

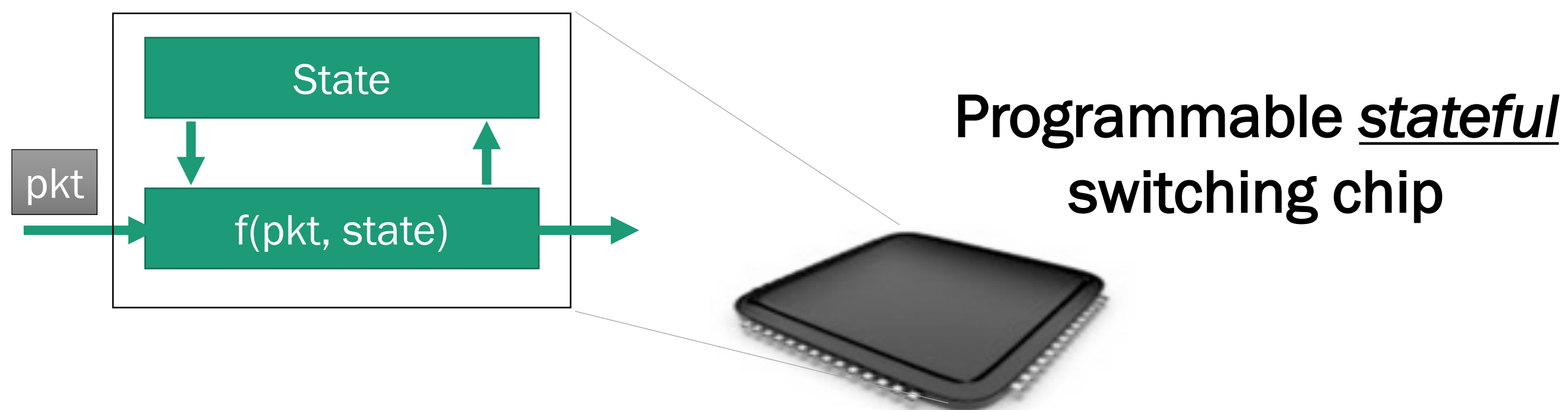
# Relaxing constraints in stateful network data plane design

Carmelo Cascone<sup>^</sup>, Roberto Bifulco<sup>\*</sup>, Salvatore Pontarelli<sup>+</sup>, Antonio Capone<sup>^</sup>  
<sup>^</sup> Politecnico di Milano (Italy), <sup>\*</sup> NEC Laboratories Europe (Germany), <sup>+</sup> Univ. Roma Tor Vergata (Italy)

## 1. INTRODUCTION

Program and run at line rate algorithms that read and modify data plane's state

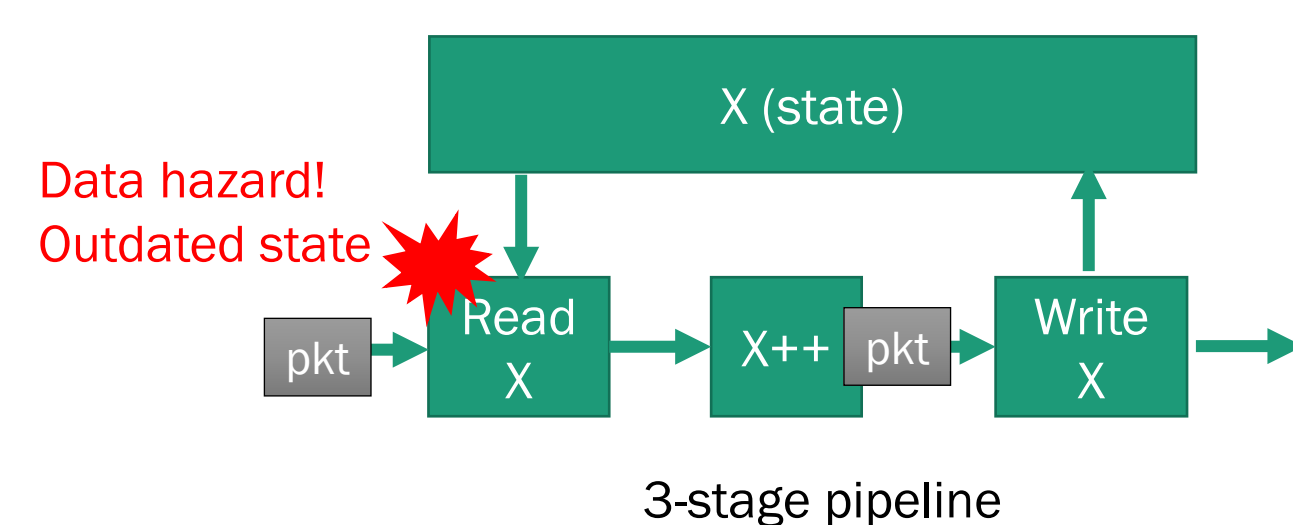
- E.g.: Stateful firewall, dynamic NAT, flowlet load balancing, AQM, measurement, etc.



## 2. PROBLEM STATEMENT

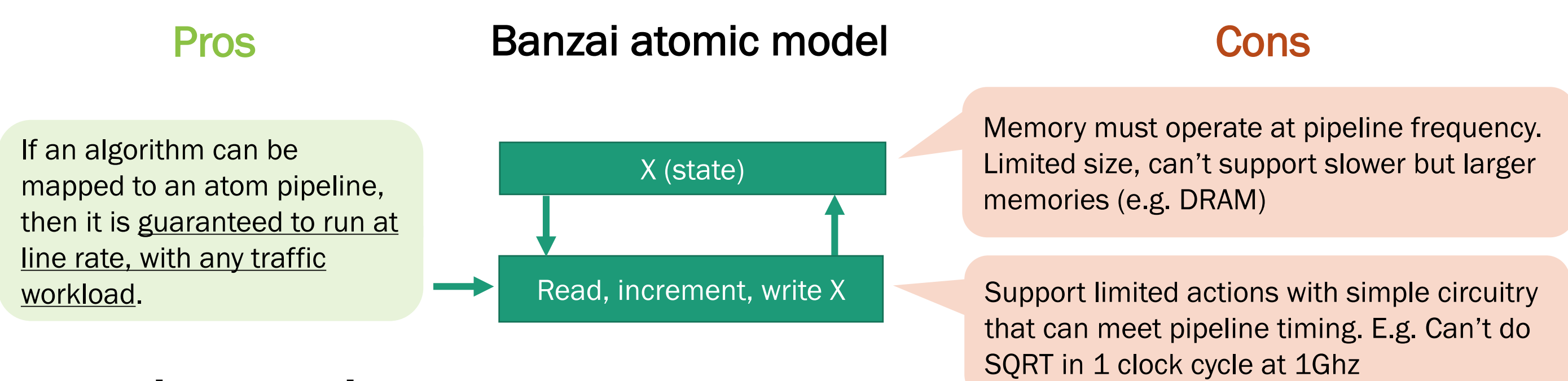
- Pipelining is the way to scale for high throughput (Tb/s)
- When pipelining, accessing state at different stages of the pipeline can cause data hazards

**Example:**  
A packet counter. For each packet increase the value of X.



## State of the art

- Reconfigurable Match Table (RMT) [SIGCOMM 2013]
  - Programmable parser, actions, table size, stateless
  - 640 Gb/s non-blocking line rate, 6.5 Tb/s Tofino Chip
- Domino/Banzai [SIGCOMM 2016]
  - Extends RMT with support for stateful actions named "Atoms"
  - Atoms can be pipelined to implement complex algorithms
  - Strict constraints on the atom execution time
  - State read-modify-write must be executed in 1 clock tick



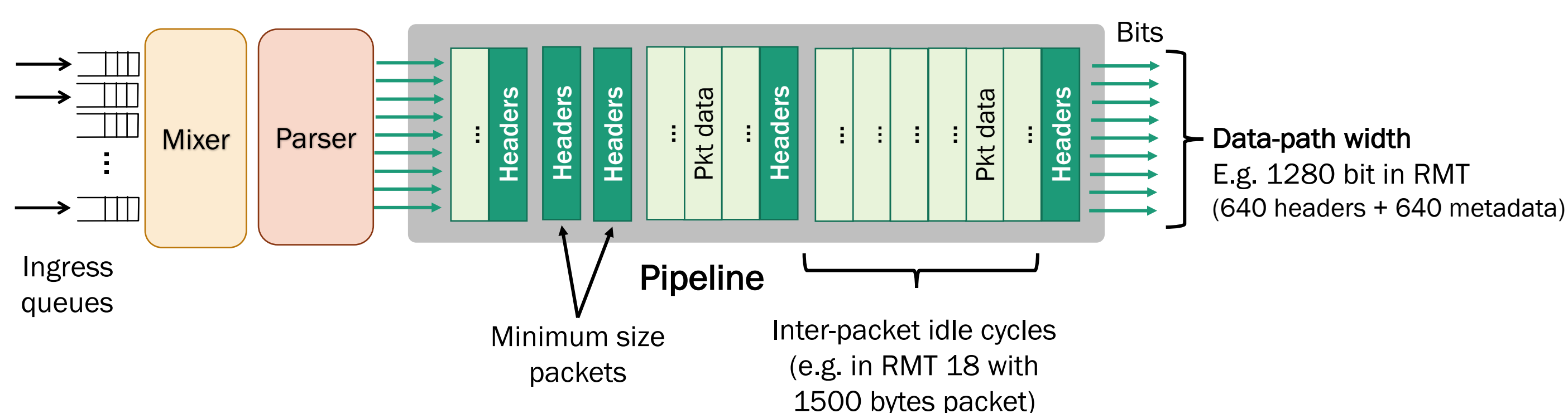
## Research question:

- If we allow for stateful processing blocks that span many clock cycles:
  - What is the risk of data hazards with realistic traffic conditions?
  - What is the throughput when using locking to access the memory?
  - How much silicon is needed to implement such a locking scheme?

## 3. OBSERVATIONS

### 1. Pipeline's header processing rate depends on packet size

- Packets are read from input ports in chunks, e.g. 80 bytes in RMT
  - 80 bytes \* 1 Ghz (chip clock freq.) = 640 Gb/s line rate
- Larger packets cause inter-packet idle cycles
  - Minimize risk of data hazards



### 2. Distinguish between per-flow and global state

- Global:** shared among all packets
- Per-flow:** shared by packets of the same flow
  - Different flows can be processed in parallel
- E.g. a stateful firewall needs only per-flow state, a DNAT both

## 4. MOTIVATING EXPERIMENTS

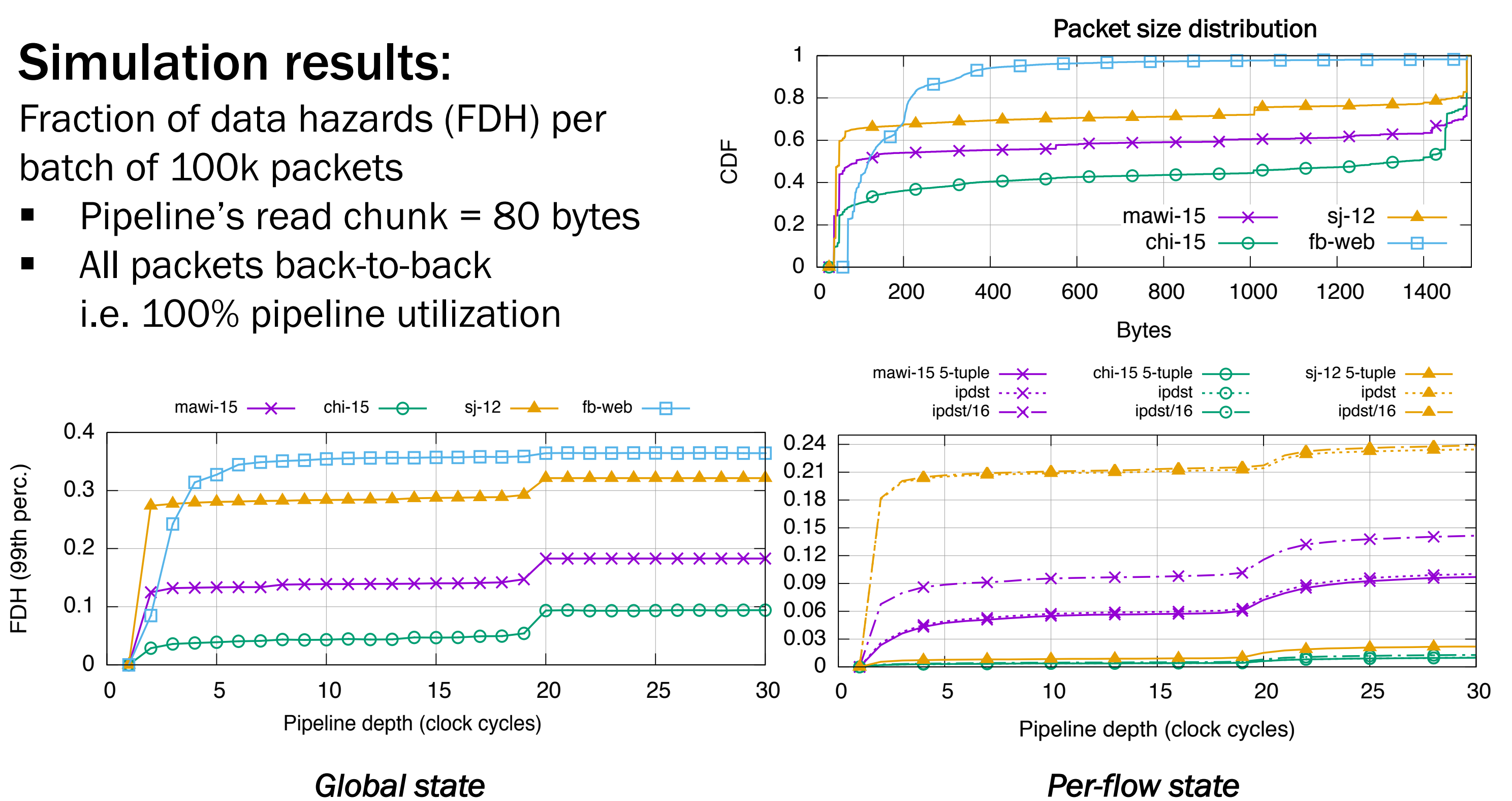
Evaluate the risk of data hazards using simulations with real traffic traces

Trace	Provider	Description	Date	Num pkts	Num flows per 1m pkts		
					5-tuple	ipdst	ipdst/16
chi-15	CAIDA	10Gb/s backbone link in Chicago. Usual conditions.	Feb 19, 2015	3.5b	100.6k	57.7k	4.6k
sj-12	CAIDA	10Gb/s backbone link in San Jose. Unusually high number of 5-tuples.	Nov 15, 2012	3.6b	249k	17k	2k
mawi-15	MAWI	1Gb/s backbone link in Japan. High volume of anomalous traffic.	Jul 21, 2015	135m	40.8k	17.3k	1.7k
fb-web	Facebook	Packet samples from 10 most active ToR switches in a web cluster	2015	447m	n/a	n/a	n/a

### Simulation results:

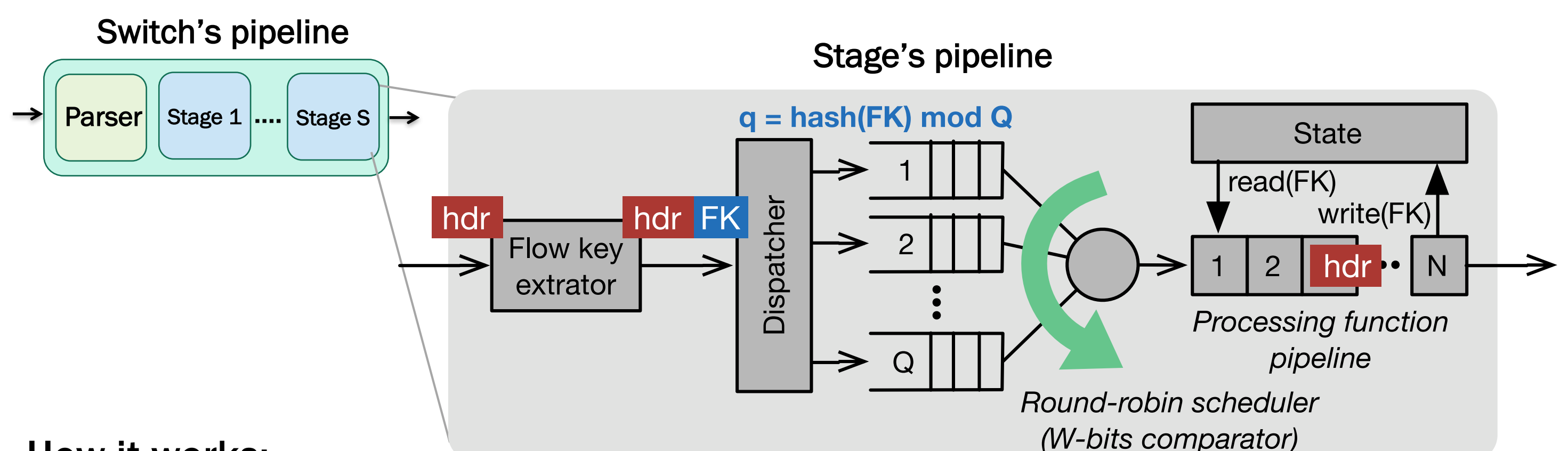
Fraction of data hazards (FDH) per batch of 100k packets

- Pipeline's read chunk = 80 bytes
- All packets back-to-back i.e. 100% pipeline utilization



## 5. APPROACH: MEMORY LOCKING

If two packets of the same flow arrive back-to-back, processing is paused for the second packet until the first one has left the stage pipeline.

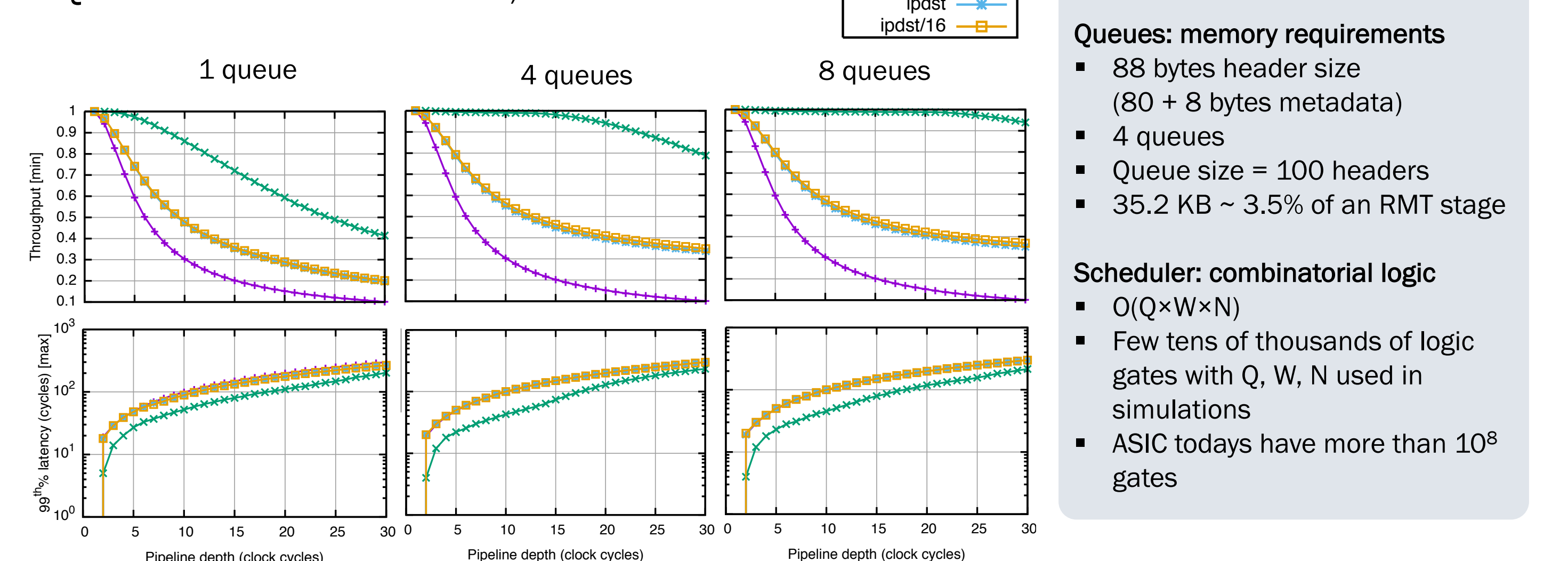


### How it works:

- Each packet is associated with a flow key (FK)
- A dispatcher enqueues packets by hashing on the FK
- A round-robin scheduler decides if a queue can be served by looking at the head-of-line's FK and comparing it to what is currently in the stage's pipeline.
- Comparison is performed by reducing the space of the FK to few bits (W)

### Simulation results:

- Trace sj-12 (worst case)
- Queue size = 10 headers, W = 4 bits



### Clock cycle budget (and latency) for all traces:

Maximum number of clock cycles (limited to 30) per processing function, to sustain a given throughput. W = 4 bits. Latency values are given for 1 Ghz clock frequency, i.e. 1 clock cycle = 1 ns.

Thrtpt	Q <sub>ten</sub>	Q	chi-15				sj-12				mawi-15				fb-web global	
			5-tuple	ipdst	ipdst/16	global	5-tuple	ipdst	ipdst/16	global	5-tuple	ipdst	ipdst/16	global		
100%	10	4														
		8														
		16														
		1	20 (174ns)	20 (190ns)	21 (230ns)	8 (282ns)	4 (49ns)					2 (18ns)	2 (20ns)	2 (35ns)		2 (10ns)
		100	8 (137ns)	30 (144ns)	30 (259ns)	8 (152ns)	1					2 (12ns)	2 (14ns)	2 (25ns)		
99.9%	10	4	8 (16ns)	8 (16ns)	8 (18ns)	4 (18ns)	2 (5ns)									
		8	14 (33ns)	14 (31ns)	14 (38ns)	2 (4ns)										
		16	17 (37ns)	18 (43ns)	16 (42ns)	2 (5ns)										
		1	27 (568ns)	27 (618ns)	26 (605ns)	8 (282ns)	6 (143ns)	2 (86ns)	2 (84ns)	1		3 (42ns)	3 (48ns)	3 (100ns)	2 (60ns)	2 (10ns)
		100	4 (175ns)	30 (192ns)	30 (320ns)	15 (52ns)	2 (79ns)	2 (77ns)				4 (41ns)	4 (52ns)	4 (135ns)		
99%	10	4	21 (80ns)	21 (89ns)	21 (89ns)	7 (60ns)	3 (14ns)				2 (13ns)	2 (14ns)	1		2 (10ns)	
		8	30 (142ns)	30 (148ns)	30 (184ns)	10 (45ns)	1				5 (34ns)	4 (31ns)	2 (18ns)			
		16	30 (129ns)	30 (138ns)	30 (184ns)	11 (47ns)					6 (42ns)	4 (31ns)	2 (18ns)			
		1	30 (116ns)	30 (122ns)	30 (180ns)	12 (47ns)					7 (52ns)	4 (31ns)	2 (18ns)			
		100	4 (175ns)	30 (192ns)	30 (320ns)	22 (1.1us)	9 (842ns)	8 (268ns)	2 (86ns)	2 (84ns)	2 (171ns)	9 (422ns)	8 (380ns)	5 (316ns)	4 (379ns)	2 (10ns)
99%	10	4	30 (137ns)	30 (144ns)	30 (259ns)	22 (1.1us)	3 (285ns)	3 (285ns)	3 (285ns)	2 (171ns)	24 (1.7us)	17 (1.1us)	7 (572ns)			
		8	30 (137ns)	30 (144ns)	30 (259ns)	30 (1.9us)	3 (290ns)	3 (292ns)	3 (292ns)	30 (2.1us)	23 (2.2us)	8 (753ns)				
		16	30 (122ns)	30 (126ns)	30 (221ns)	25 (74ns)	2 (79ns)	2 (72ns)			30 (1.3us)	25 (2.5us)	8 (759ns)			
		1	21 (80ns)	21 (89ns)	21 (89ns)	7 (60ns)	3 (14ns)	1			2 (13ns)	2 (14ns)	1		2 (10ns)	
		100	4 (175ns)	30 (192ns)	30 (320ns)	22 (1.1us)	3 (285ns)	3 (285ns)	3 (285ns)	3 (285ns)	30 (2.1us)	23 (2.2us)	8 (753ns)			