FAB: Towards Flow-aware Buffer Sharing in Programmable Switches Maria Apostolaki





Joint work with Laurent Vanbever & Manya Ghobadi

An old story...

An old story... Fan-in causing queue built-up



Senders

An old story... Fan-in causing drops





An old story... **Drops increase Flow Completion Times**























Buffer management: the algorithm according to which ports/queues of a device share a common buffer

Buffer management: the algorithm according to which ports/queues of a device share a common buffer

Most of today's devices have a shared buffer to absorb bursts

How many packets can each port store in the common buffer?



How many packets can each port store in the common buffer?



remaining excessive packets will be dropped

How many packets can each port store in the common buffer?



remaining excessive packets will be dropped



Let's give... ...half of the buffer to each port! ...small fraction to each port!



FAB: Towards Flow-aware Buffer Sharing in Programmable Switches





Joint work with Laurent Vanbever & Manya Ghobadi

Outline

Background

FAB

Initial Results

Practicality

Outline

Background

FAB

Initial Results

Practicality

Complete Partitioning:

statically allocated buffer space per port

Complete Partitioning:

statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Complete Partitioning:

statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port

Back to our story...



Senders



Shared buffer can host up to 180 packets



Senders



Long flows will consume as much buffer as there is available





Senders



Short flows will attempt to instantaneously store at most 75 packets in the shared buffer



Complete Partitioning: statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port

Complete Partitioning: static buffer space per port = 10 packets





8.0

Time (s)

1

Receiver of long flows buffers up to 10 packets





0.8

Receiver of long flows buffers up to 10 packets





Receiver of long flows buffers up to 10 packets





used by port 1

Time (s)

1

Unused 95%



Receiver of short flows buffers up to 10 packets





Receiver of short flows buffers up to 10 packets





Buffer is 90% empty





used by port 2





Buffer is 90% empty, yet the burst is not absorbed



Pkts in Buffer



90%

Complete Partitioning: statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port

works for balanced traffic wastes buffer otherwise

Complete Partitioning:

statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port
Complete Sharing: unrestricted buffer space per port



0

Long flows use all buffer



0

Pkts in Buffer

0.8

1 Time (s)

Long flows use all buffer



Pkts in Buffer

0.8 1 Time (s)

Long flows use all buffer



Pkts in Buffer



3 1 Time (s)



Upon arrival, the burst finds the buffer fully occupied







Three common buffer management techniques with pros & cons

Complete Partitioning:

statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port



Three common buffer management techniques with pros & cons

Complete Partitioning:

statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port

Dynamic Sharing: fraction (α) of remaining buffer per port



n: Number of congested ports N: Buffer Size (packets/Bytes) α: per port/queue parameter



Dynamic Sharing: a fraction (α) of remaining shared buffer per port



6 0.8 1 Time (s)



180 = 90

0.8 1 Time (s)



- -



50%

- -





- -







Three common buffer management techniques with pros & cons

Complete Partitioning:

statically allocated buffer space per port

Complete Sharing:

unrestricted buffer space per port

Dynamic Sharing: fraction (α) of remaining buffer per port

adapts fairly to loadignores queue content

Why should we care about queue content?



Why should we care about queue content?

Short flows would benefit more from buffer



Outline

Background

FAB

Initial Results

Practicality

FABULOUS Sharing (FAB) improves dynamic sharing by

using multiple α per queue/port
Two packets of same ingress and egress port
which arrived together, might see different limits

using multiple α per queue/port Two packets of same ingress and egress port which arrived together, might see different limits

mapping packets to an α according to priority

- α is proportionate to the per-flow expected benefit from buffering

using multiple α per queue/port Two packets of same ingress and egress port which arrived together, might see different limits

deciding packet's priority directly in the data plane e.g. prioritizing packets based on flow size

- mapping packets to an α according to their priority in buffering α is proportionate to the per-flow expected benefit from buffering

Flow-aware Buffer Sharing (FAB)

$\alpha = 0.1$ for long flows $\alpha = 10$ for short flows

Flow-aware Buffer Sharing (FAB)



Time (s)

. .

FAB maps long flows to $\alpha = 0.1$



Time (s)

. .

















FAB maps short flows to $\alpha = 10$


















Pkts in Buffer

Outline

Background

FAB

Initial Results

Practicality

Simulation Results

Short-Flows Workload

Mixed Workload

Short-Flows Workload



Senders



Receiver

Dynamic Sharing does not allow burst to be fully buffer, resulting in increased tail FCT



Dynamic Sharing does not allow burst to be fully buffer, resulting in increased tail FCT



50	70	90	99	
	Ē٨	\B_		

Mixed Workload



Senders



FAB limits long flows and allows the burst to use the buffer





FAB limits long flows and allows the burst to use the buffer



Outline

Background

FAB

Initial Results

Practicality

Is FAB practical?

approximates flow size with flow arrival time Use bloom filter to store flows that started in discrete time windows

enables dropping at the ingress based on FAB

Dropping decisions at ingress based on buffer occupancy and flow information

configures complete sharing

Disallow traffic manager to drop any packet as long as there is space the buffer

Outline

Background

FAB

Initial Results

Practicality

FAB: Flow-aware buffer sharing







Significantly decreases flow completion time

Splits the buffer space according to the expected benefit from using it for each flow.

Allocates more buffer to short flows, which are distinguished in the data plane



Questions?





