

# Extending SDN to the Data Plane

Anirudh Sivaraman, Keith Winstein, Suvinay Subramanian,  
Hari Balakrishnan

M.I.T.

<http://web.mit.edu/anirudh/www/sdn-data-plane.html>

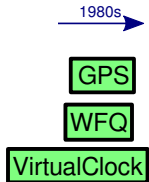
November 7, 2013

## Switch Data Planes today

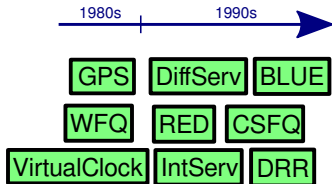
Two key decisions on a per-packet basis:

- ▶ Scheduling: Which packet to transmit next?
- ▶ Queue Management: How long can queues grow? Which packet to drop?

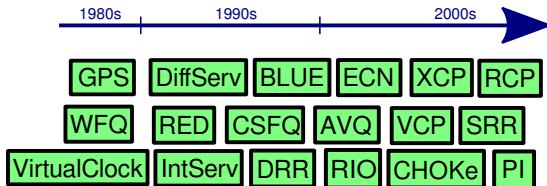
# The long lineage of in-network algorithms



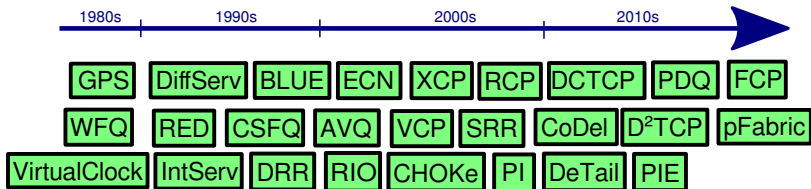
# The long lineage of in-network algorithms



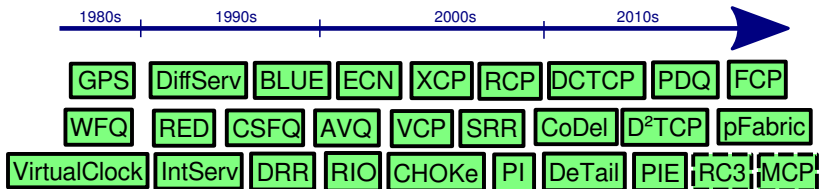
# The long lineage of in-network algorithms



# The long lineage of in-network algorithms



# The long lineage of in-network algorithms



# The Data Plane is continuously evolving

- ▶ Each scheme wins in its own evaluation.
- ▶ Some believe in a “silver bullet” knobless in-network method.



## We disagree: There is no silver bullet!

- ▶ Different applications care about different objectives.
- ▶ Applications use different transport protocols.
- ▶ Networks are heterogeneous.

## Our work:

- ▶ Quantify non-universality of in-network methods.
- ▶ Extend SDN to the Data Plane to handle in-network diversity.

## Quantifying “No Silver Bullet”: Network Configurations

<b><u>Configuration</u></b>	<b><u>Description</u></b>
<b>CoDel+FCFS</b>	One shared FCFS queue with CoDel
<b>CoDel+FQ</b>	Per-flow fair queueing with CoDel on each queue
<b>Bufferbloat+FQ</b>	Per-flow fair queueing with deep buffers on each queue

## Quantifying “No Silver Bullet”: Workloads and Objectives

<u>Workload</u>	<u>Description</u>	<u>Objective</u>
<b>Bulk</b>	Long-running bulk transfer flow	Max. throughput
<b>Web</b>	Switched flow with ON and OFF periods	Min. 99.9 %ile flow completion time
<b>Interactive</b>	Long-running interactive flow	Max. $\frac{\text{throughput}}{\text{delay}}$ (power)

# Quantifying “No Silver Bullet”

CoDel+FCFS

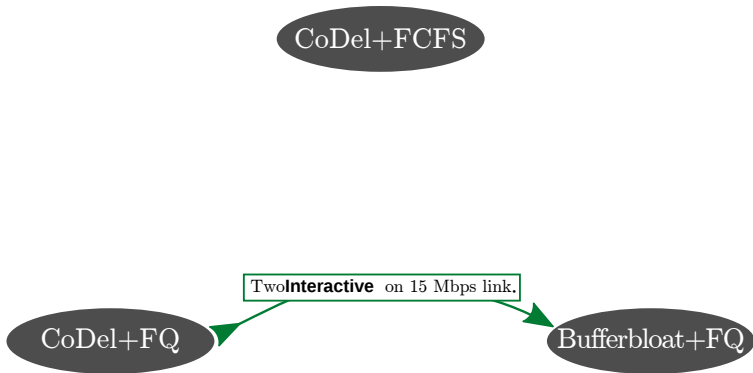
CoDel+FQ

Bufferbloat+FQ

## Quantifying “No Silver Bullet”

Nwk con-fig.	Avg. through-put, delay	Power
Bufferbloat+FQ	7.47 Mbps, 62165 ms	0.12 $Mbit/s^2$
CoDel+FQ	6.55 Mbps, 76.5 ms	85.6 $Mbit/s^2$

# Quantifying “No Silver Bullet”

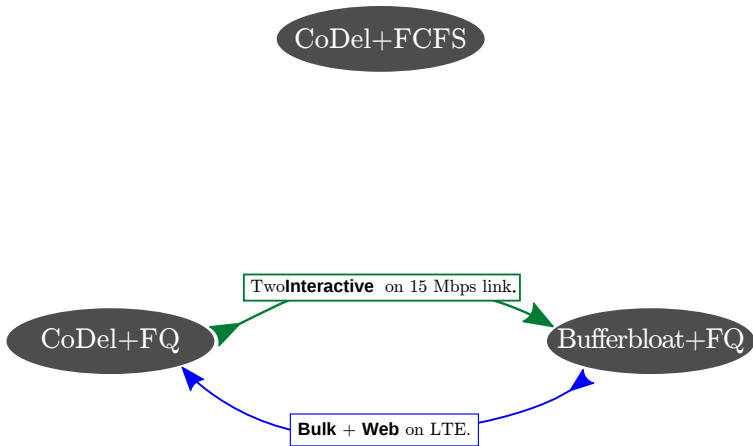


## Quantifying “No Silver Bullet”

Nwk config.	Bulk Throughput	Web Tail FCT
Bufferbloat+	11.22 Mbps	20.94 secs
CoDel+FQ	3.92 Mbps	43.72 secs



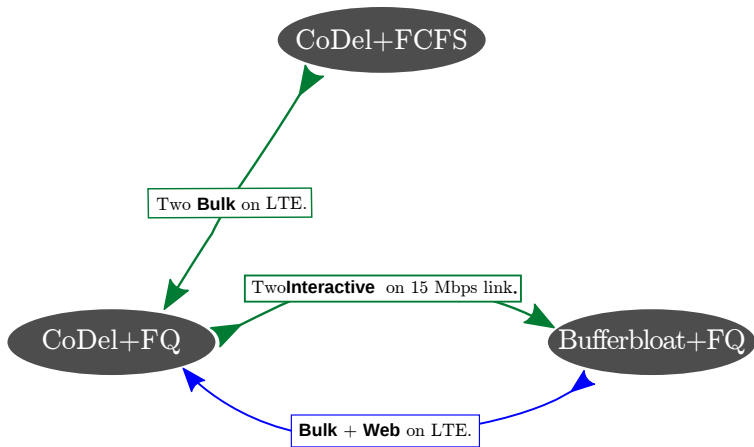
# Quantifying “No Silver Bullet”



## Quantifying “No Silver Bullet”

Nwk con-fig.	Avg. throughput
CoDel+FC	2.00 Mbps
CoDel+FQ	1.90 Mbps

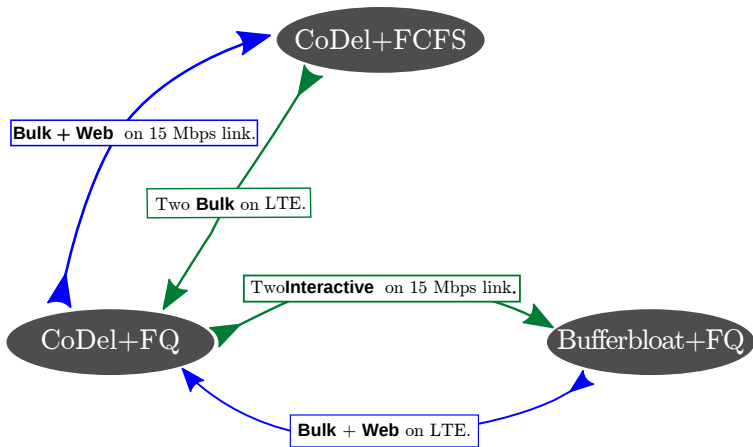
# Quantifying “No Silver Bullet”



## Quantifying “No Silver Bullet”

Nwk con-fig.	Bulk Throughput	Web Tail FCT
CoDel+FCFS	9.48 Mbps	22.25 secs
CoDel+FQ	9.48 Mbps	18.71 secs

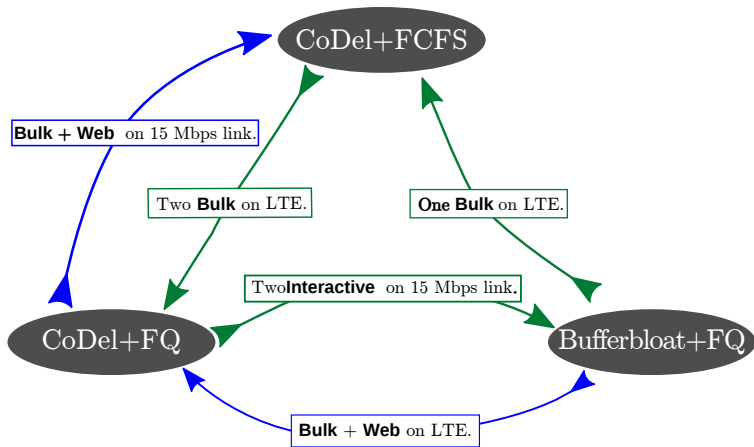
# Quantifying “No Silver Bullet”



## Quantifying “No Silver Bullet”

Nwk config.	Bulk throughput
Bufferbloat+FG	11.96 Mbps
CoDel+FCFS	4.35 Mbps

# Quantifying “No Silver Bullet”

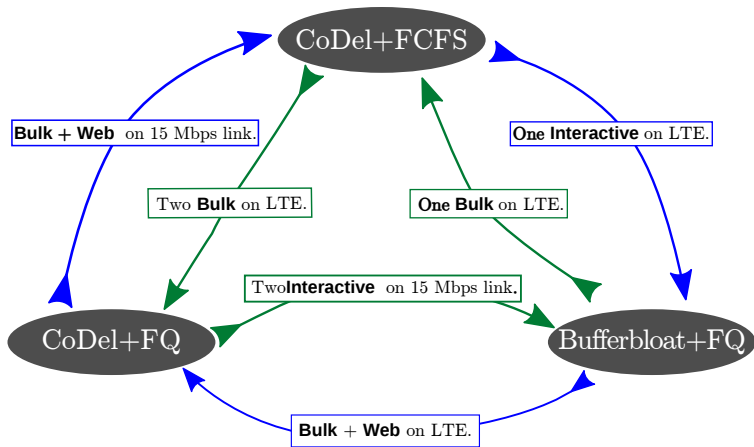


## Quantifying “No Silver Bullet”

Nwk con-fig.	Interactive through-put, delay	Power
Bufferbloat+EQ96	4.35 Mbps, 46028 ms	0.26 <i>Mbit/s<sup>2</sup></i>
CoDel+FCFS	4.35 Mbps, 83.2 ms	52.28 <i>Mbit/s<sup>2</sup></i>



# Quantifying “No Silver Bullet”



# Why is no single data plane configuration the best?

- ▶ Bufferbloat on variable-rate links helps throughput!
  - ▶ Variable-rate links have an inherent delay-throughput tradeoff
- ▶ FCFS is preferable to Fair Queuing in some cases
  - ▶ When equally aggressive flows compete, they don't need protection from each other
  - ▶ Helps reduce tail packet delay
- ▶ Fair Queuing is required in some cases
  - ▶ When competing flows aren't equally aggressive, isolation helps

## So what should the network designer do?

### Architect a flexible data plane

- ▶ Programmable queue management and scheduling
- ▶ Not just for selecting among pre-built choices, but to change behavior in the field
- ▶ Because there is no silver bullet and innovation will continue!

## Controlled flexibility: Want performance, security

- ▶ Provide interfaces only to the head and tail of queues
- ▶ Operators specify only queue-management/scheduling logic
- ▶ No access to packet payloads.

# Building such a data plane in four parts

- ▶ Hardware gadgets
  - ▶ Random number generators (RED, BLUE)
  - ▶ Binary tree of comparators (pFabric, SRPT)
- ▶ I/O interfaces
  - ▶ Drop/mark head/tail of queue
  - ▶ Interrupts for enqueue/dequeue
- ▶ State maintenance
  - ▶ Per-flow (WFQ, DRR)
  - ▶ Per-dst address (PF)
- ▶ A domain-specific instruction set
  - ▶ Expresses control flow
  - ▶ Implements new functions unavailable in hardware

## Feasibility study: CoDel

Synthesis numbers on Xilinx Kintex-7:

Resource	Usage	Fraction
Slice logic	1,256	1%
Slice logic dist.	1,975	2%
IO/GTX ports	27	2%
DSP slices	0	0%
Maximum speed	12.9 $\times 10^6$ pkts/s $\sim 10$ gbps	

- ▶ Small fraction of the FPGA's resources.
- ▶ Can be improved by pipelining or parallelizing.

## Limitations and Practical Considerations:

- ▶ Cannot express several network functions that need payloads.
- ▶ How do applications signal objectives to the network?
- ▶ Feasibility at 10G on high port-density switches.
- ▶ Energy and Area overheads.

## Related Work:

- ▶ Active Networking, e.g., ANTS
- ▶ Software Routers, e.g., Click
- ▶ Software-Defined Networking, e.g., OpenFlow



## Conclusion

- ▶ There is no silver bullet to in-network resource control because of application and network diversity
- ▶ Algorithms will continue to evolve: the data plane should help
- ▶ Directions to reproduce results:  
<http://web.mit.edu/anirudh/www/sdn-data-plane.html>