

Fastpass

A Centralized “Zero-Queue” Datacenter Network

Jonathan Perry



Amy Ousterhout



Hari Balakrishnan



Devavrat Shah

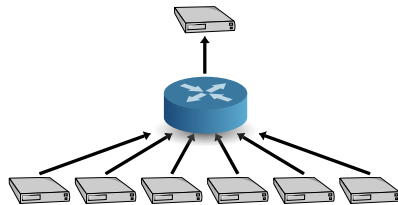


Hans Fugal

Ideal datacenter network properties

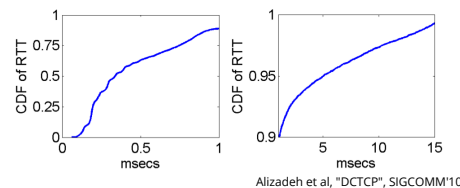
No current design satisfies
all these properties simultaneously

Scaling Memcache at
Facebook,
Fine-grained TCP
retransmissions



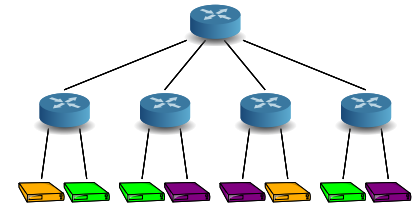
Burst Control

Datacenter TDMA,
Tail at scale,
pFabric, PDQ, DCTCP,
D3, Orchestra



Low Tail Latency

EyeQ, Seawall, Oktopus,
Hedera, VL2, Mordia,
SWAN, MATE, DARD

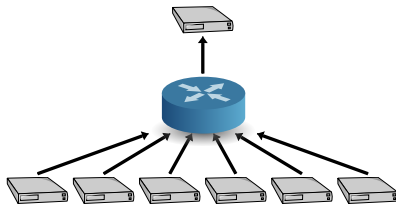


Multiple Objectives

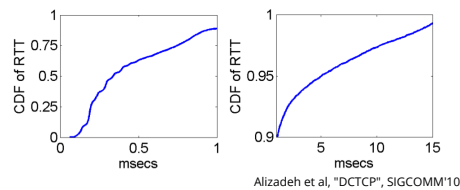
Fastpass goals

Is it possible to design a network that provides

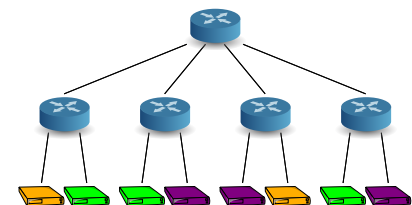
1. Zero network queues
2. High Utilization
3. Multiple app and user objectives



Burst Control



Low Tail Latency



Multiple Objectives

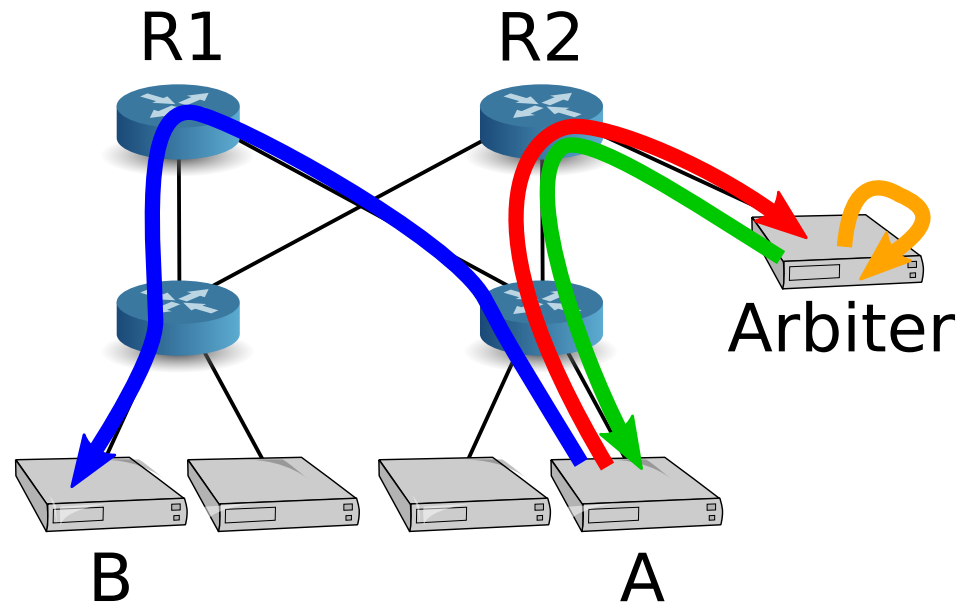
Centralized *arbiter* schedules
and assigns paths to all packets

Concerns with centralization:

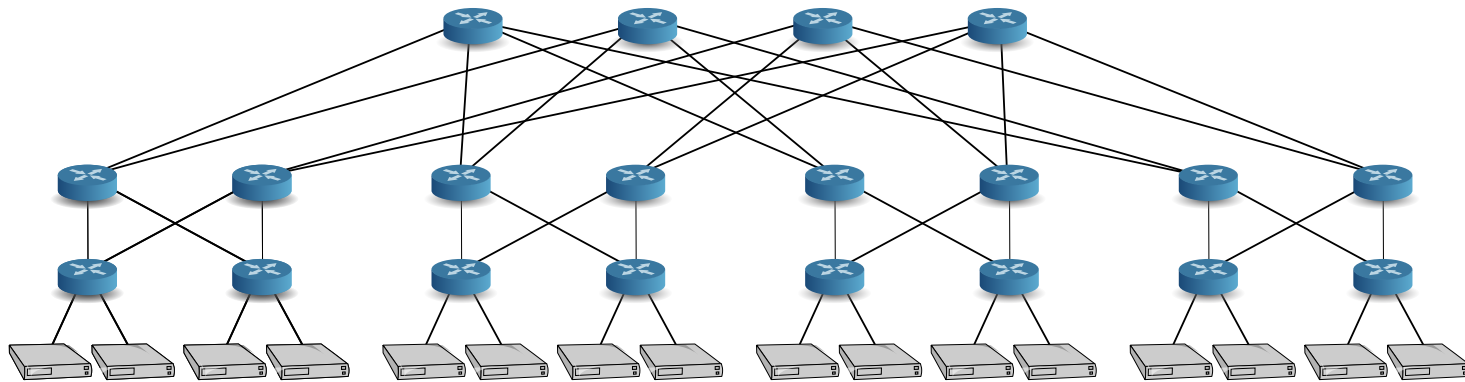
- Latency
- Scaling
- Fault tolerance

Example: Packet from A to B

5 μ s	A \rightarrow Arbiter	"A has 1 packet for B"
1-20 μ s	Arbiter	timeslot allocation & path selection
15 μ s	Arbiter \rightarrow A	"@t=107: A \rightarrow B through R1"
no queuing	A \rightarrow B	sends data



Scheduling and selecting paths



Timeslot = $1.2 \mu\text{s}$

Step 1: Timeslot Allocation

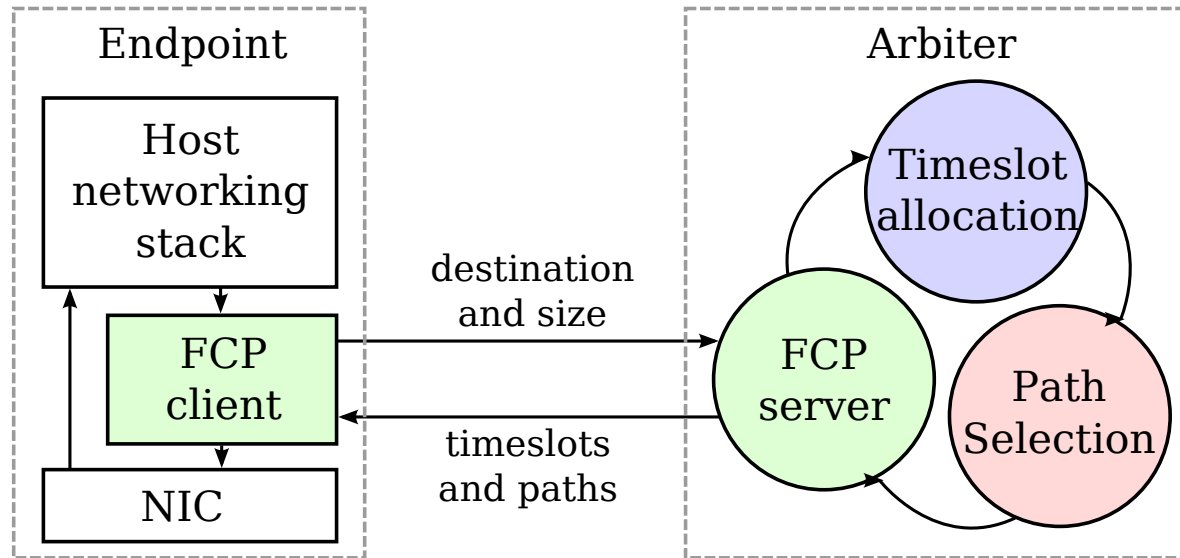
Choose a matching

Step 2: Path selection

Map matching onto paths

Arbiter treats network as a big switch

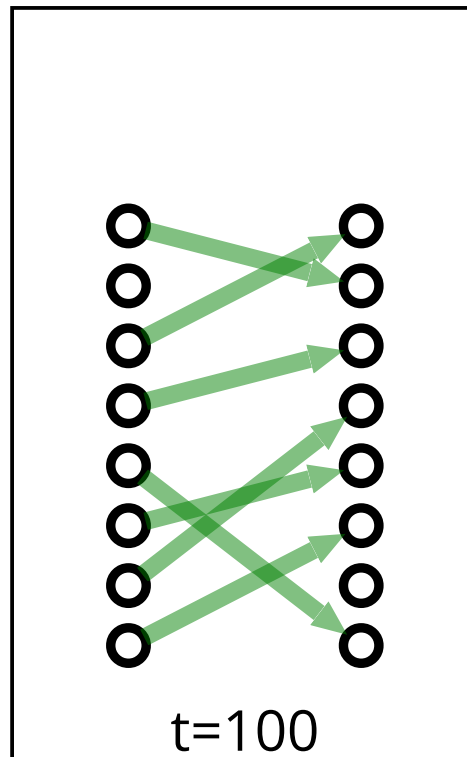
System structure



Challenges:

- Latency
- Scaling
- Fault tolerance

Timeslot allocation = maximal matching



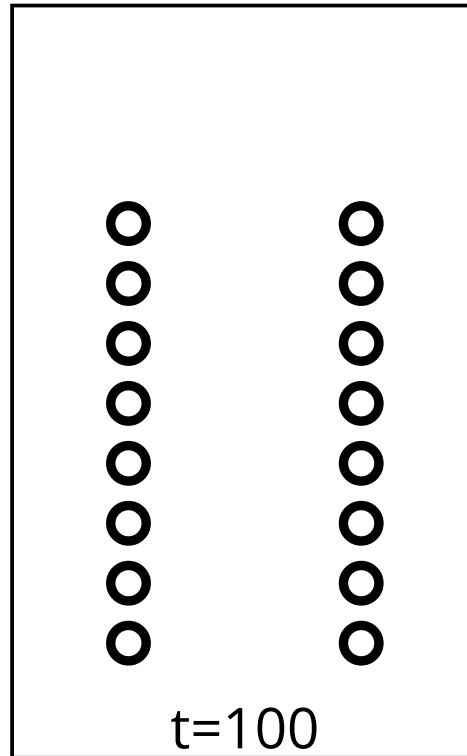
src	dst	pkts
1	2	3
src	dst	pkts
3	1	3
src	dst	pkts
7	4	1
src	dst	pkts
5	8	2
src	dst	pkts
4	3	4
src	dst	pkts
1	3	1
src	dst	pkts
8	6	3

~10ns per demand

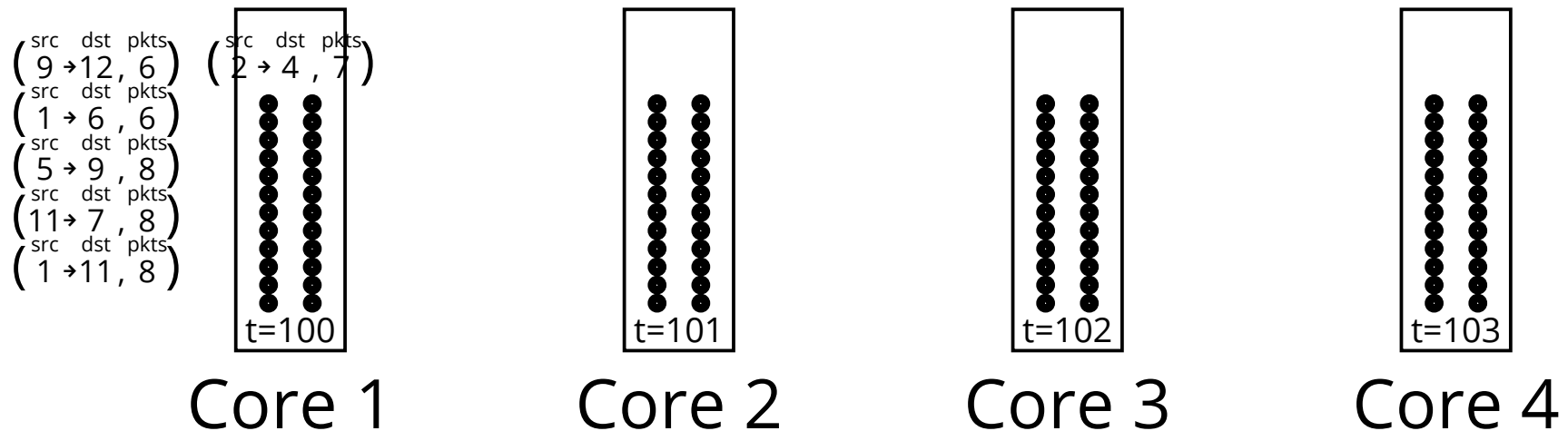
How to support different objectives?

Order matters!

src	→	dst	,	pkts
1	→	4	,	1
6	→	2	,	2
4	→	3	,	2
1	→	7	,	3
8	→	5	,	4
6	→	7	,	5
6	→	1	,	5
2	→	5	,	6



How to scale timeslot allocation?



Can pipeline timeslot allocation

2211.8 Gbits/s on 8 cores

Are maximal matchings good matchings?

Maximal Matching
2C network capacity

Optimal scheduler
C network capacity

Dai-Prabhakar '00:
Finite average latency

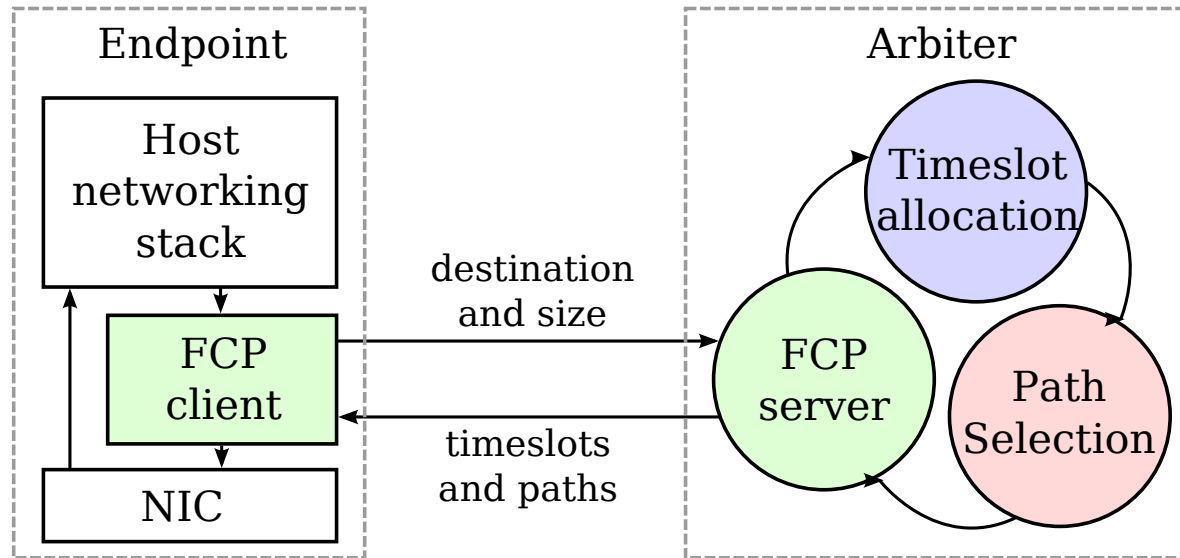


Finite average latency

Our theorem:

Average latency $\leq 2 \times$ Average latency

System structure



Challenges:

- Latency
- Scaling
- Fault tolerance

Fault-tolerance

Arbiter failures

Hot backups , TCP as last resort

Switch failures

Packet loss to arbiter

Experimental results

Timeslot allocation	2.21 Terabits/s with 8 cores
Path selection	>5 Terabits/s with 10 cores

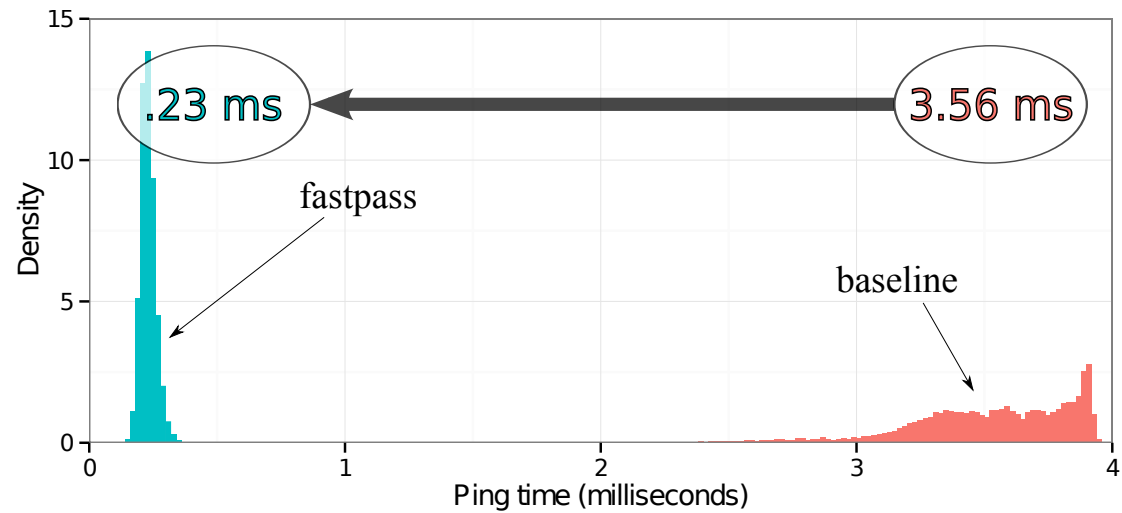
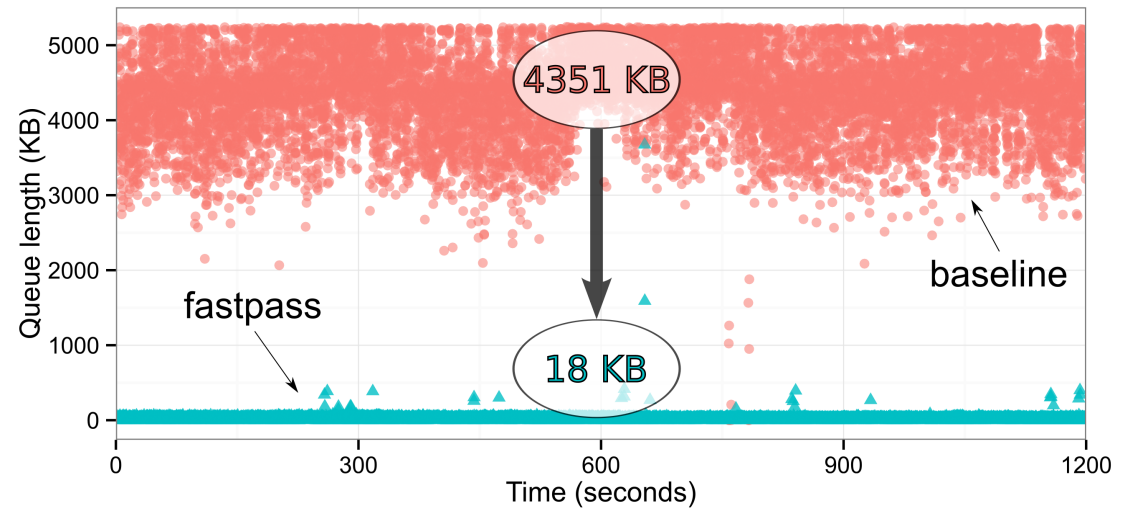
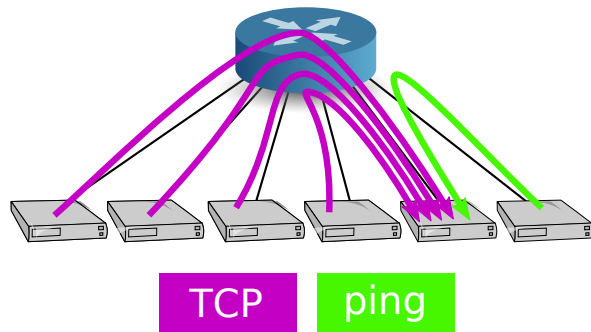
Facebook experiments:

Switch queue length, RTT

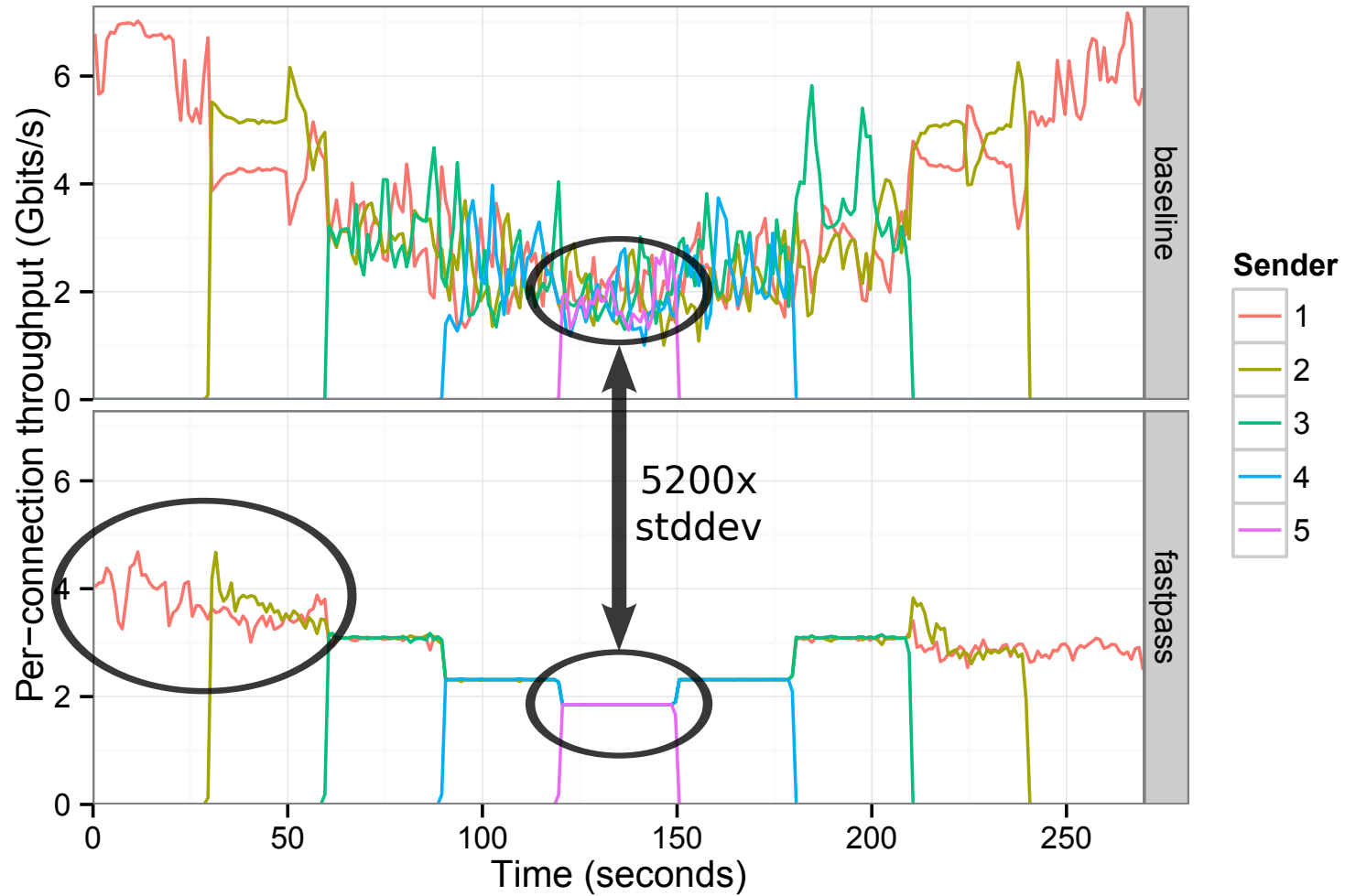
Convergence to network share

Reducing retransmission in production

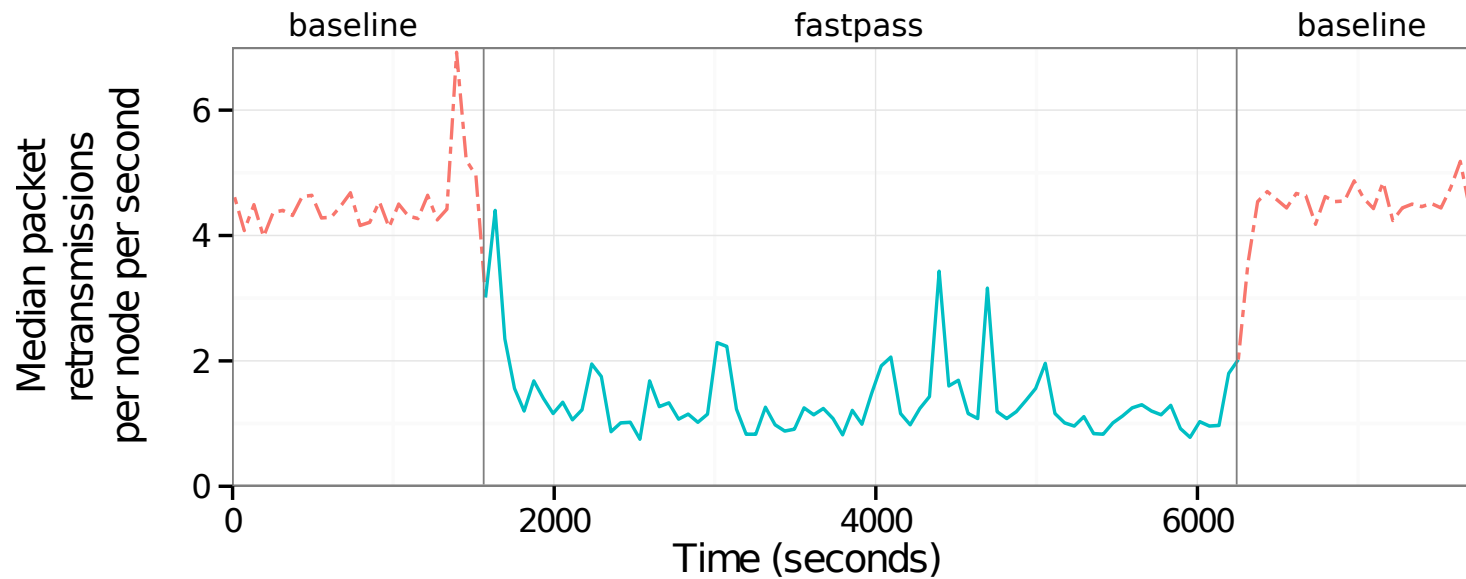
Queues & RTT



Convergence to network share



Reducing retransmissions in production



Each server: ~50k QPS

Benefits



- A: "Now I can see pictures of other's people's food and children so much more quickly...can't wait..>.>"
- B: "You forgot about [...] cats. I will say, faster pics of cats is probably worth some merit."

Benefits

Low user latency

Stronger network semantics

No packet drops, predictable latency, deadlines, SLAs

Developer productivity

Less dealing with bursts, tail latency, hotspots

Simplify building complex systems

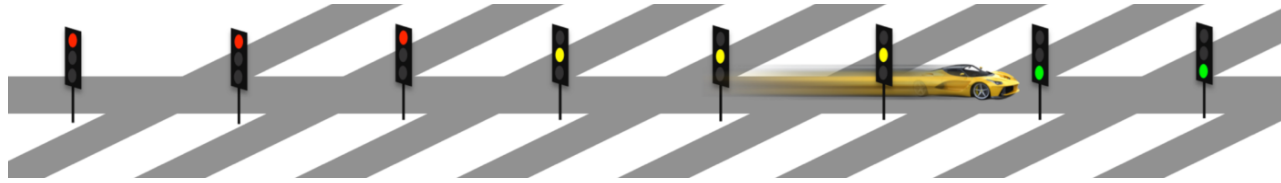
Lower infrastructure cost

Less over-provisioning

Fastpass enables new network designs

Traditional	Flow control	Congestion control	Update routing tables	Scheduling & Queue management	Packet forwarding
SDN	Flow control	Congestion control	Update routing tables	Scheduling & Queue management	Packet forwarding
Fastpass	Flow control	Congestion control	Per-packet path selection	Scheduling & Queue management	Packet forwarding
	Endpoint		Centralized		Switch

Fastpass: centralizes control at packet granularity
Switches can become even simpler and faster



Conclusion

Zero network queues
High Utilization
Multiple app and user objectives

Pushes centralization to a logical extreme
Opens up new possibilities for even faster networks

Code (MIT licensed): <http://fastpass.mit.edu>