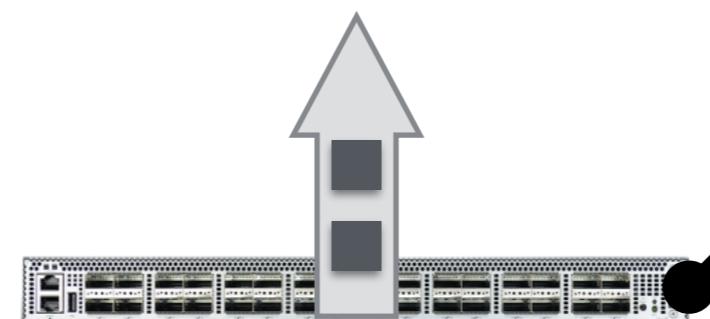
Custom Hardware Processing  
(e.g., NetFlow) Cluster

## Flexibility

*Applications can't define custom statistics*



## Efficiency



**Measurement Tput:**  
 $< 1 \text{ M}$  packets / second per core

**Network Tput:**  
 $> 100 \text{ M}$  packets / second per switch

# Prior Measurement Systems

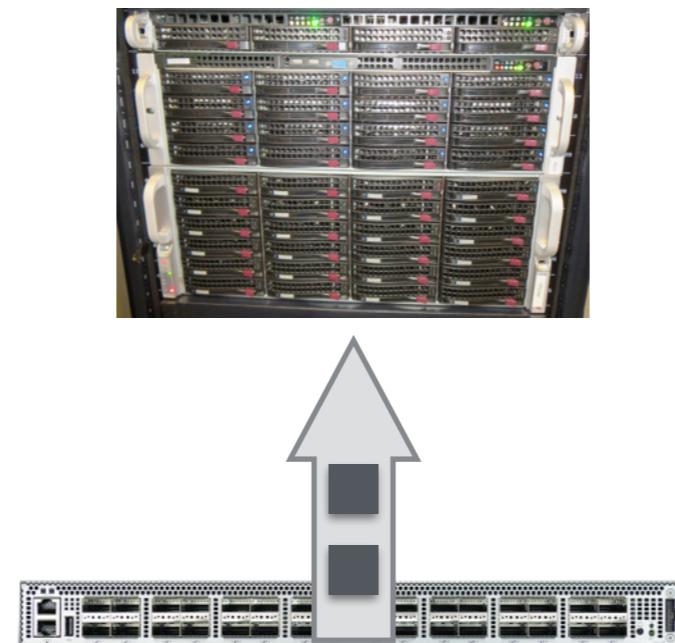
Custom Hardware  
(e.g., NetFlow)

Flexibility



Processing Cluster

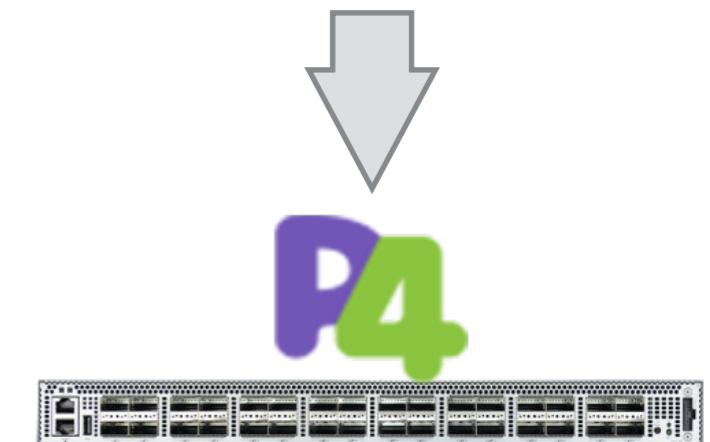
Efficiency



Reconfigurable  
Switch ASICs

[Marple SIGCOMM 17, TurboFlow  
EuroSys 18, Sonata SIGCOMM 18]

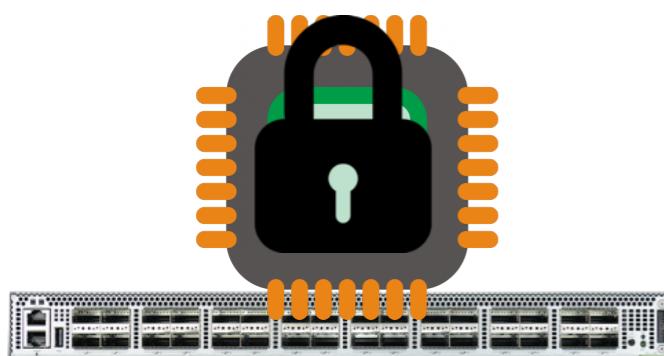
```
SELECT packet.length  
GROUP BY tcp flow  
AGGREGATE average
```



# Prior Measurement Systems

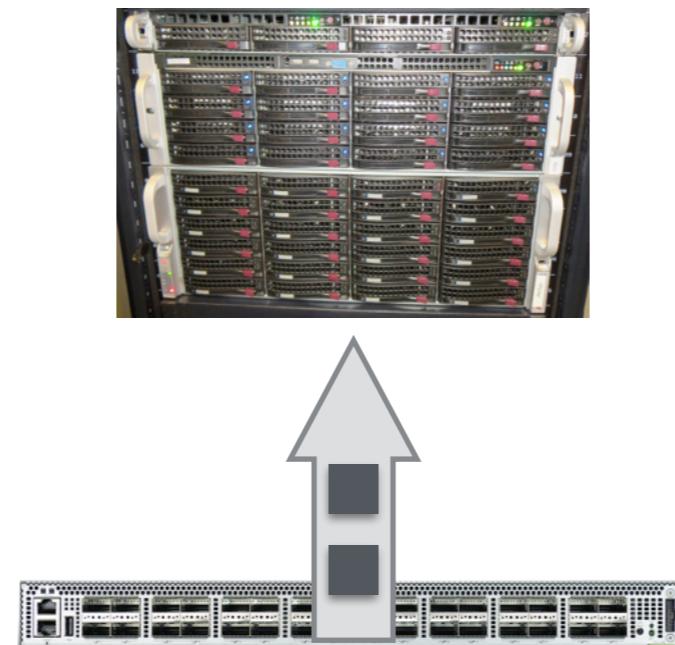
Custom Hardware  
(e.g., NetFlow)

Flexibility



Processing Cluster

Efficiency

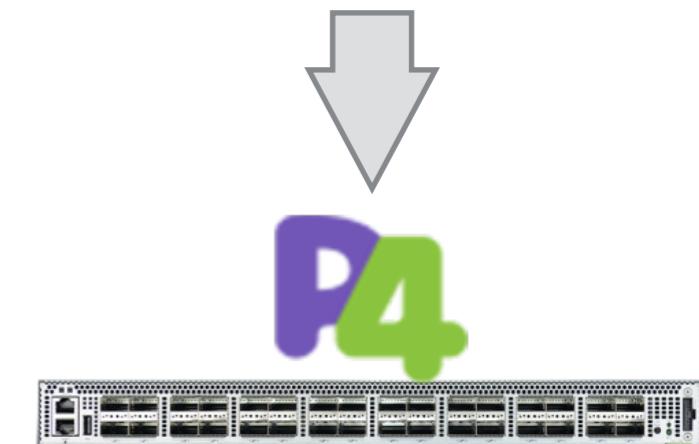


Reconfigurable Switch ASICs

[Marple SIGCOMM 17, TurboFlow EuroSys 18, Sonata SIGCOMM 18]

Concurrency

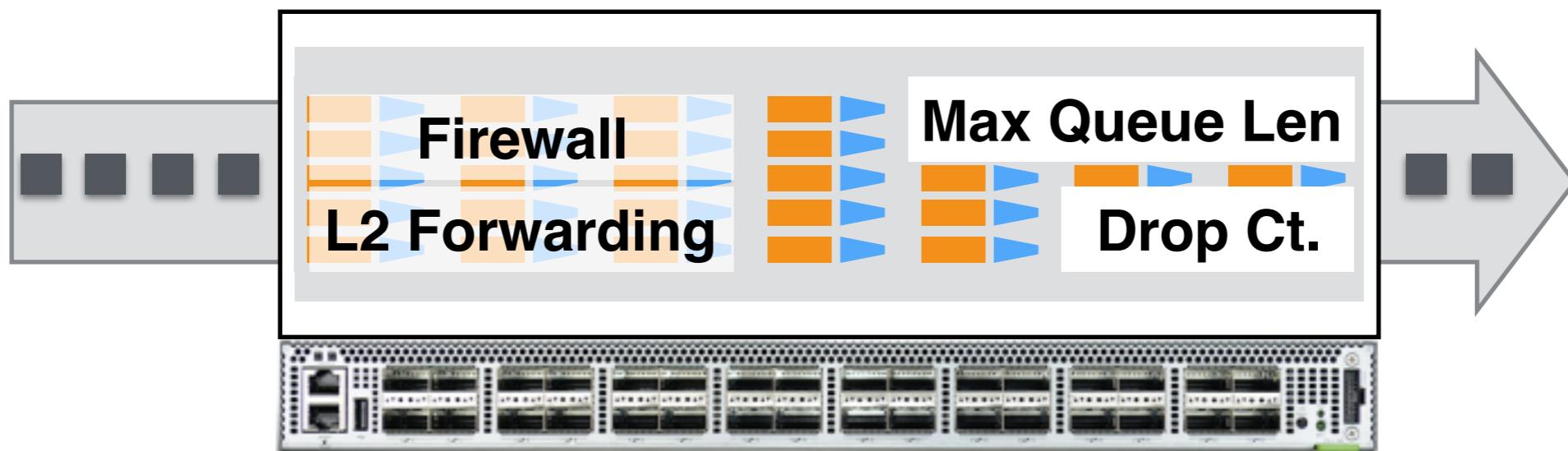
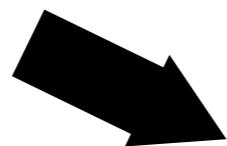
```
SELECT packet.length
GROUP BY tcp flow
AGGREGATE average
```



# Concurrency Challenges with Reconfigurable ASICs



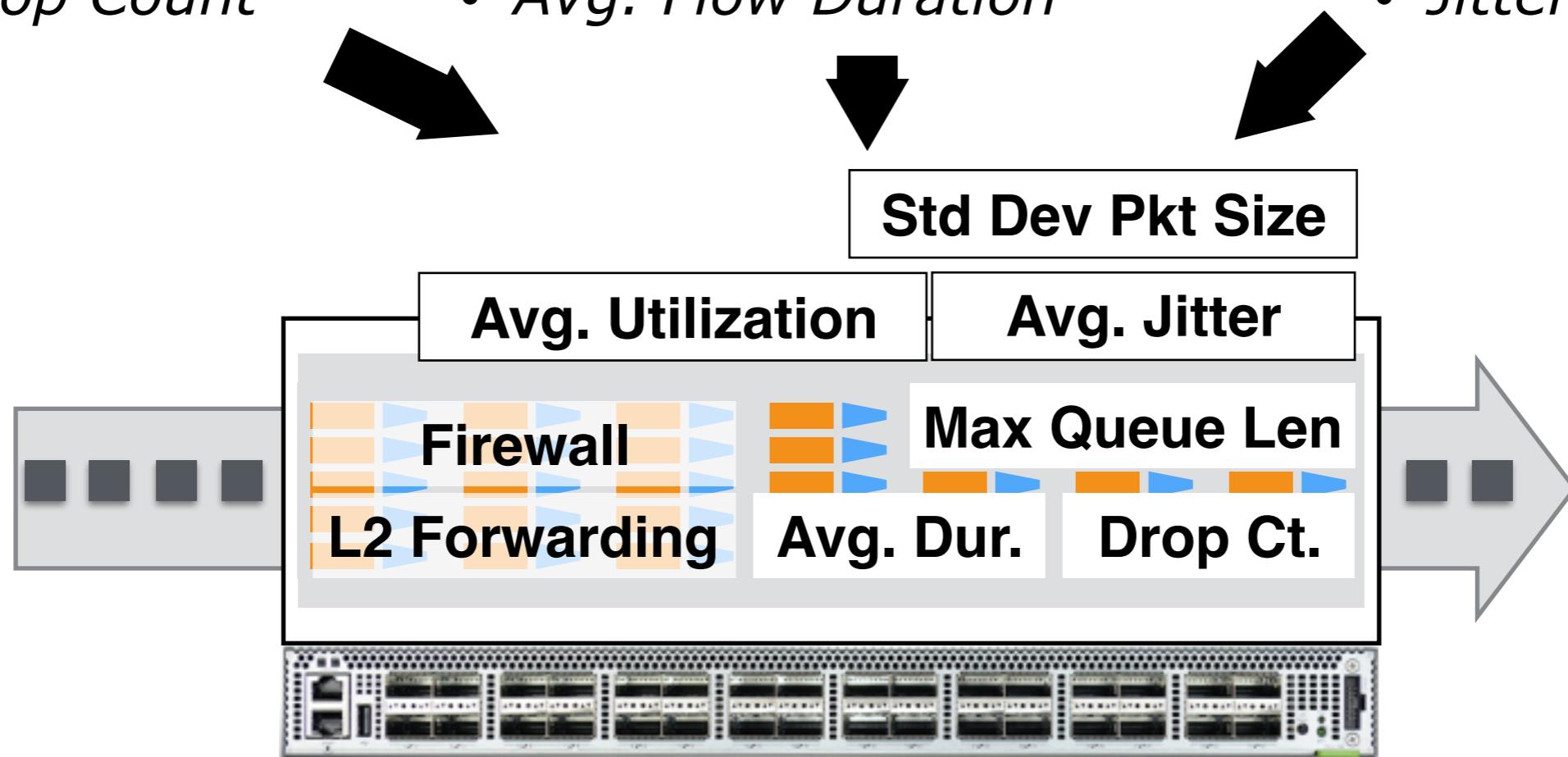
- *Max Queue Length*
- *Total Drop Count*



# Concurrency Challenges with Reconfigurable ASICs



- *Max Queue Length*
- *Total Drop Count*
- *Std Dev. Packet Size*
- *Avg. Flow Duration*
- *Avg. Utilization*
- *Jitter*



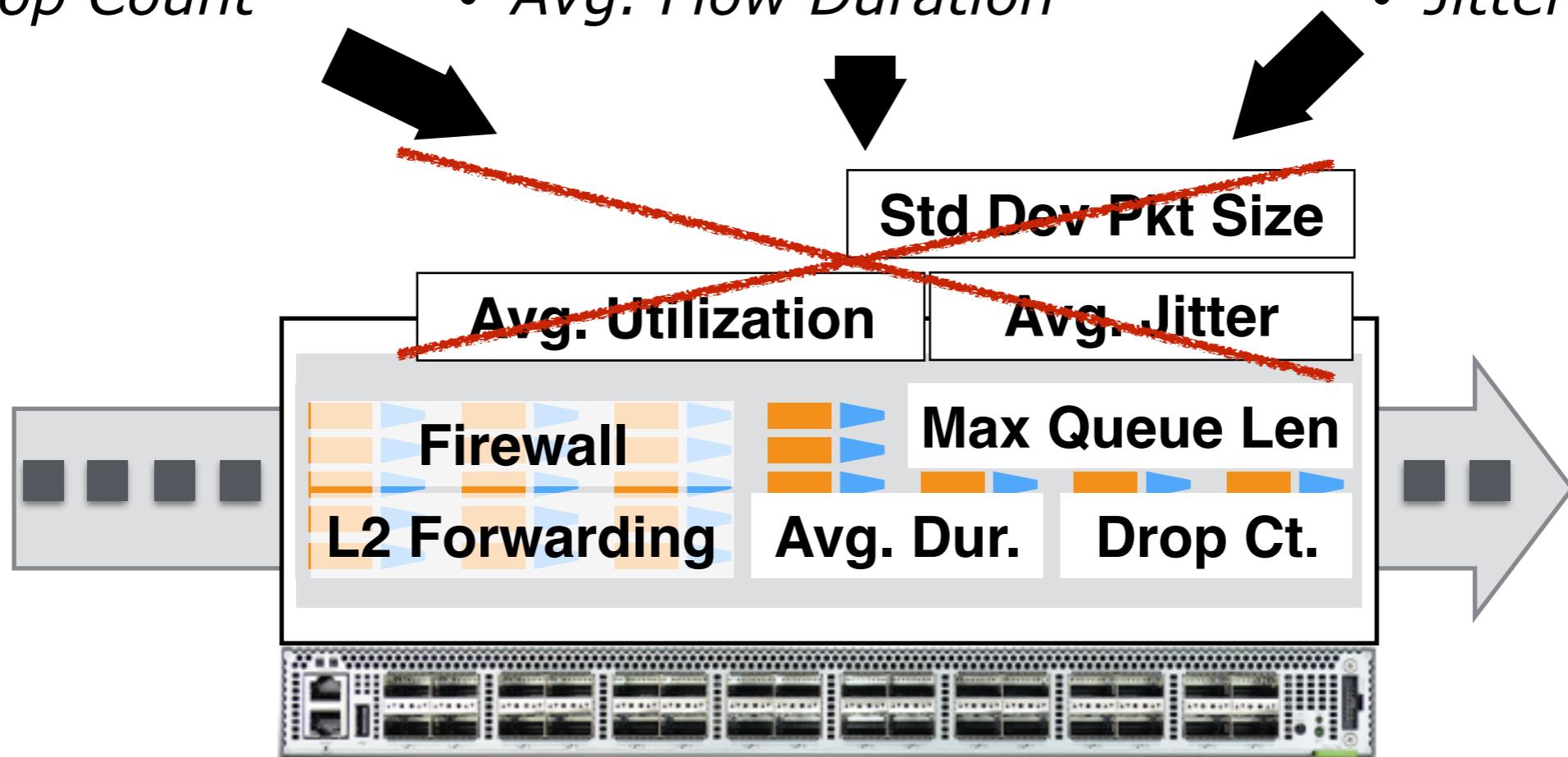
# Concurrency Challenges with Reconfigurable ASICs



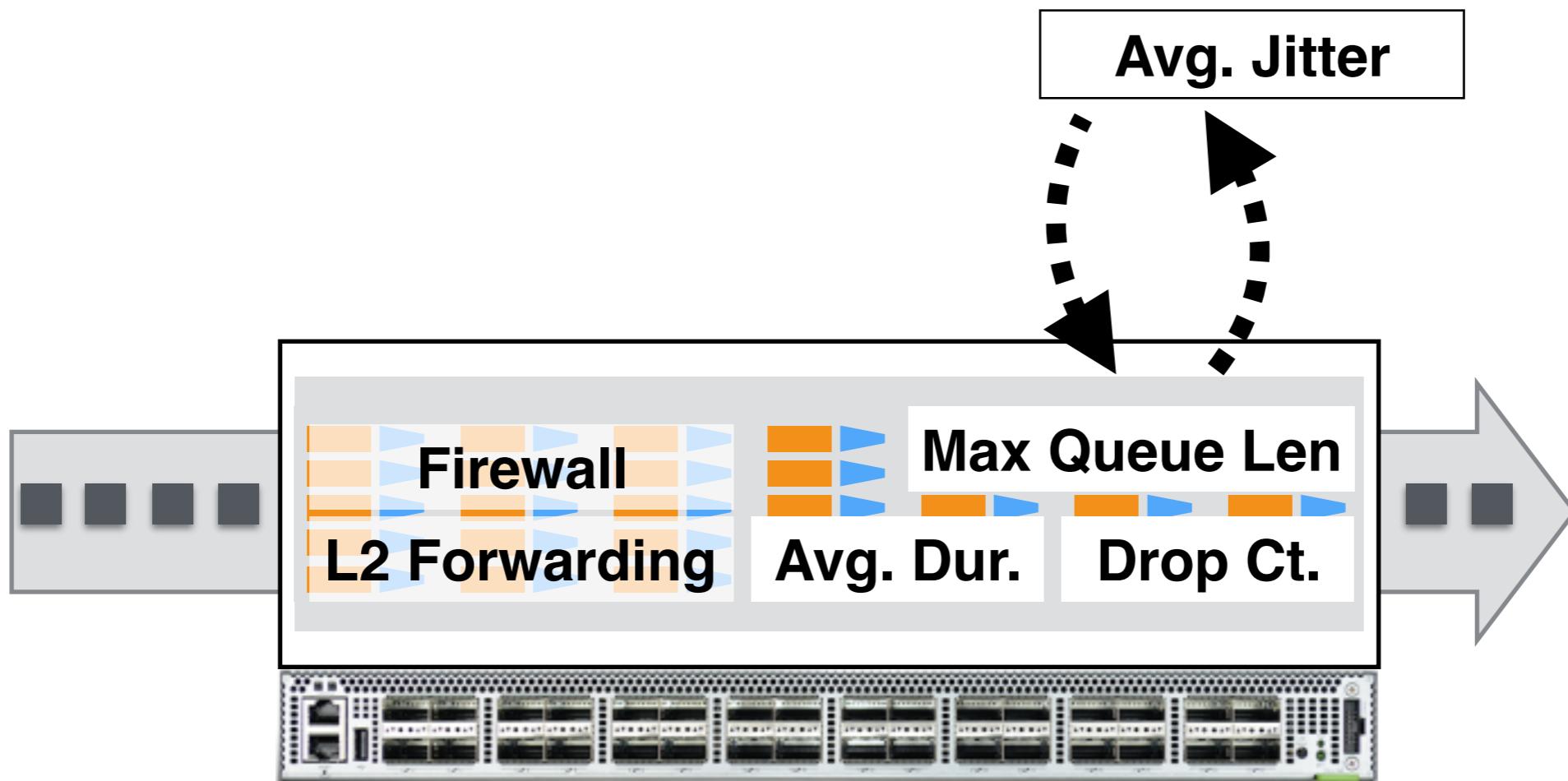
*Not Enough room in the datapath for all the queries.*



- Max Queue Length
- Total Drop Count
- Std Dev. Packet Size
- Avg. Flow Duration
- Avg. Utilization
- Jitter

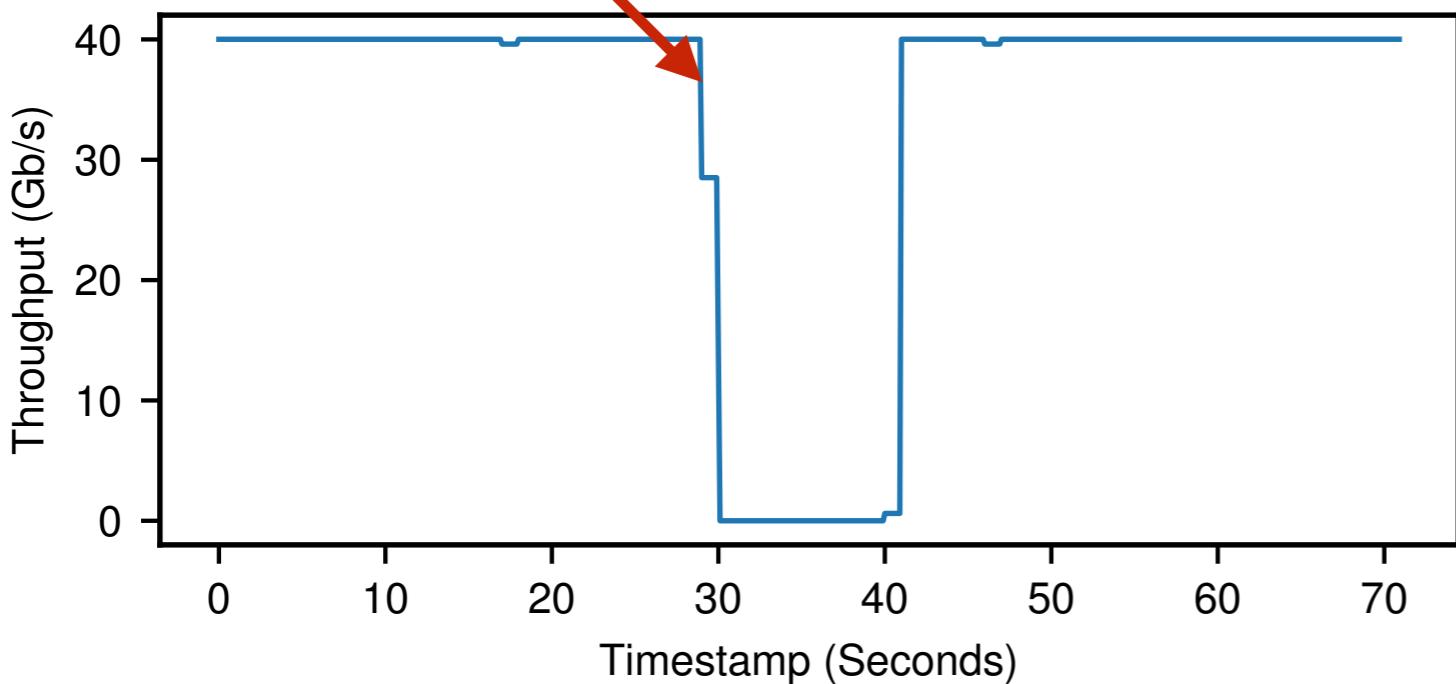


# Concurrency Challenges with Reconfigurable ASICs

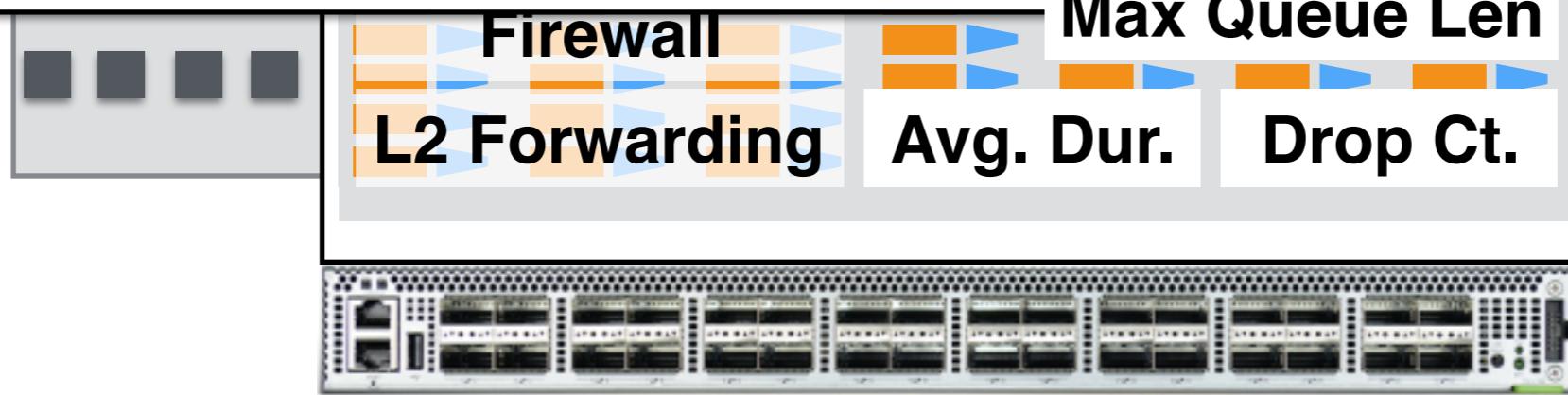
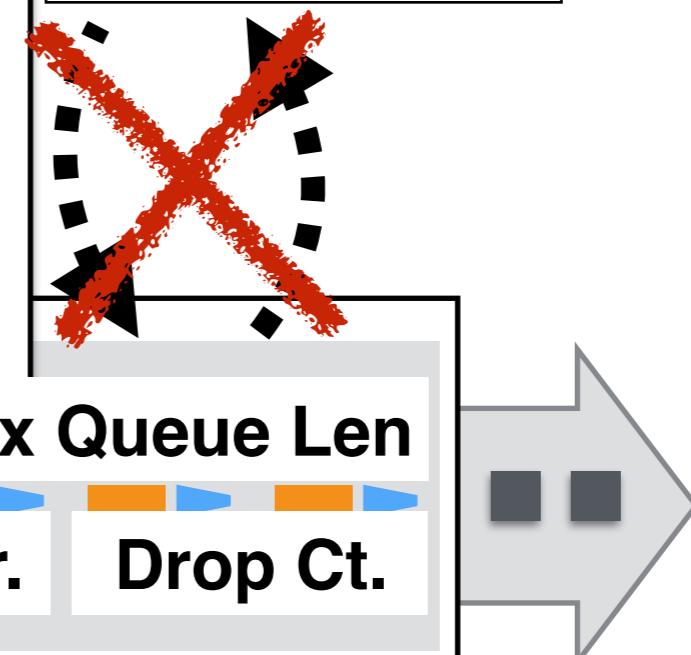


# Concurrency Challenges with Reconfigurable ASICs

Recompile P4 ASIC at 29 seconds



Recompiling the datapath disrupts measurement **and forwarding**.



# \*Flow: Efficiency and Concurrency?

**Question:** can we leverage reconfigurable ASICs for concurrent monitoring?

# \*Flow: Efficiency and Concurrency?

**Question:** can we leverage reconfigurable ASICs for concurrent monitoring?

**Insight:** concurrency challenges are caused by trying to do *too much in the ASIC*.

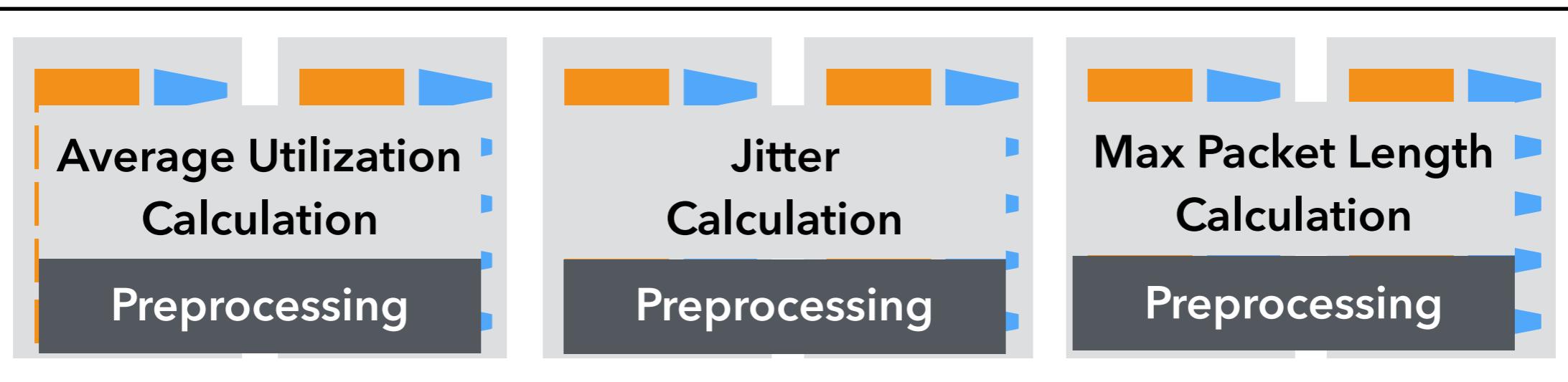
# The Main Idea

# The Main Idea

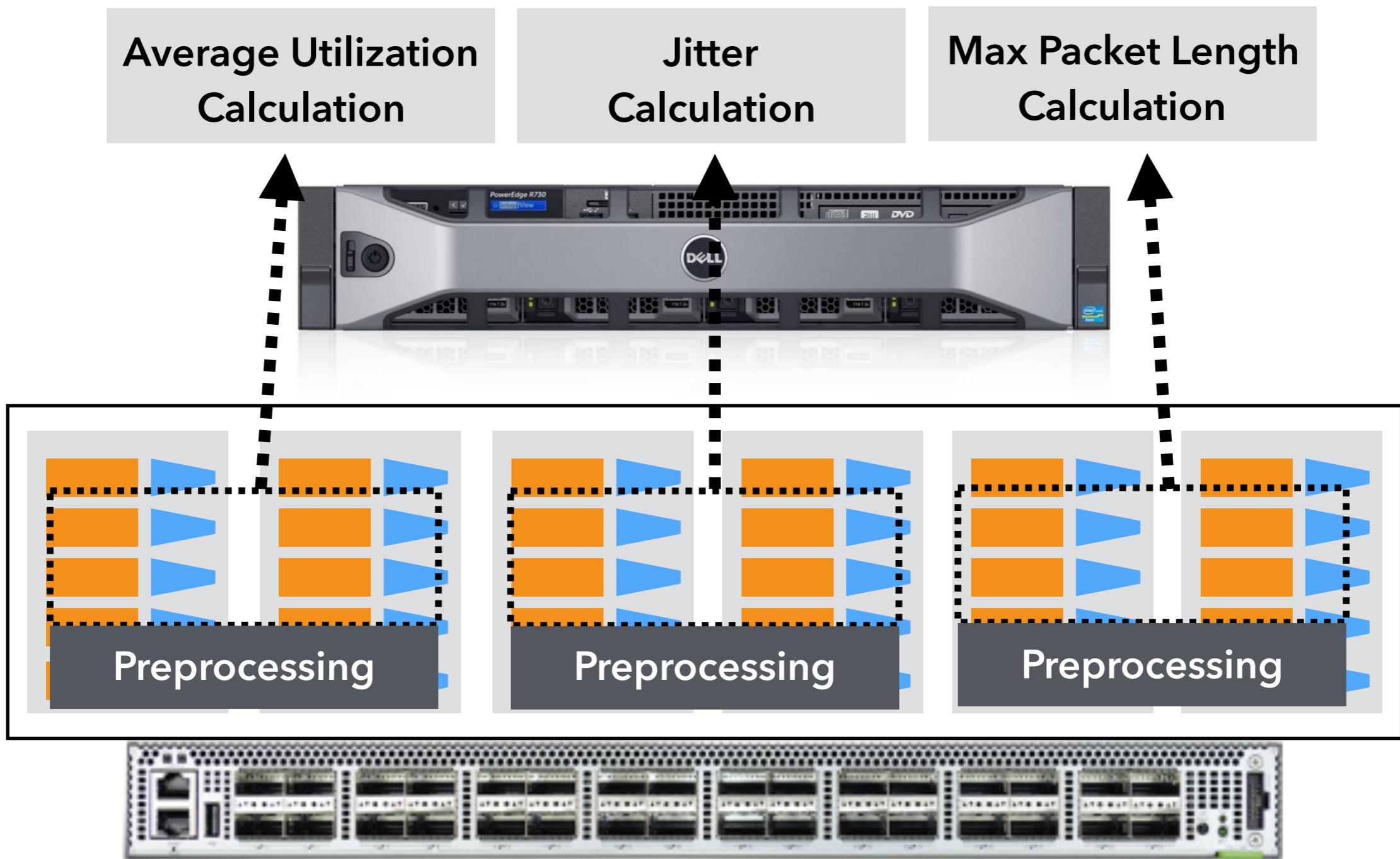
Preprocessing

```
SELECT packet.length  
GROUP BY tcp flow  
AGGREGATE max
```

Calculation



# The Main Idea

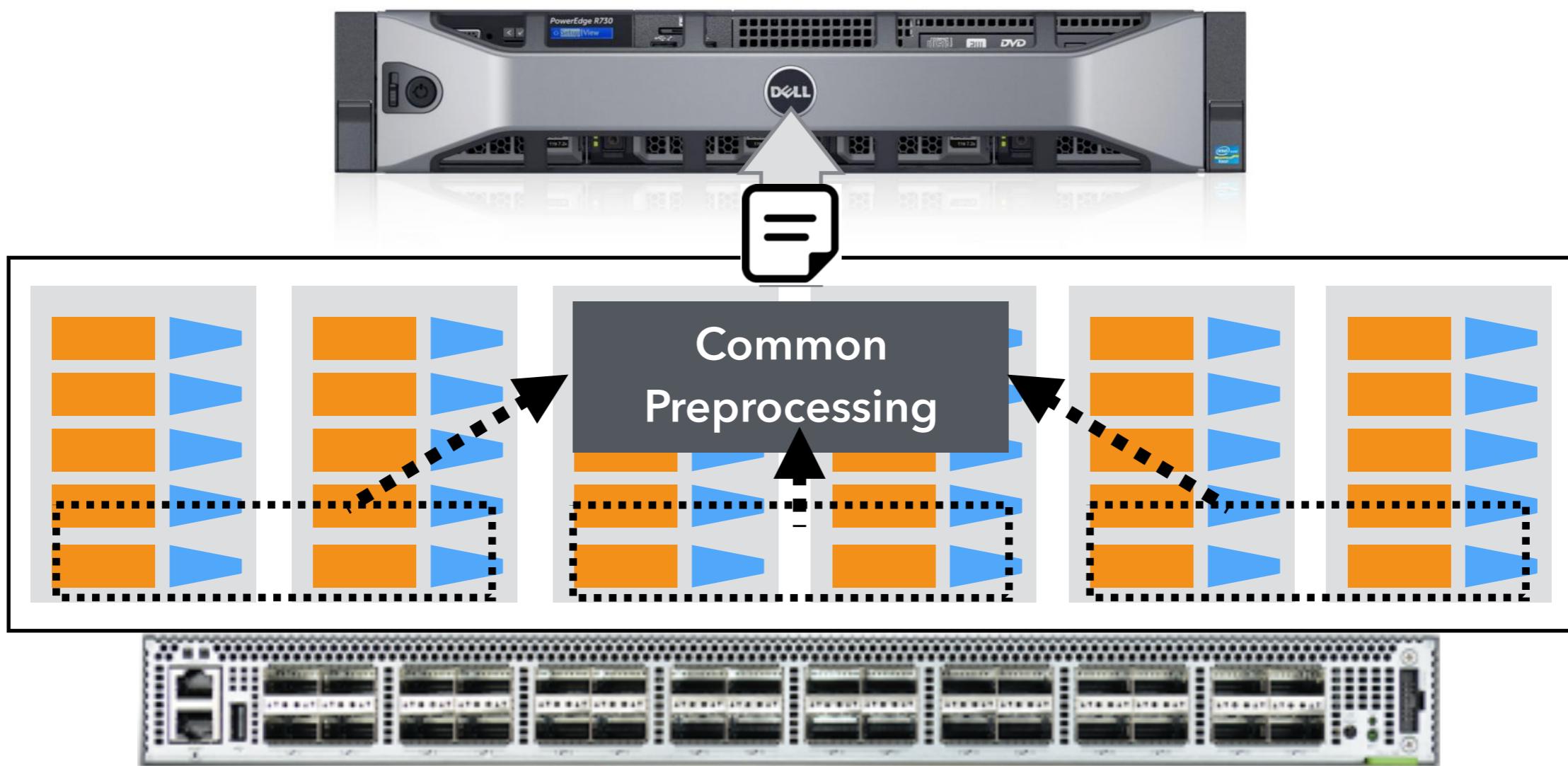


# The Main Idea

Average Utilization  
Calculation

Jitter  
Calculation

Max Packet Length  
Calculation



# The Main Idea



**Concurrency**

Average Utilization  
Calculation

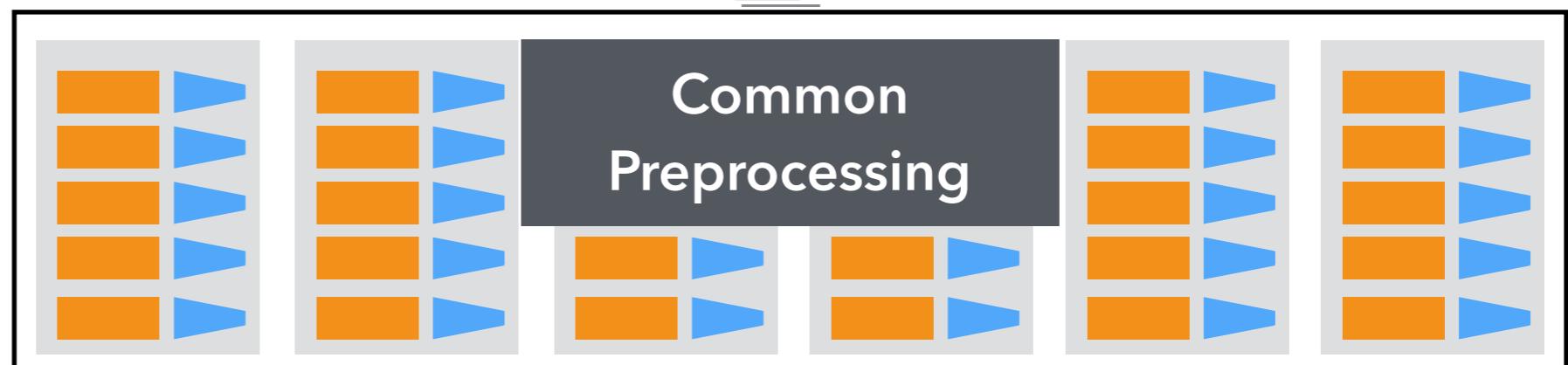
Max Packet Length  
Calculation

Jitter  
Calculation

...



Common  
Preprocessing



# The Main Idea



**Concurrency**

Average Utilization  
Calculation

Max Packet Length  
Calculation

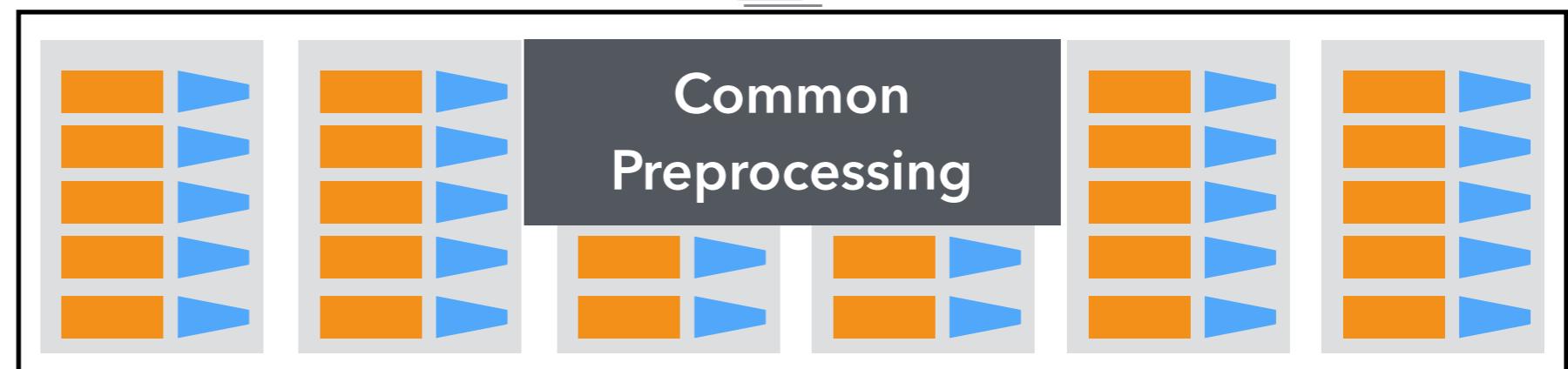
Jitter  
Calculation

...



**Efficiency**

Common  
Preprocessing



# The Main Idea



**Concurrency**

Average Utilization  
Calculation

Max Packet Length  
Calculation

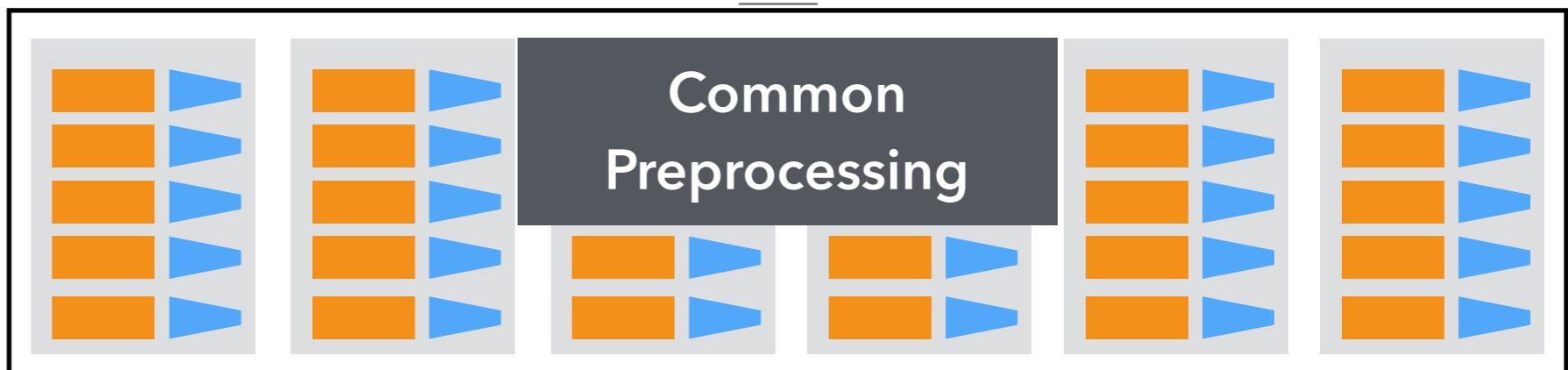
Jitter  
Calculation

...

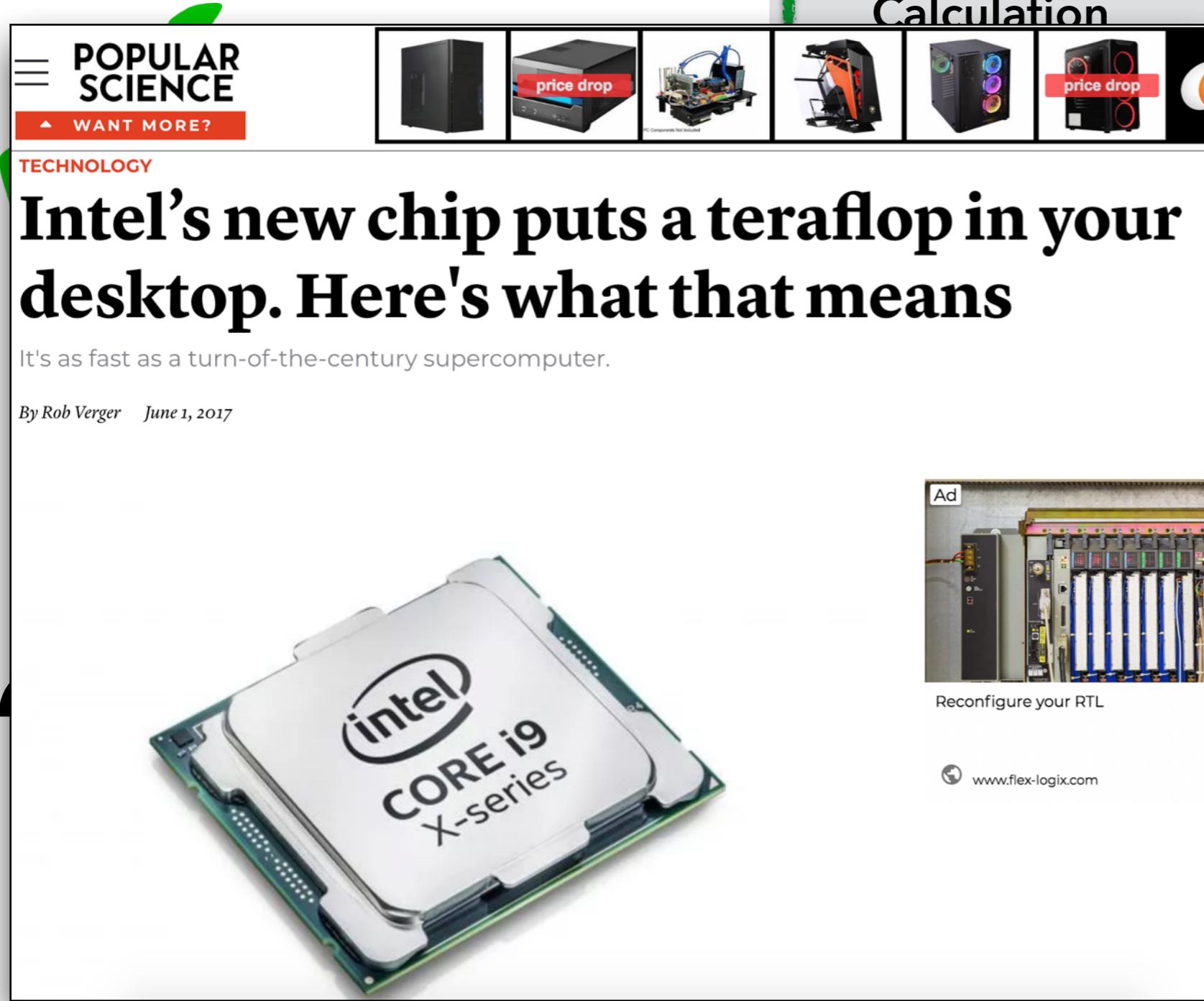


**Efficiency**

Common  
Preprocessing



# The Main Idea



A screenshot of a Popular Science article. The header reads "POPULAR SCIENCE" and "WANT MORE?". Below it is a "TECHNOLOGY" section. The main title is "Intel's new chip puts a teraflop in your desktop. Here's what that means". A subtext says "It's as fast as a turn-of-the-century supercomputer." The author is Rob Verger and the date is June 1, 2017. At the bottom is a large image of an Intel Core i9 X-Series processor.

Average Utilization  
Calculation

Max Packet Length  
Calculation

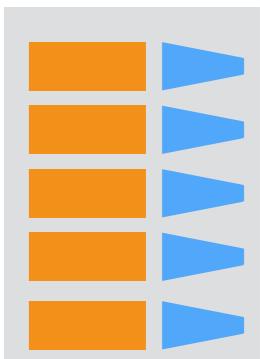
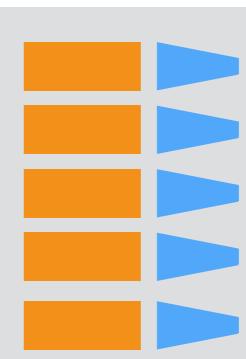
...



Reconfigure your RTL

[www.flex-logix.com](http://www.flex-logix.com)

Common  
Processing



# The Main Idea



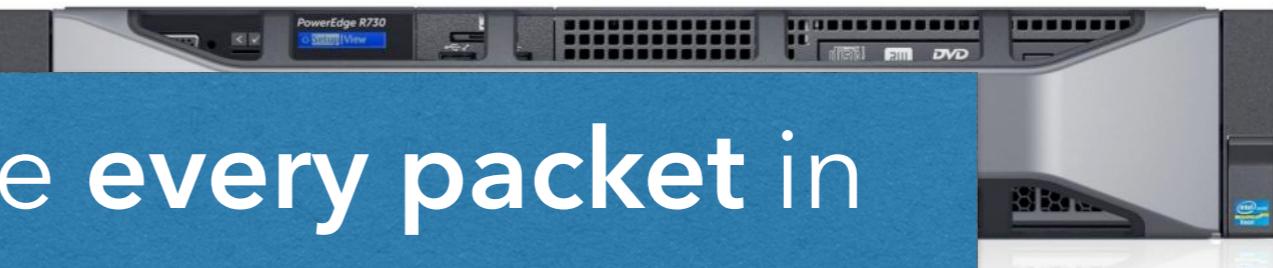
**Concurrency**

Average Utilization  
Calculation

Max Packet Length  
Calculation

Jitter  
Calculation

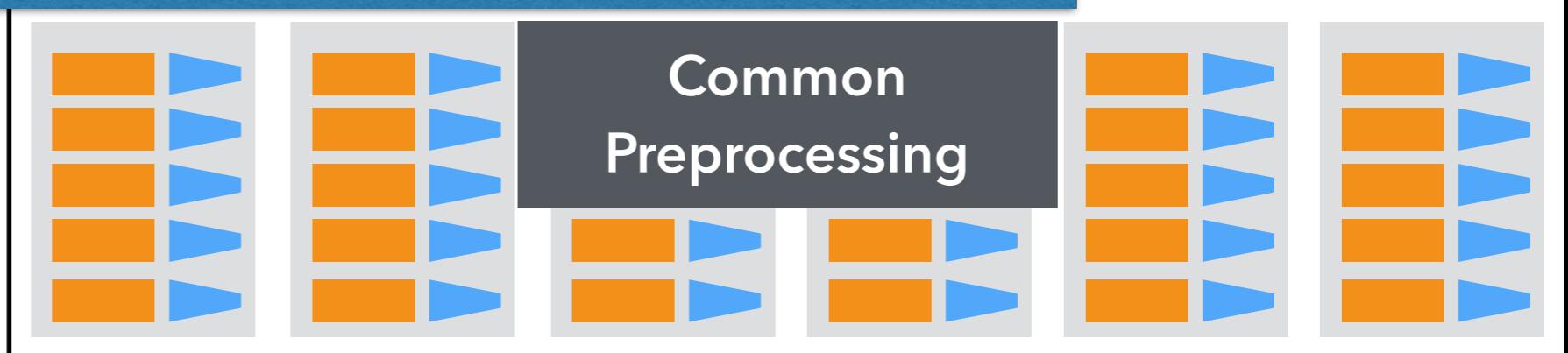
...



Measure **every packet** in  
 $> 1 \text{ Tb/s}$  traffic with 1 server.



**Efficiency**



# The Main Idea

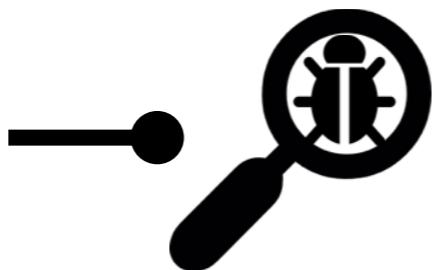


**Concurrency**



**Efficiency**

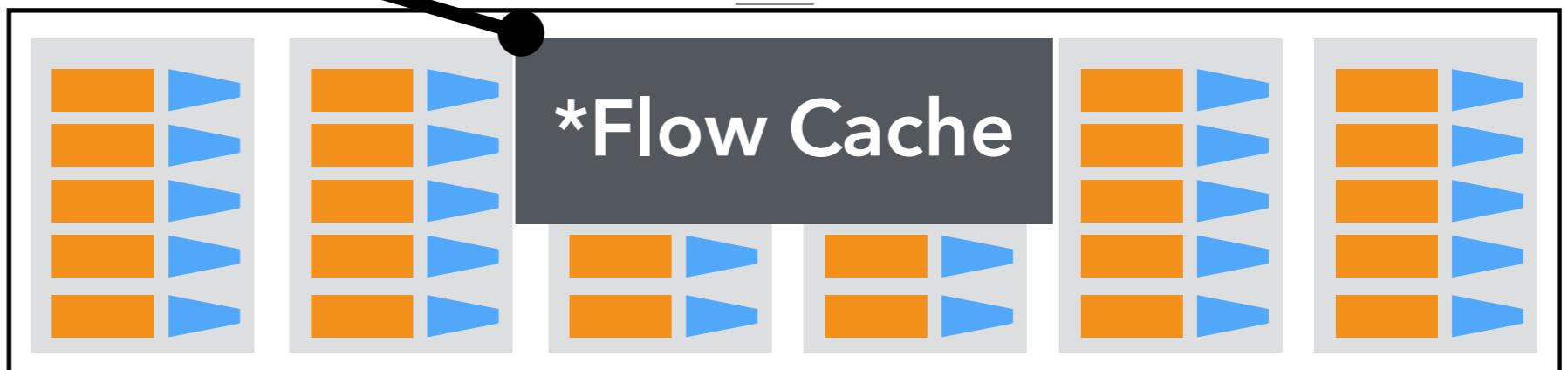
Application  
specific  
calculation



**\*Flow Agent**



Preprocessing



# Outline

- Motivation
- **Design**
- Implementation
- Evaluation

# \*Flow Design

```
SELECT ip length  
GROUP BY tcp flow  
AGGREGATE sum
```

1. Decoupling calculation
2. Generalize selection & grouping

Common Preprocessing



App. Specific Calculation



# Decoupling Calculation

```
SELECT ip length  
GROUP BY tcp flow  
AGGREGATE sum
```

| <i>Flow Key Length Sum</i> |      |
|----------------------------|------|
| A -> B                     | 1000 |
| C -> D                     | 250  |
| E -> F                     | 9416 |
| ...                        | ...  |

**Key:** A -> B  
**Length:** 200

Preprocessing



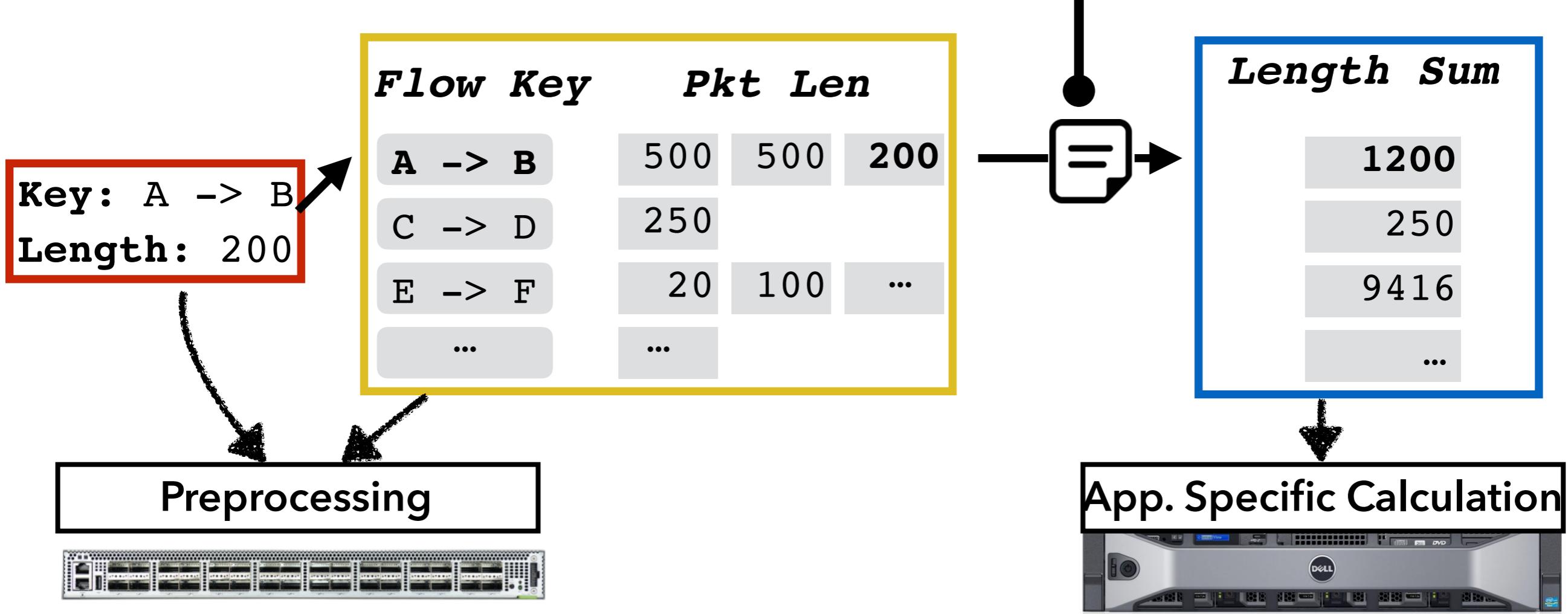
App. Specific Calculation



# Decoupling Calculation

```
SELECT ip length  
GROUP BY tcp flow  
AGGREGATE sum
```

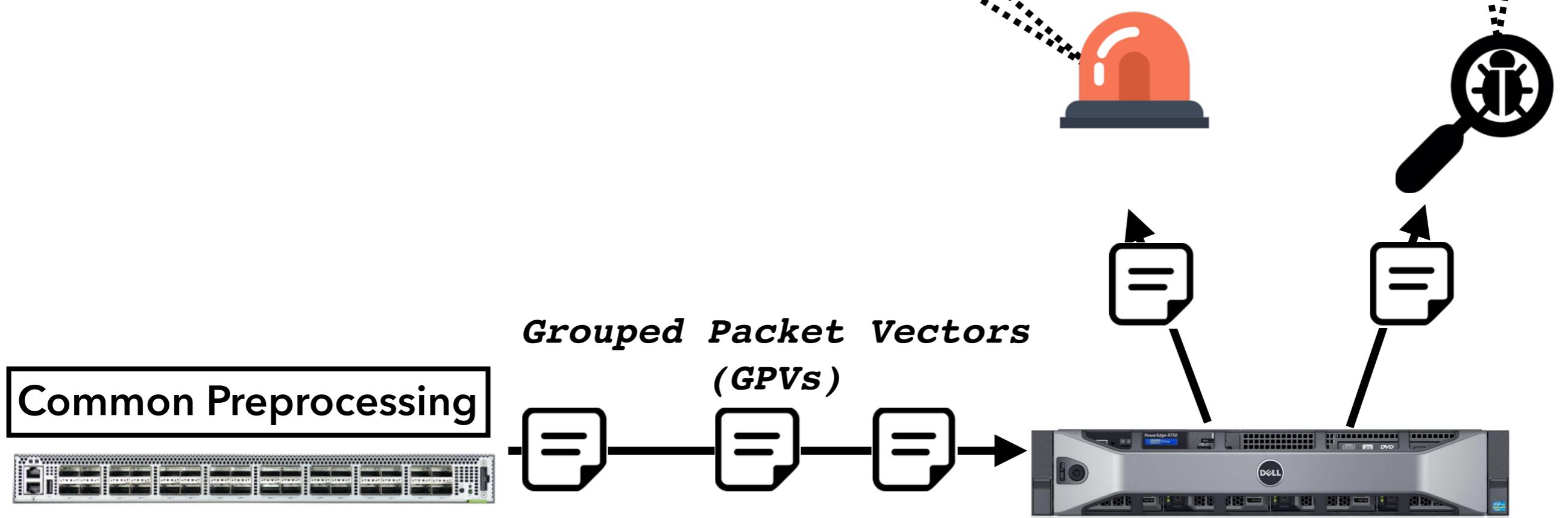
*Grouped Packet Vector (GPV):*  
 $\{flowKey : packetLength\}$



# Generalizing Selection and Grouping

```
SELECT ip length  
GROUP BY host  
AGGREGATE average
```

```
SELECT timestamp  
GROUP BY tcp flow  
AGGREGATE variance
```



# Generalizing Selection and Grouping

```
SELECT ip length  
GROUP BY host  
AGGREGATE average
```

```
SELECT timestamp  
GROUP BY tcp flow  
AGGREGATE variance
```

$< \text{flowKey} : (\text{packetLengths}, \text{packetTimestamps}, \dots) >$

Grouped Packet Vectors  
(GPVs)

Common Preprocessing



# Generalizing Selection and Grouping

```
SELECT ip.length  
GROUP BY host  
AGGREGATE average
```

```
SELECT timestamp  
GROUP BY tcp.flow  
AGGREGATE variance
```

Regroup to any sub-key for 1  
KV op. and copy per GPV.

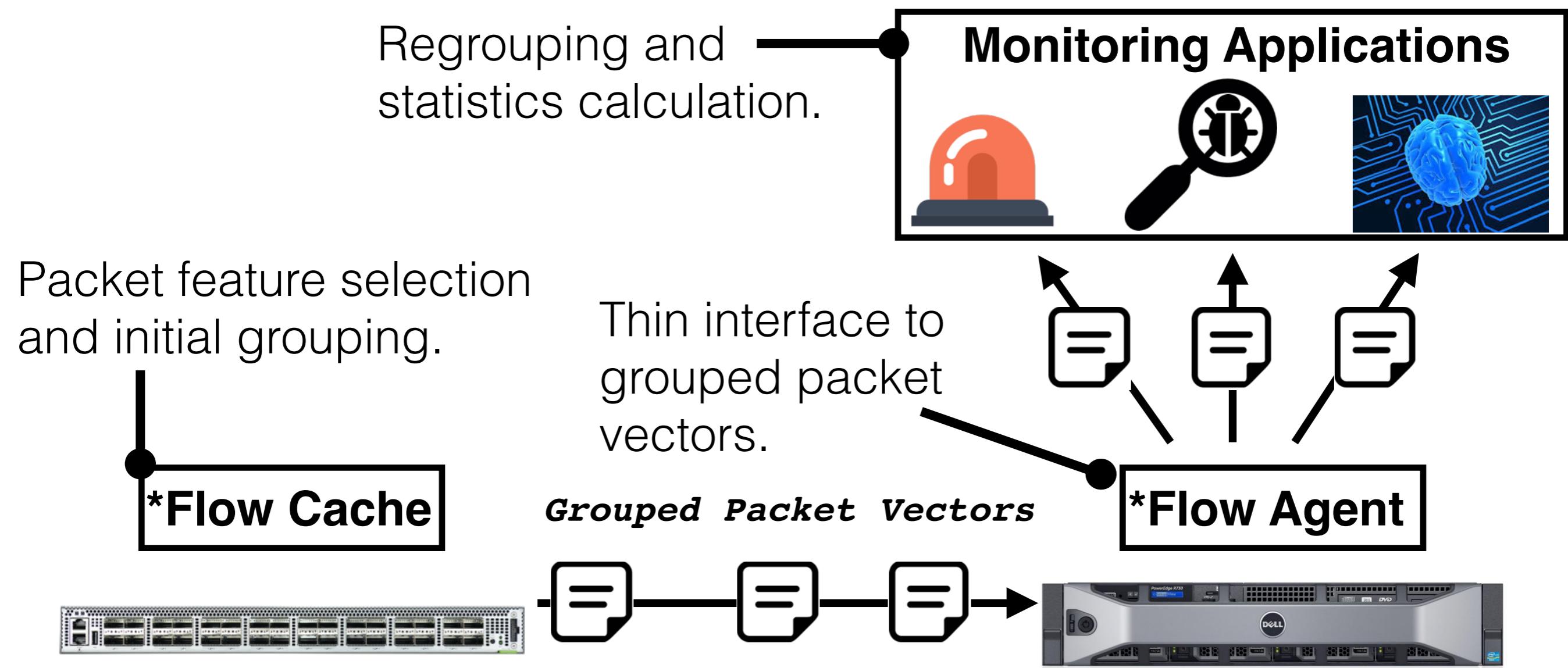
export key = <src IP, dst IP,  
src port, dst port, protocol,  
link ID>

*Grouped Packet Vectors  
(GPVs)*

Common Preprocessing



# \*Flow Architecture



# Outline

- Introduction
- Design
- **Implementation**
- Evaluation

# \*Flow Implementation

Proof-of-concept implementation.

P4 implementation for 3.2 Tb/s Barefoot Tofino.

**\*Flow Cache**

*Grouped Packet Vectors*



**\*Flow Agent**

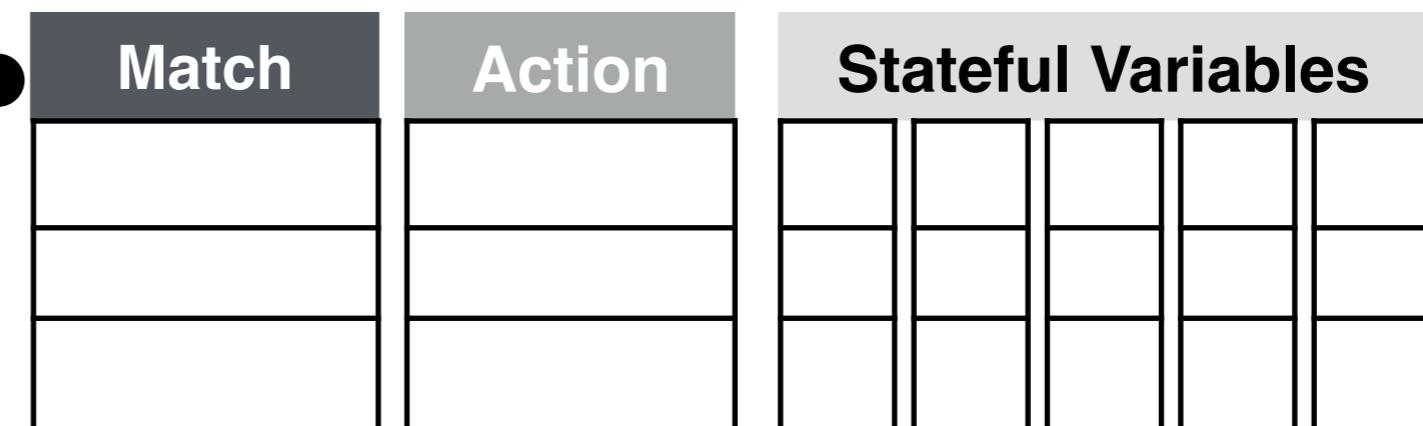


# Background: Stateful Match Action ASICs

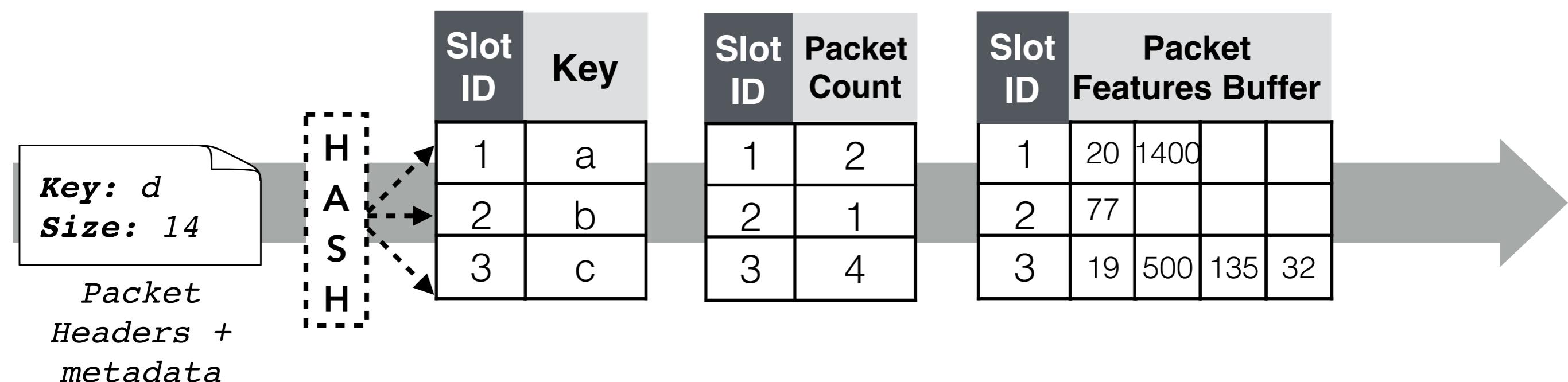
Restrictive  
computational  
model

Line rate  
processing  
(1 packet per cycle)

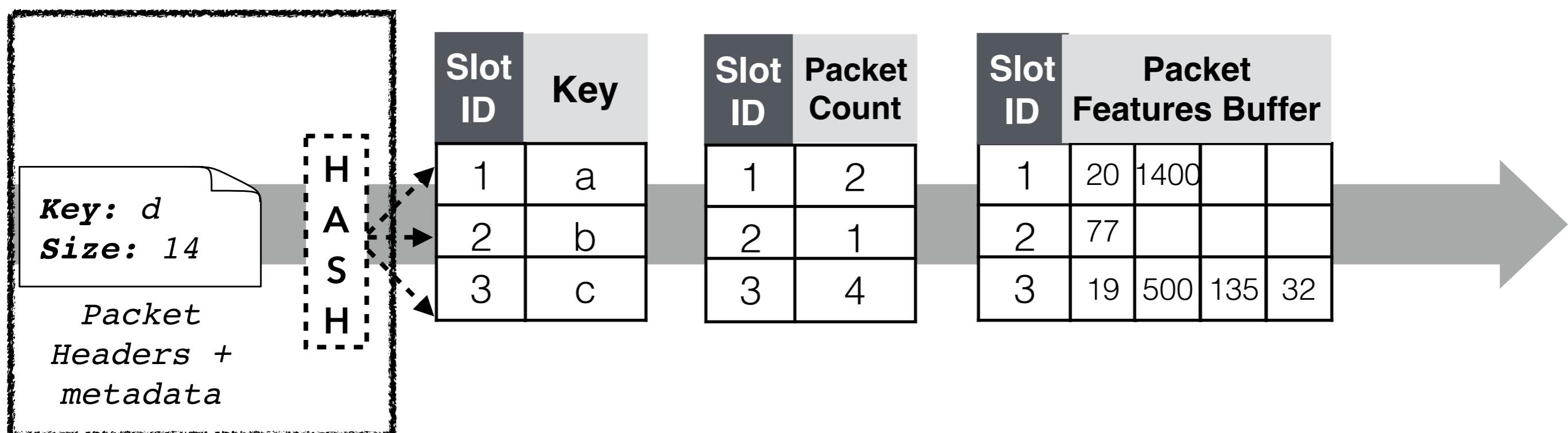
*Packet  
Headers +  
metadata*



# \*Flow as a Stateful Match Action Pipeline



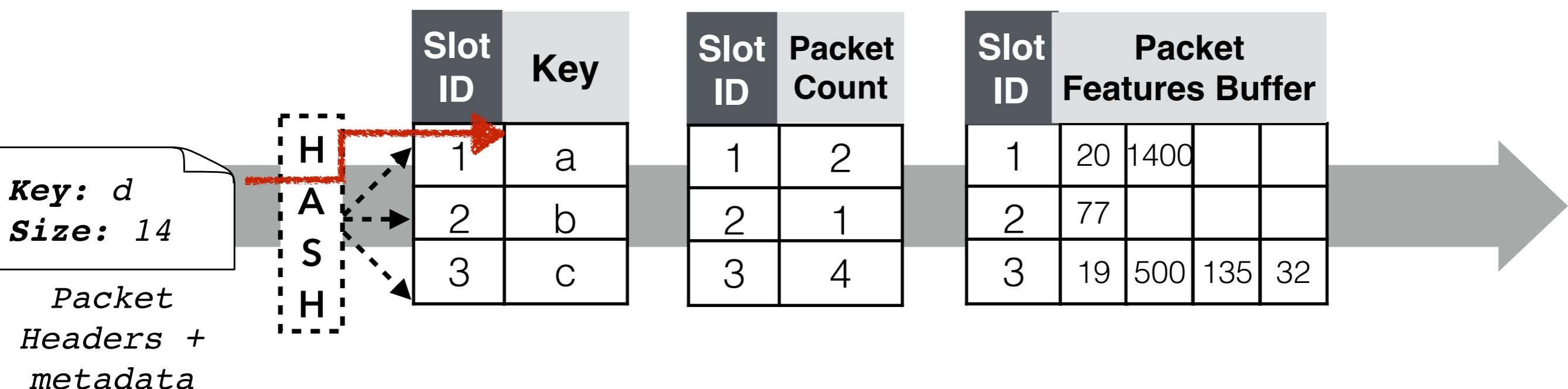
# \*Flow as a Stateful Match Action Pipeline



# \*Flow as a Stateful Match Action Pipeline

## Untracked Flow:

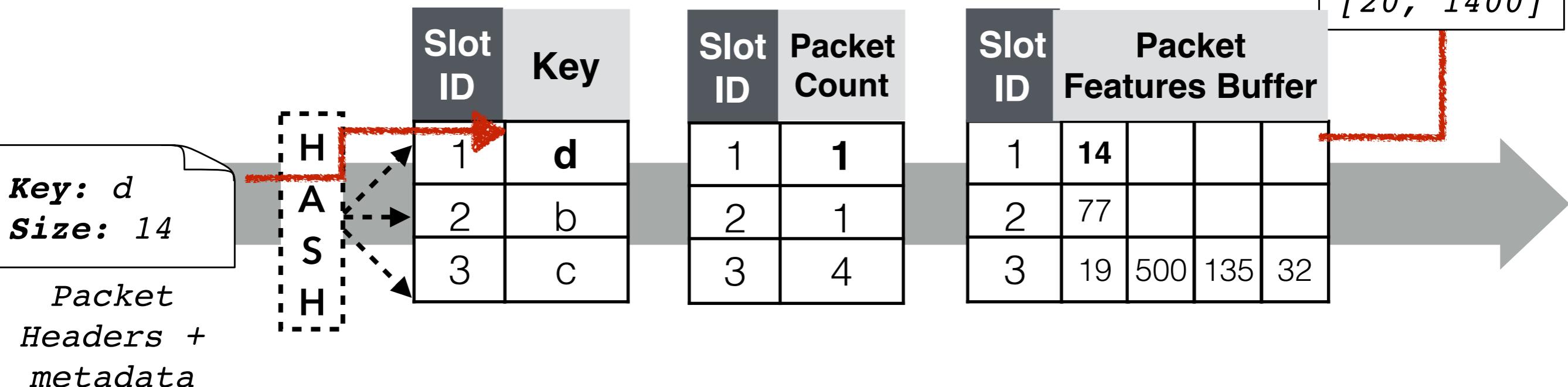
Replace and copy to software



# \*Flow as a Stateful Match Action Pipeline

## Untracked Flow:

Replace and copy to software



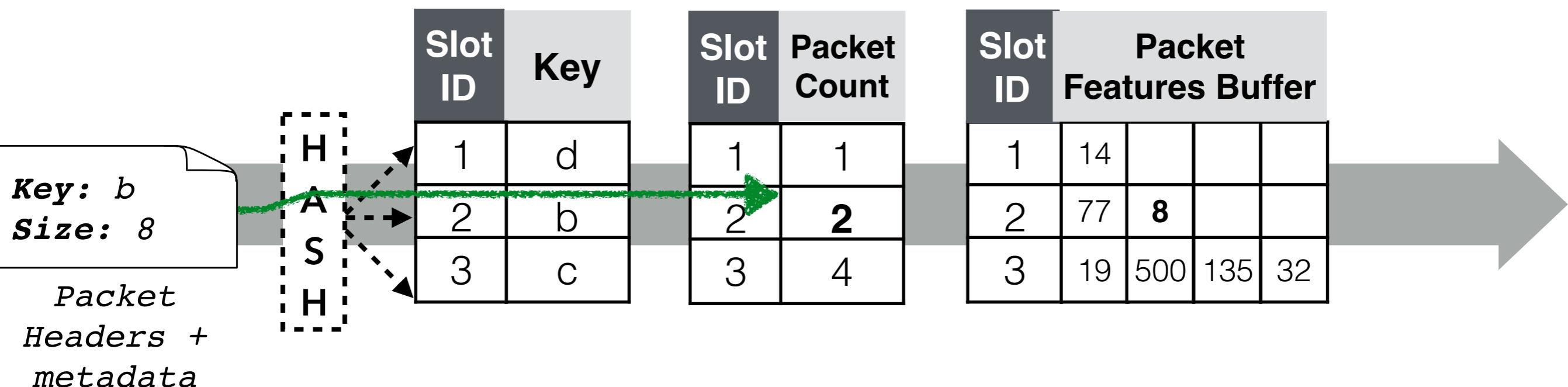
# \*Flow as a Stateful Match Action Pipeline

**Untracked Flow:**

Replace and copy to software

**Tracked Flow:**

Append tuple



# \*Flow as a Stateful Match Action Pipeline

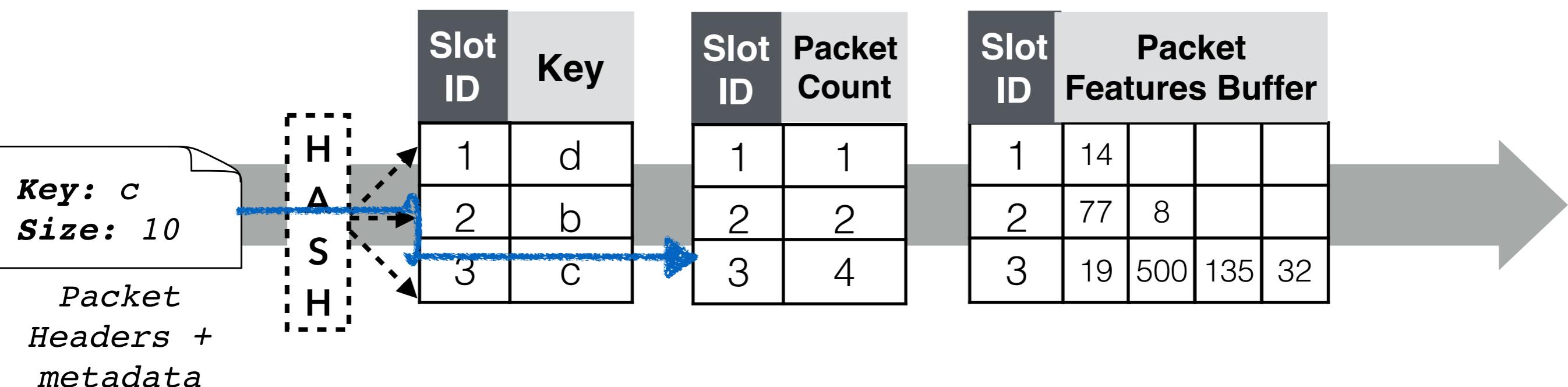
## Untracked Flow:

Replace and copy to software

## Tracked Flow:

Append tuple

Tracked Flow,  
Buffer Full: rollover  
buffer



# \*Flow as a Stateful Match Action Pipeline

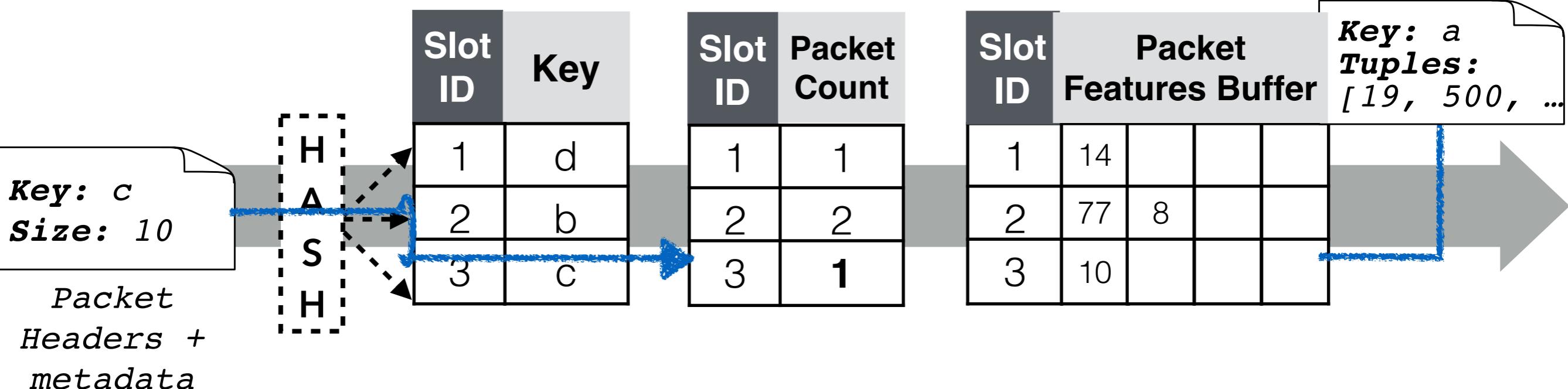
**Untracked Flow:**

Replace and copy to software

**Tracked Flow:**

Append tuple

**Tracked Flow,  
Buffer Full:** rollover  
buffer

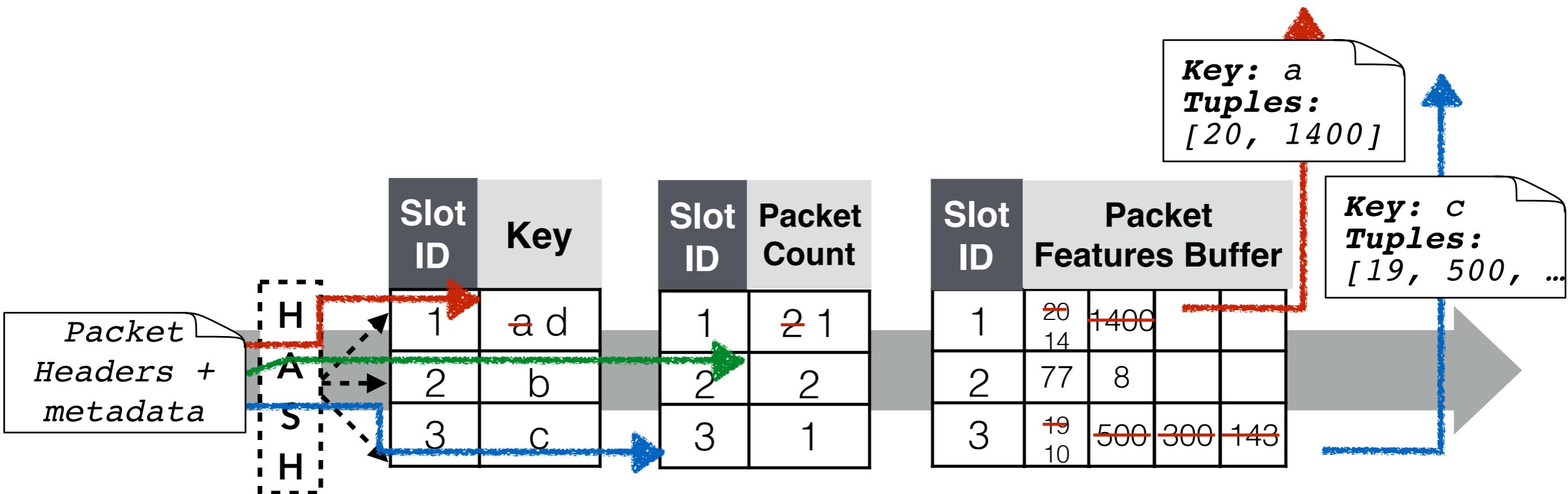


# \*Flow as a Stateful Match Action Pipeline

**Untracked Flow:**  
Replace and copy to software

**Tracked Flow:**  
Append features

**Tracked Flow,  
Buffer Full:** rollover  
buffer

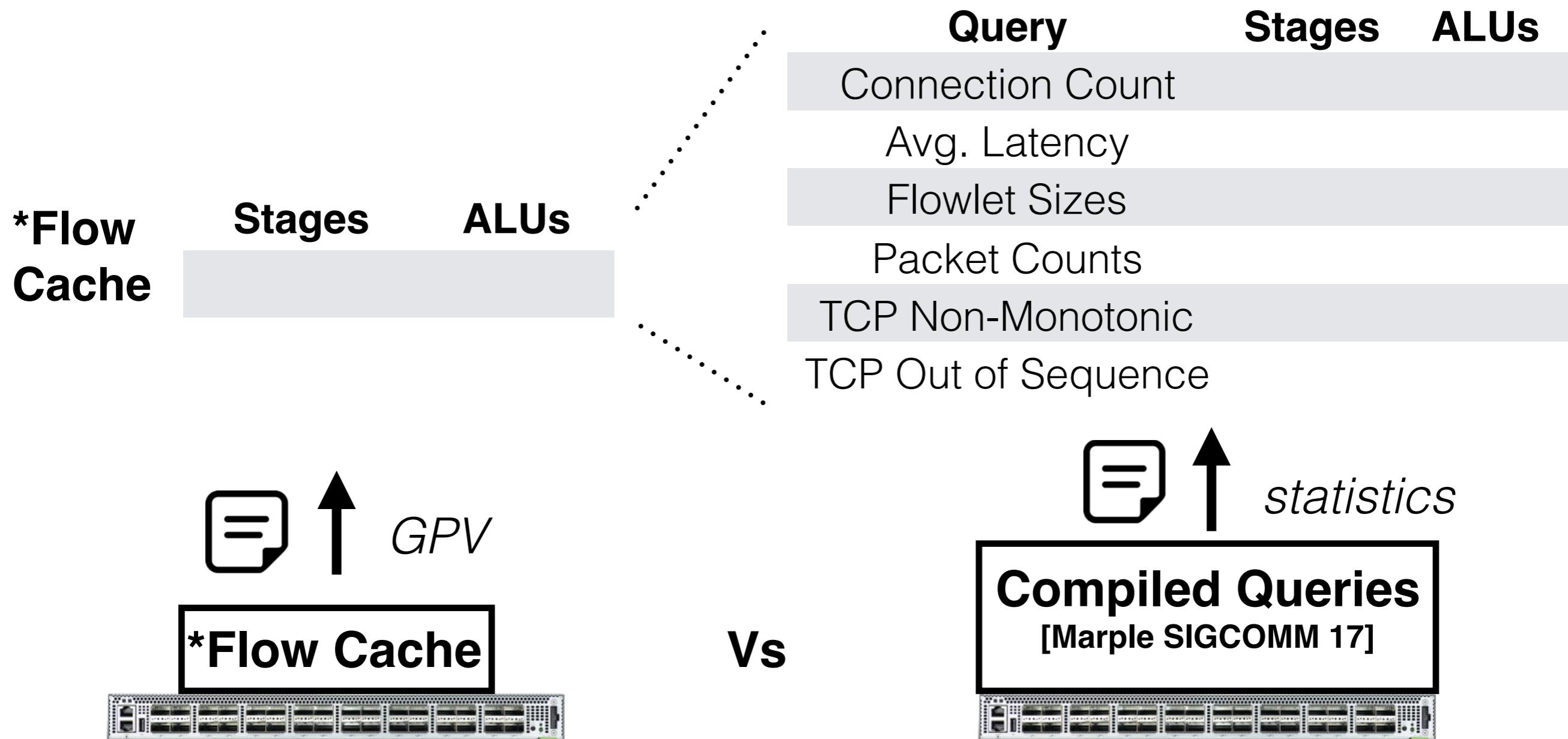


# Outline

- Motivation
- Design
- Implementation
- **Evaluation**

# Evaluation: Concurrency and Efficiency

# Evaluation: Concurrency and Efficiency

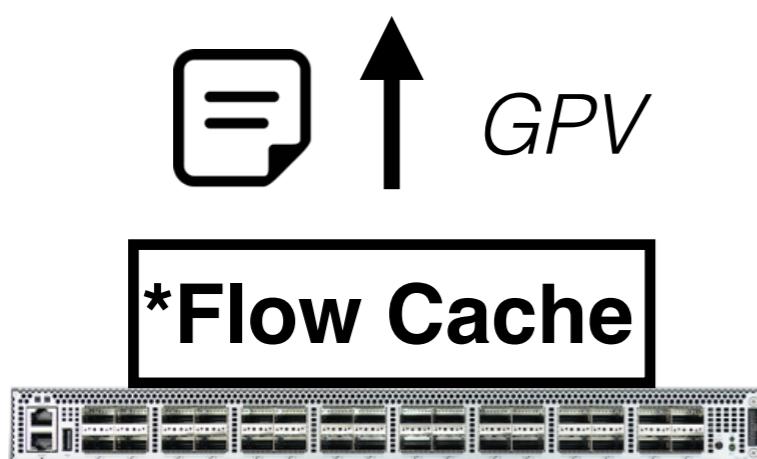


# Evaluation: Concurrency and Efficiency

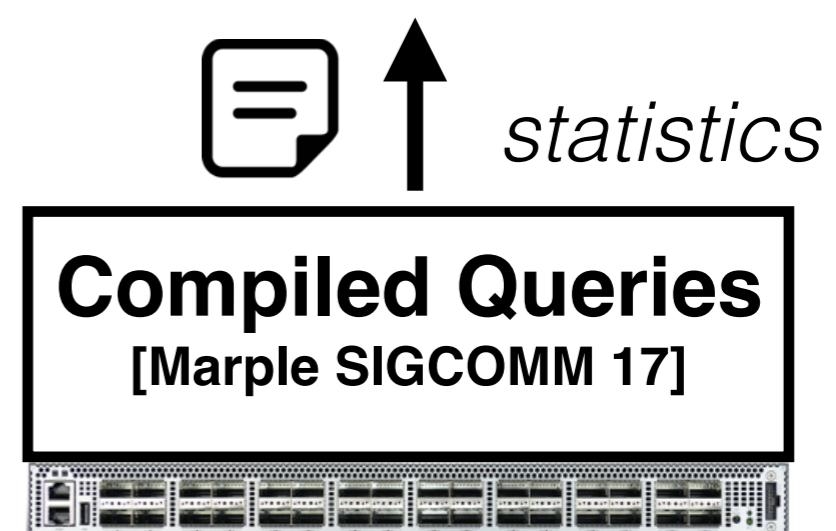
The \*Flow Cache can service **any** query but only requires as many hardware resources as  $\sim 1$  compiled query.

| *Flow Cache | Stages | ALUs |
|-------------|--------|------|
|             | 11     | 33   |

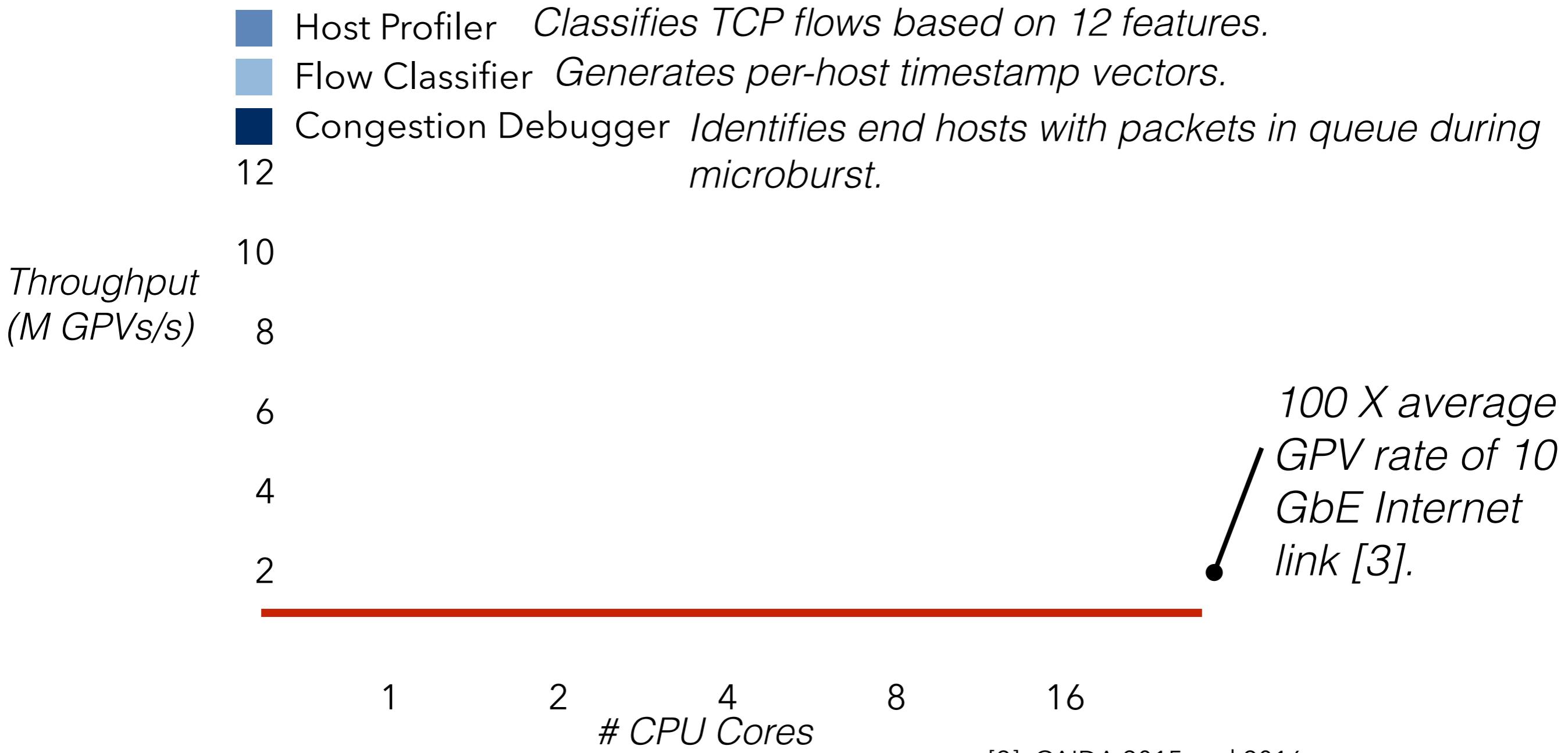
| Query               | Stages | ALUs |
|---------------------|--------|------|
| Connection Count    | 4      | 10   |
| Avg. Latency        | 6      | 11   |
| Flowlet Sizes       | 11     | 31   |
| Packet Counts       | 5      | 7    |
| TCP Non-Monotonic   | 5      | 6    |
| TCP Out of Sequence | 7      | 14   |



Vs

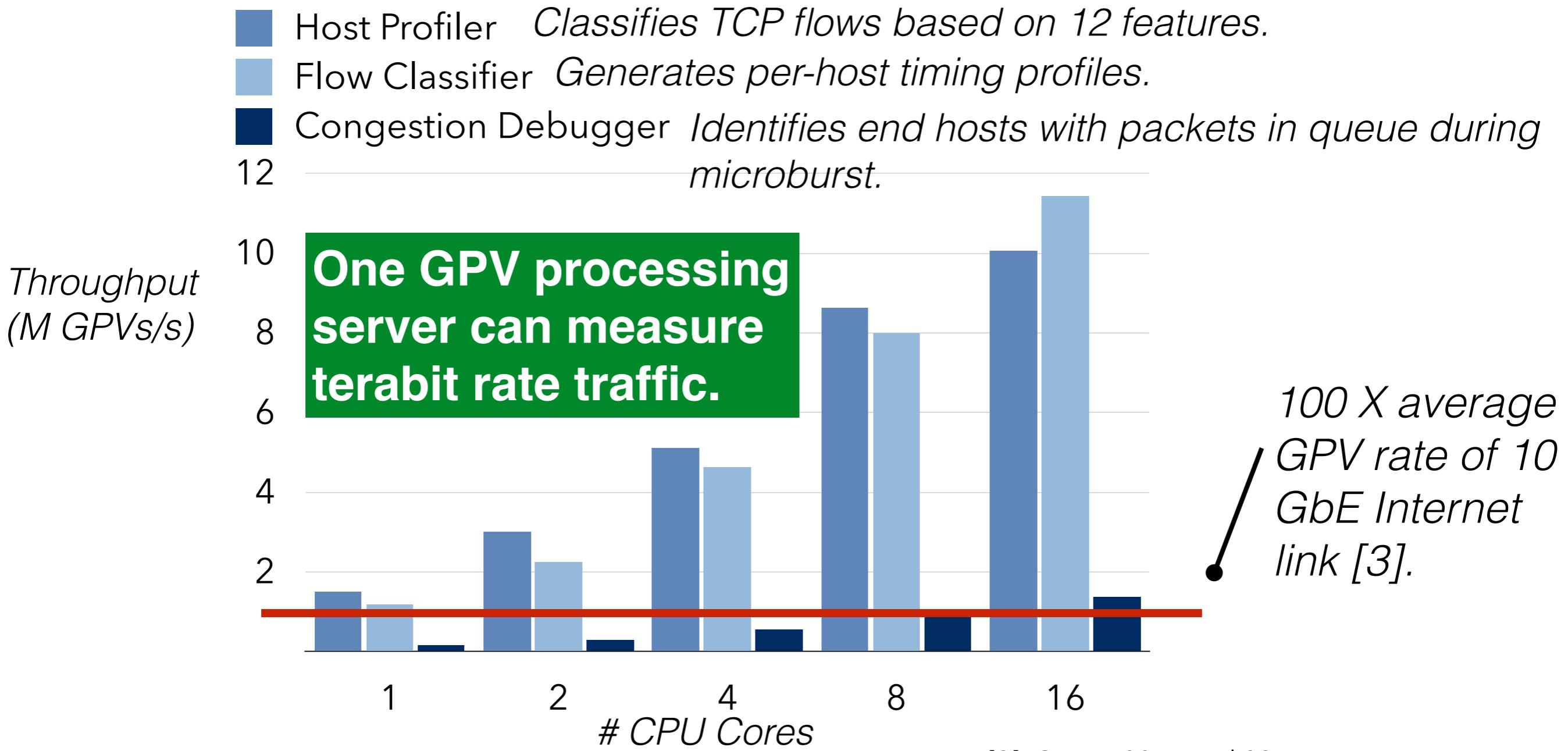


# Evaluation: Concurrency and Efficiency



[3]: CAIDA 2015 and 2016 trace stats  
[https://www.caida.org/data/passive/trace\\_stats/](https://www.caida.org/data/passive/trace_stats/)

# Evaluation: Concurrency and Efficiency



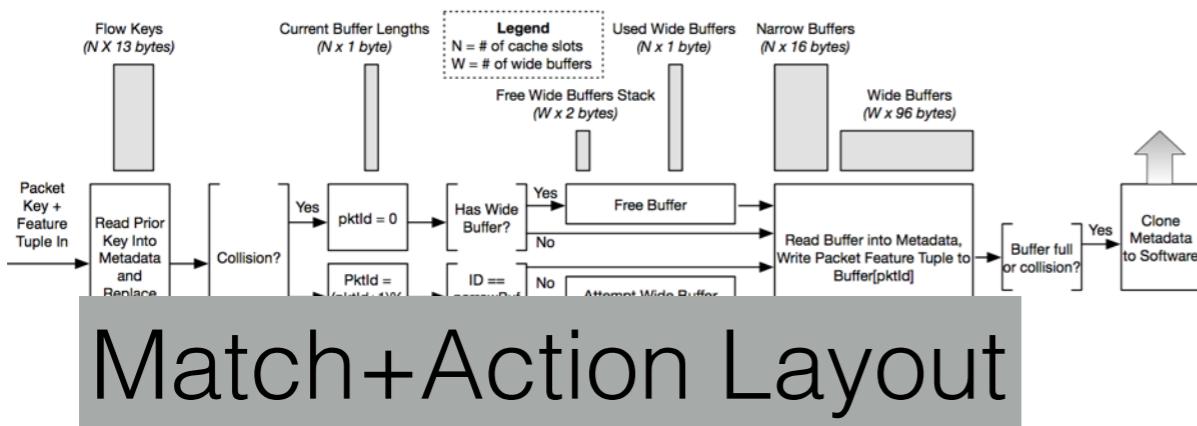
[3]: CAIDA 2015 and 2016 trace stats  
[https://www.caida.org/data/passive/trace\\_stats/](https://www.caida.org/data/passive/trace_stats/)

# Outline

- Motivation
- Design
- Implementation
- Evaluation

# In The Paper

# On Github



## Match+Action Layout

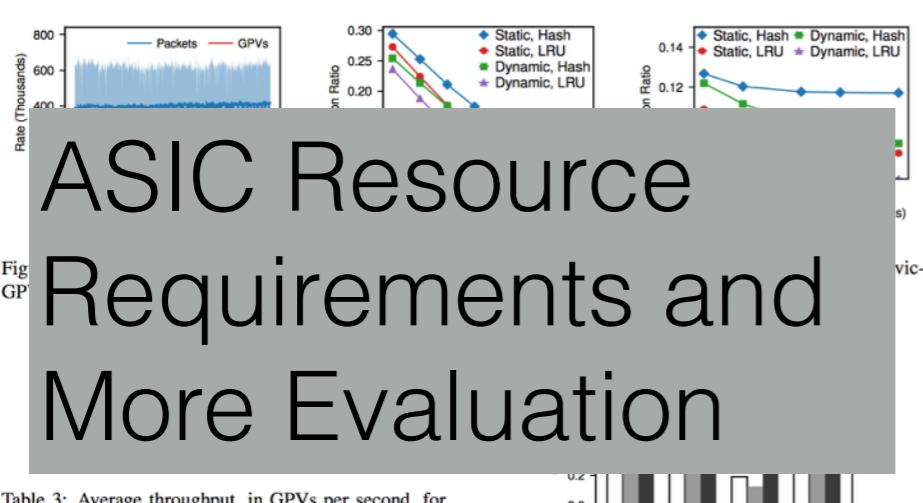


Table 3: Average throughput, in GPVs per second, for \*Flow agent and applications.

|             | Efficient | Flexible | Concurrent | Dynamic |
|-------------|-----------|----------|------------|---------|
| Netflow     | ✓         | ✗        | ✓          | ✓       |
| Software    | ✗         | ✓        | ✓          | ✓       |
| PFE Queries | ✓         | ✓        | ✗          | ✗       |
| *Flow       |           |          |            | ✓       |

## Background

[www.github.com/jsonch/starflow](http://www.github.com/jsonch/starflow)

```
// Tables.
table UpdateKey { default_action :UpdateKeyAction(); }
table UpdateFeatures { default_action
    :UpdateFeaturesAction(); }
table ResetFeatures { default_action
    :ResetFeaturesAction(); }
```

```
// Actions.
// Update key for every packet.
```

3.2 Tb/s  
\*Flow Cache  
Prototype for  
Barefoot Tofino

```
register_write(evictBufArr, evictBufArr, 0);
register_write(evictBufArr, 0, tempMfr.evictBufPos+1);
register_write(evictBufKey, tempMfr.evictBufPos,
    tempMfr.key);
register_write(evictBufPktCt, tempMfr.evictBufPos,
    tempMfr.pktCt);
```

}

# Conclusion (and Thank You for Listening!)

*What are **the right** measurement subtasks to do in hardware?*

✓ **Concurrency**

✓ **Efficiency**

[www.github.com/jsonch/starflow](http://www.github.com/jsonch/starflow)

