# ANALYSIS AND EVALUATION OF TCP CONGESTION ALGORITHMS

Lawrence Brakmo

Facebook

brakmo@fb.com

# INTRODUCTION

- Congestion control is a very hard problem
- People have been working at it for many decades
- Algorithm needs to utilize available bandwidth
  - Fairly
  - When many unrelated flows are competing
- Consider start up time. We do not know available bw
  - What if we did? What could we do with that information?
  - Jump to that rate? NO – there may be other flows starting up and getting the same information

## FOCUS

- I will focus on the following congestion algorithms
- *Reno* – the grandfather of all, although it has been improved
- *Cubic* – the default in Linux. Better than Reno for WAN traffic. Has hystart.
- *DCTCP* – Uses ECN markings to achieve congestion avoidance. Much better than TCP's default ECN behavior. Only good for Data Centers
- *BBR* – The new player in town. Still lots of questions about it.
- *NV* – A follow up to Vegas (my babies). Only tuned for Data Centers using TCP-BPF to set baseRTT to 80us
- *TCP-BPF* – Cubic using TCP-BPF to clamp cwnd. Only for D. C.

- Then in the last couple of days I added the following for WAN tests
- *BIC*
- *Yeah*
- *HighSpeed*
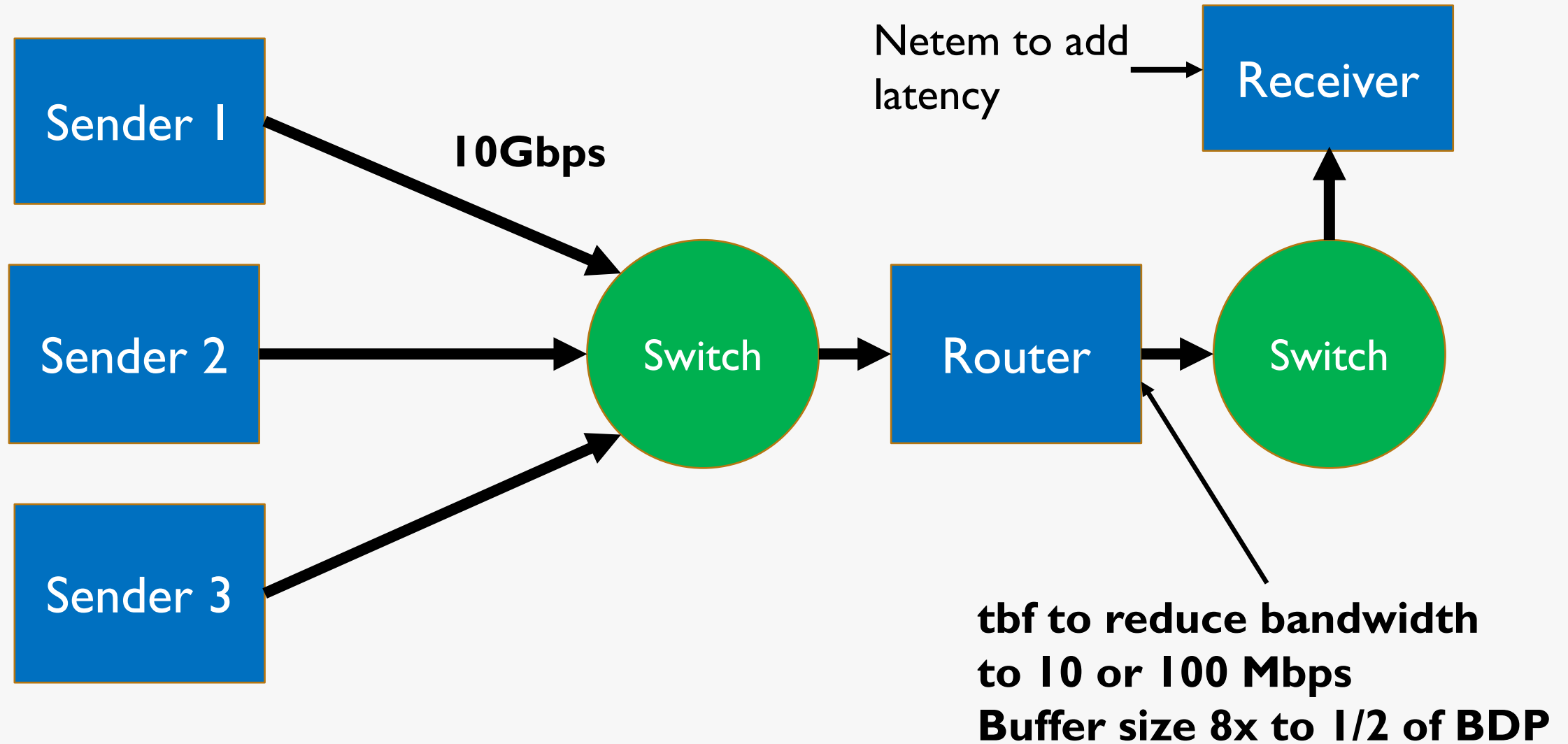- *H-TCP*
- *Westwood*

# CONGESTION VS. AVOIDANCE

- *Reno* and *Cubic* do not avoid congestion. On the contrary, they periodically create congestion and losses. It is the only way they can know they have reached full bandwidth use.

- *DCTCP, BBR\** and *NV* do congestion avoidance. They detect, or try to, congestion before losses occur. And in many cases they can keep buffers quite small improving latency

- No losses => better high percentile latencies.

EXPERIMENTAL SETUP FOR 10 AND 100 MPBS TESTS

Sender 1

Sender 2

Sender 3

10Gbps

Switch

Router

Switch

Netem to add latency

Receiver

tbf to reduce bandwidth to 10 or 100 Mbps
Buffer size 8x to 1/2 of BDP

- Scenarios
  - LAN with 20us RTT, 10 Gbps - servers in same rack.
  - Fast WAN with 10ms RTT, 10 Gbps
  - WAN with 40ms RTT, 10 and 100 Mbps
- Tests
  - *Fairness & Stability* - consists of 2 or 3 stream flow tests (each from a different server) to one receiver
  - *Size Fairness* – consists of a combination of streaming, 1MB and 10KB RPCs (8MB and 1MB for 10G-10ms scenario)

# EXPERIMENTAL SETUP

- Netesto is used to run the experiments, collect the data and create graphs and tables

  - Graphs of goodput, cwnd, RTTs, minRTTs, retransmissions

  - Tables with all the details (Goodputs, RTTs, cwnd, latencies, retransmissions, etc.

- Used Linux kernel 4.14.0-rc5

- Used mq and fq_codel queuing disciplines.

- For DCTCP and NV switch has 2 queues, one for DCTCP with ECN enable or NV, one for everything else

# RESULTS

# 10G LAN 2 FLOWS

# CUBIC

# DCTCP

# BBR



### Goodputs

| | |
|---|---|
| **Flows** | |
| — | Sum |
| — | bbr |
| — | bbr |

### cwnd

| | |
|---|---|
| **Flows** | |
| — | bbr |
| — | bbr |

# NV W/BASERTT OF 80US

# CUBIC W/TCP-BPF

# 3-FLOWS, 1-CUBIC VS. 2..

# 1-CUBIC VS. 2-BBR

# 1-CUBIC VS. 2-DCTCP

# 1-CUBIC VS. 2-BBR WITH TCP-BPF CLAMP
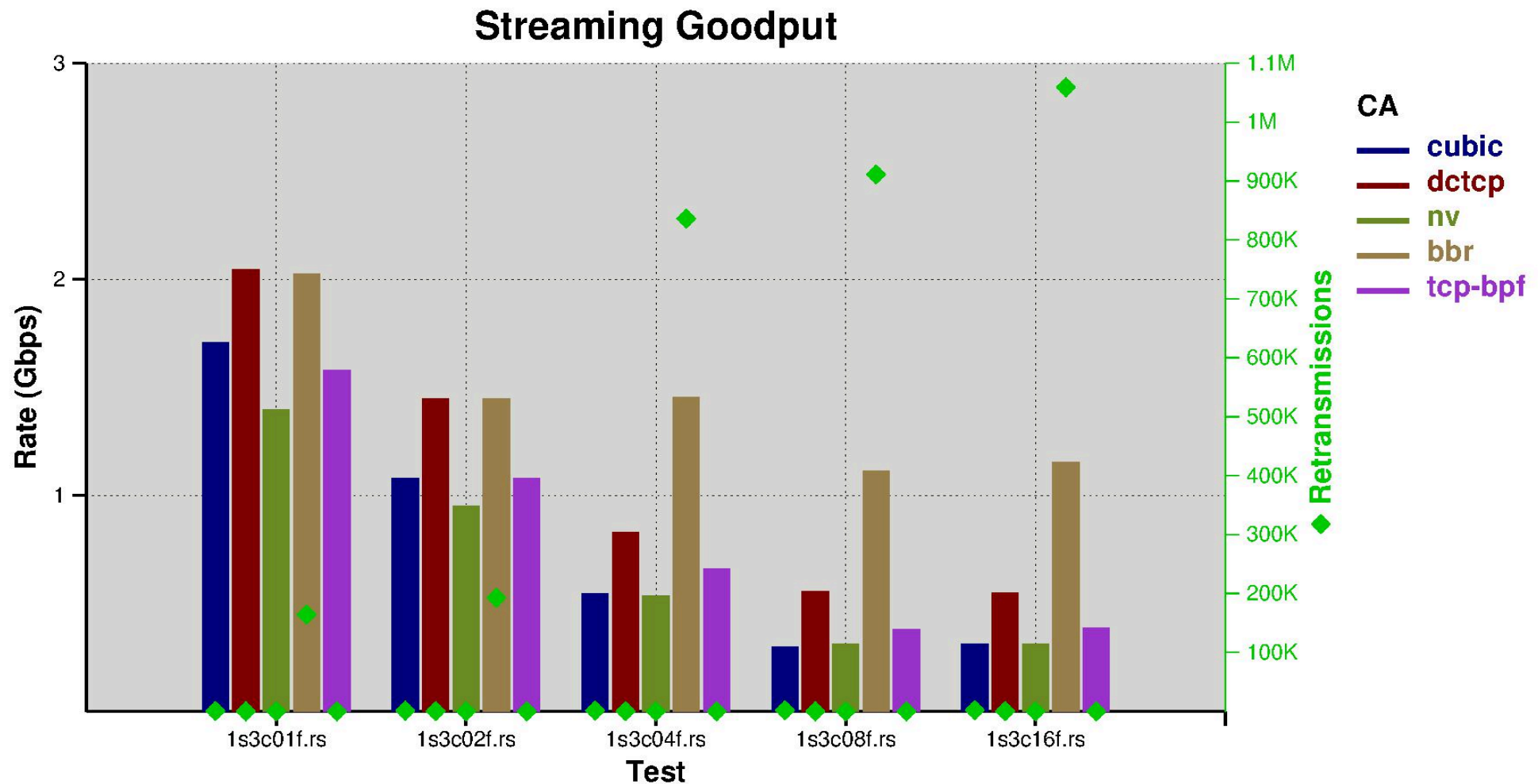
# SIZE FAIRNESS AND MANY FLOWS

Streaming Goodput

1MB RPC Goodput

10KB RPC Latencies

# 10G-10MS SCENARIOS

2 and 3 Flow Aggregate Goodput and Retransmissions

Goodputs

# 3-FLOW BBR

## Goodputs
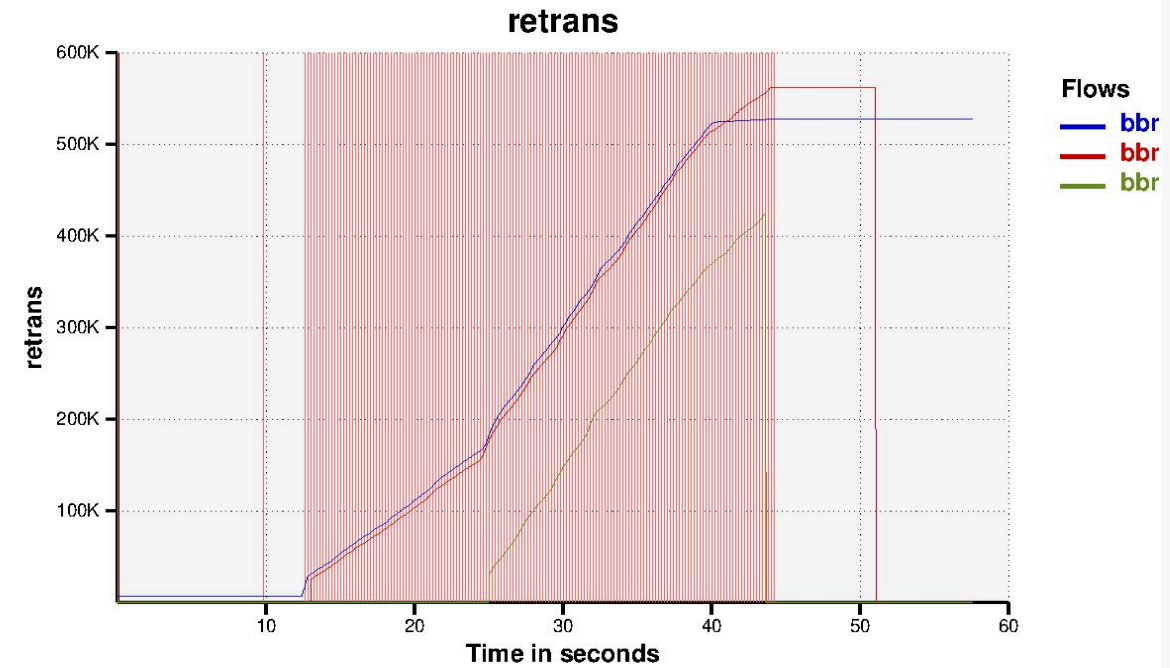
Flows
— Sum
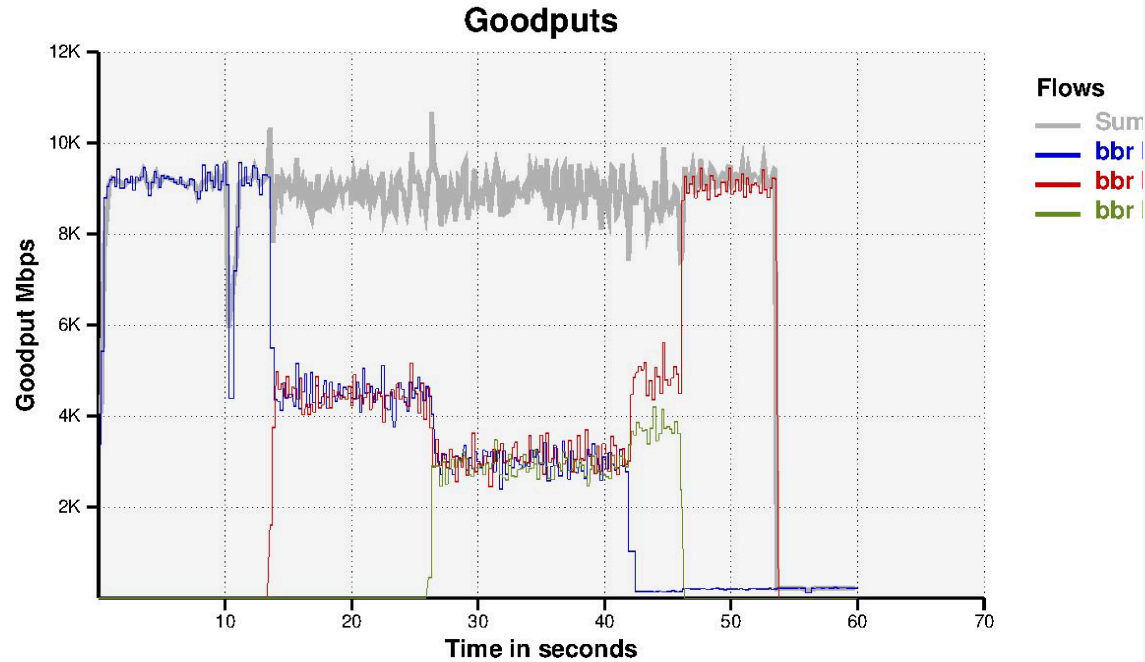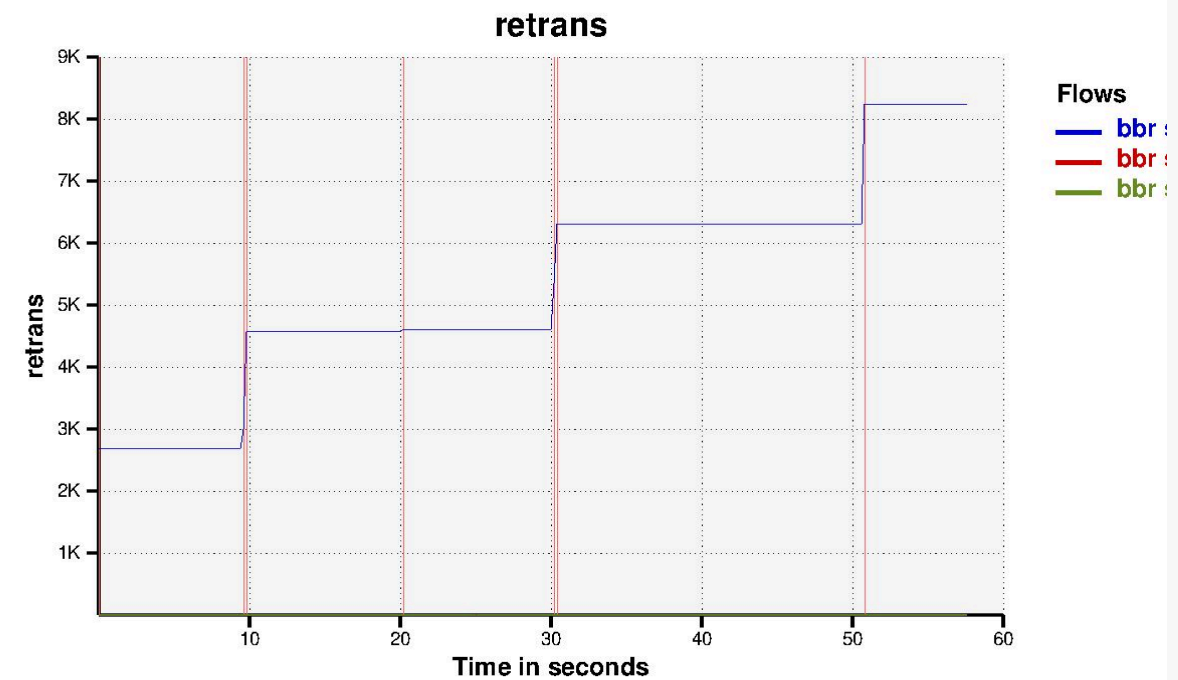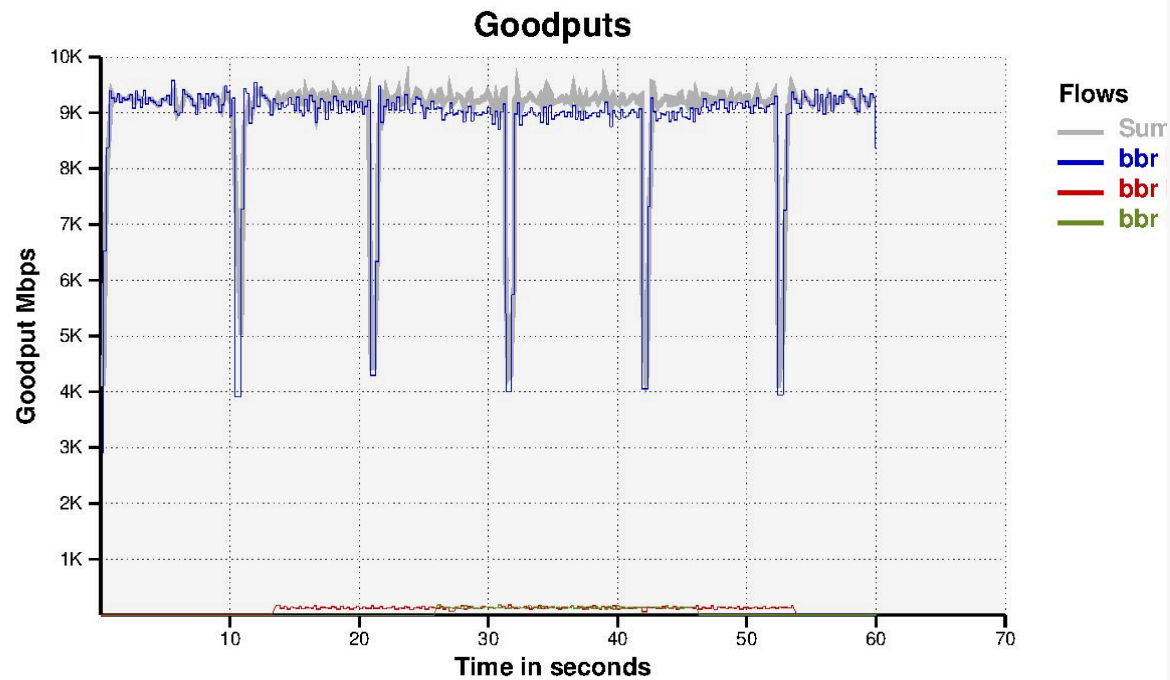— bic
— bic
— bic

## retrans

Flows
— bic s
— bic s
— bic s

# BBR FLOW COLLAPSE

# 3 FLOW BBR BAD COLLAPSE
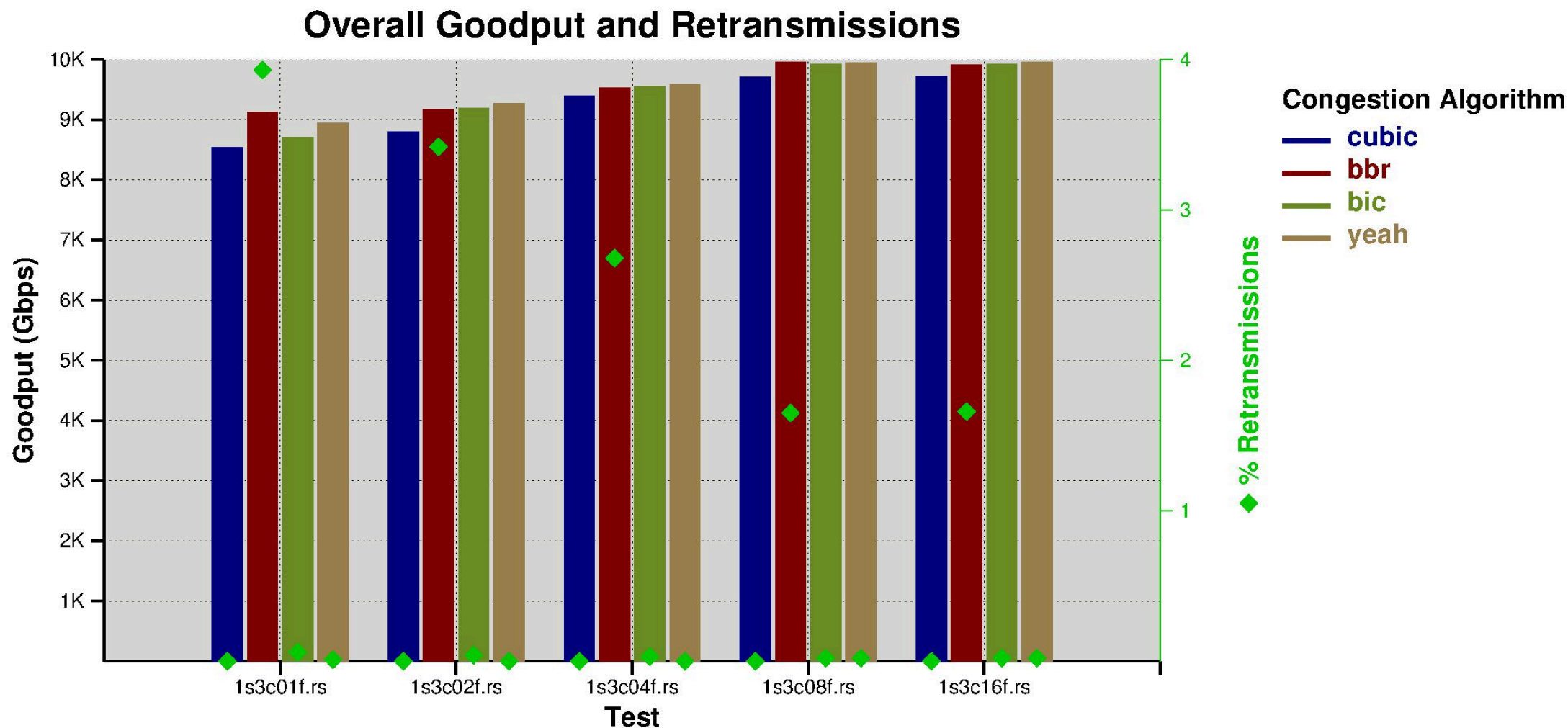
# BBR COLLAPSES

- Collapses seen in 10% of 2 flow BBR tests
- Collapses seen in 20% of 3 flow BBR tests


- It is possible that some netem is causing collapse, but
  - Other people have seen it without using netem
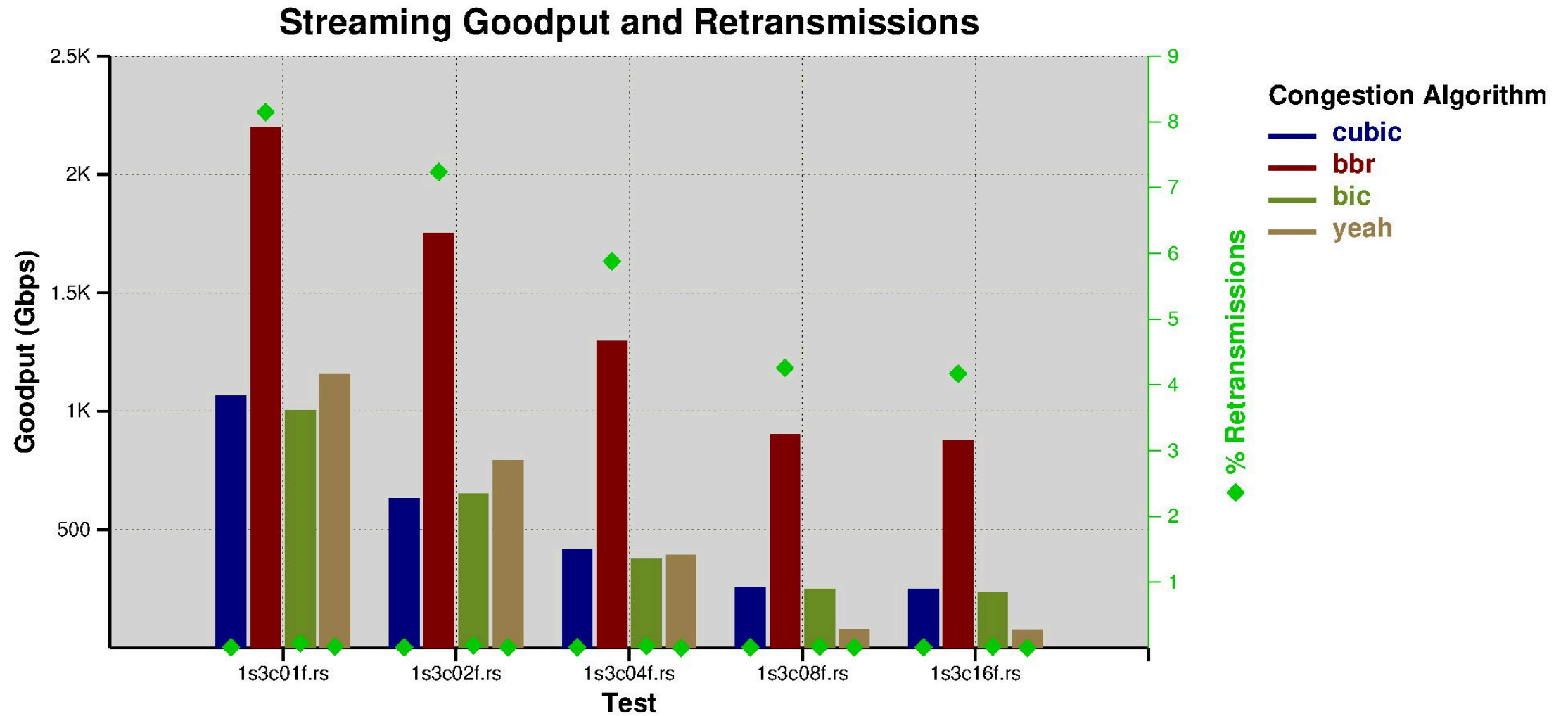  - It should be able to recover

# FAIRNESS AGAINST CUBIC

- Cubic looses against BIC and BBR

- Yeah losses against Cubic

- Cubic and Reno even

# SIZE FAIRNESS AND MANY FLOWS
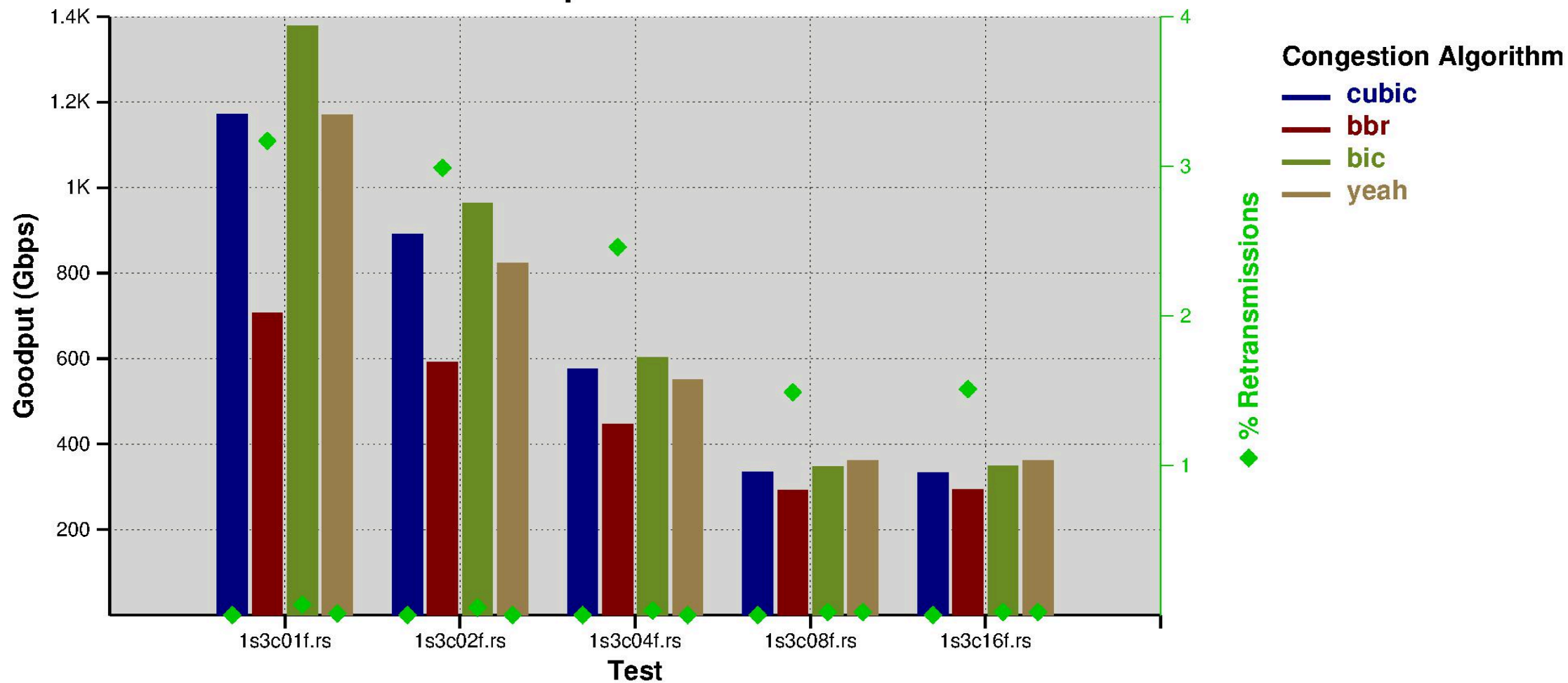
Overall Goodput and Retransmissions
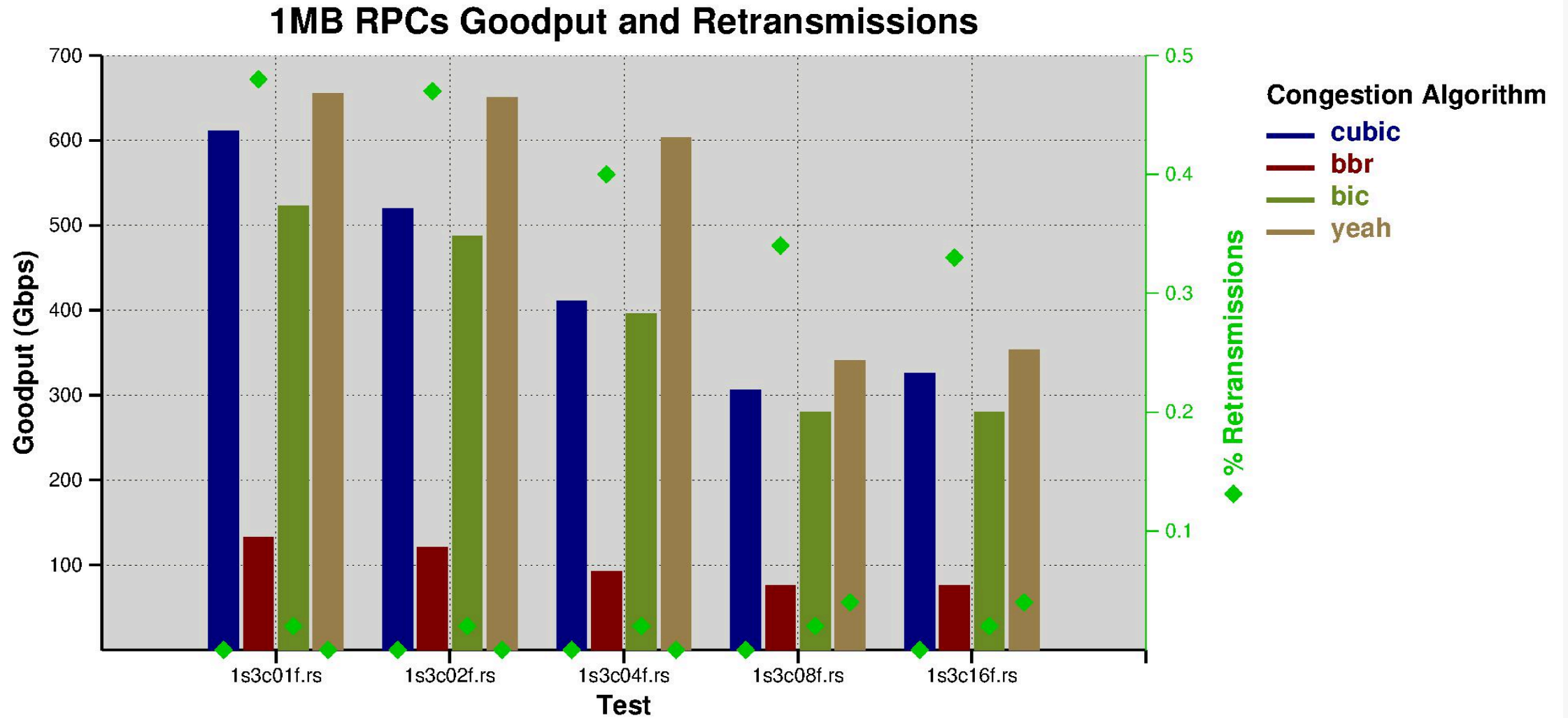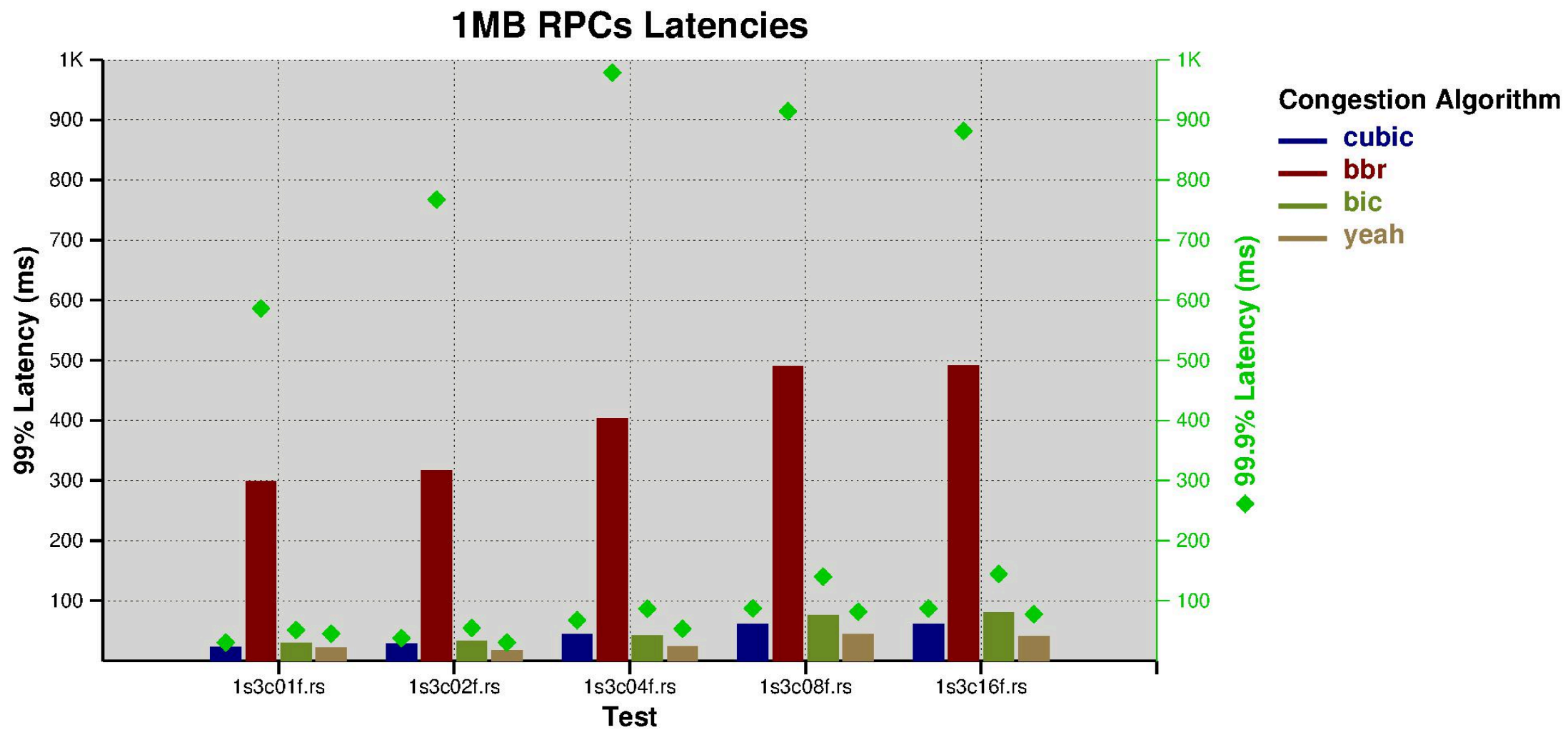
Streaming Goodput and Retransmissions

**8MB RPCs Goodput and Retransmissions**

**1MB RPCs Goodput and Retransmissions**

1MB RPCs Latencies
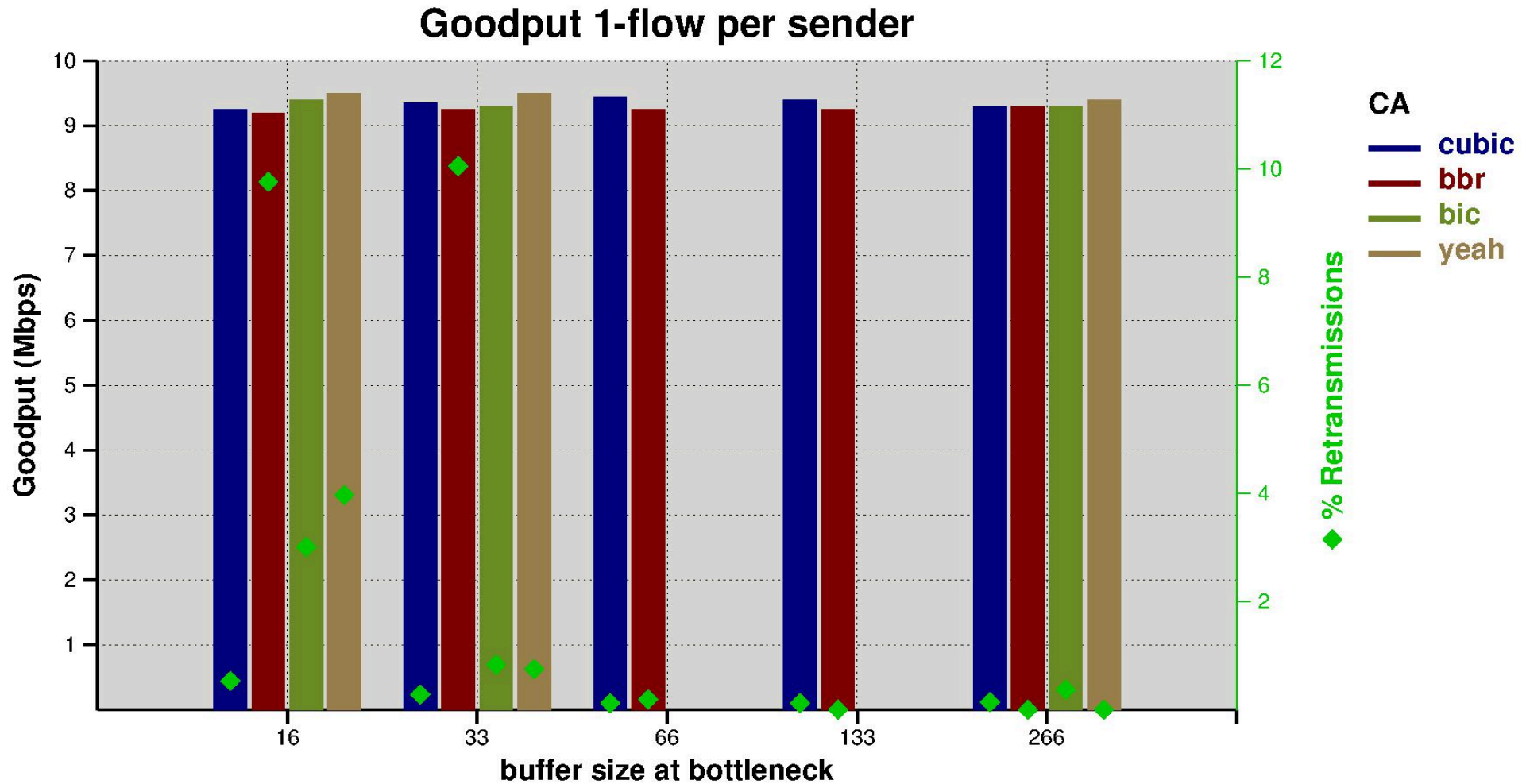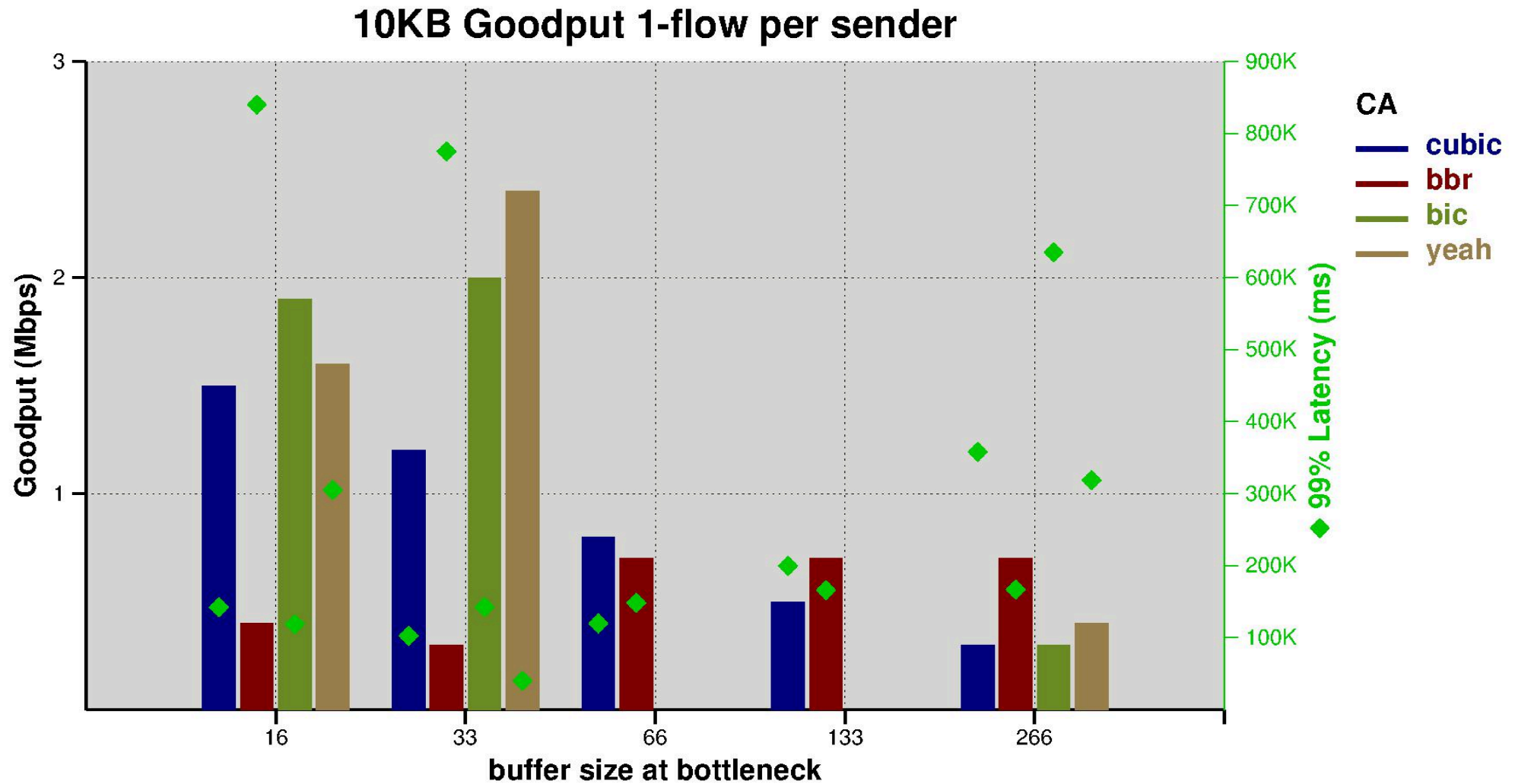
- No CA is perfect
- Yeah suffers against Cubic
- BBR and BIC hurt Cubic
- BBR is good at using available bandwidth
- BBR does well when it is the only flow
- BBR hurts itself
- BBR has a lot of retransmissions
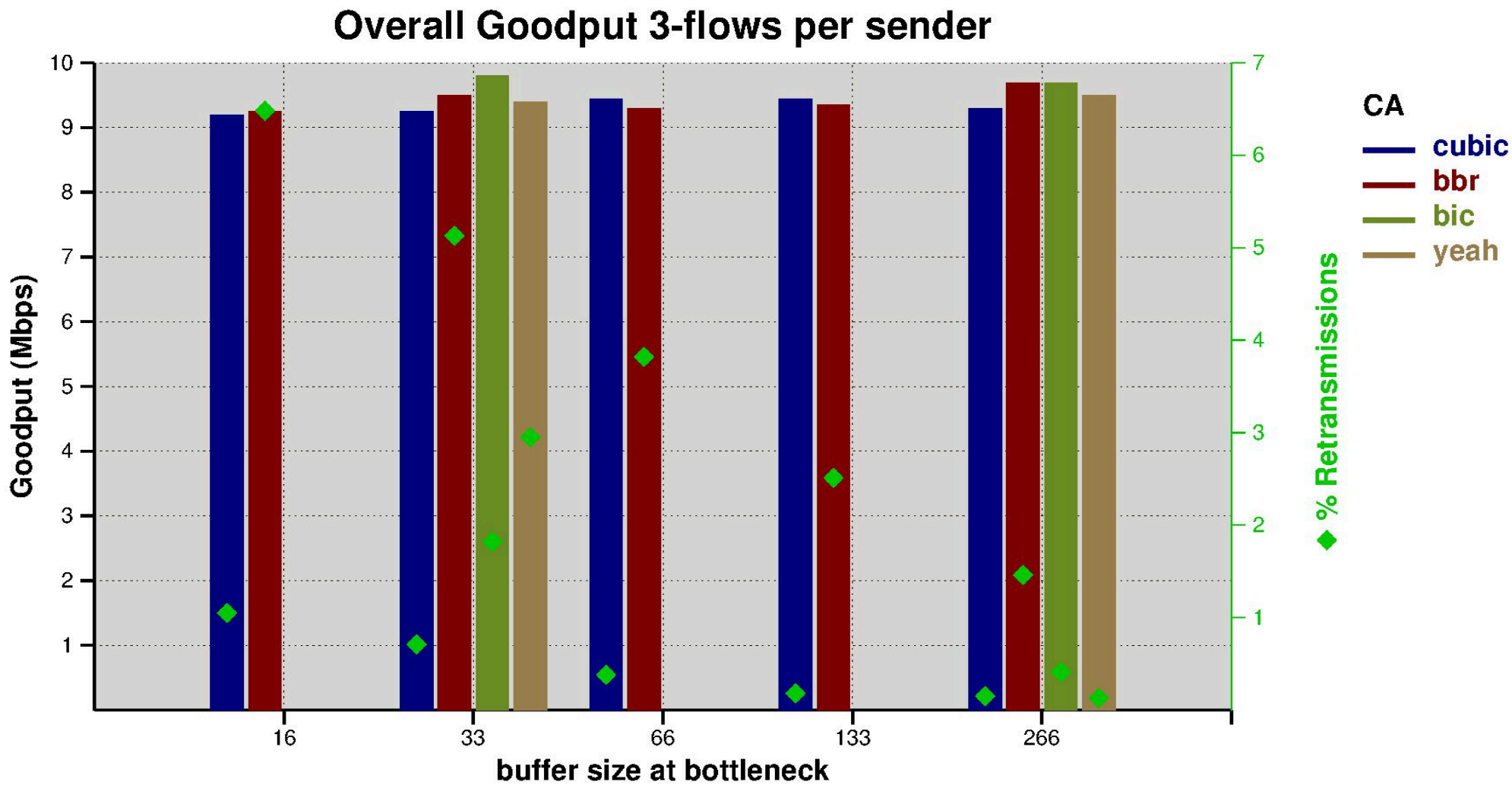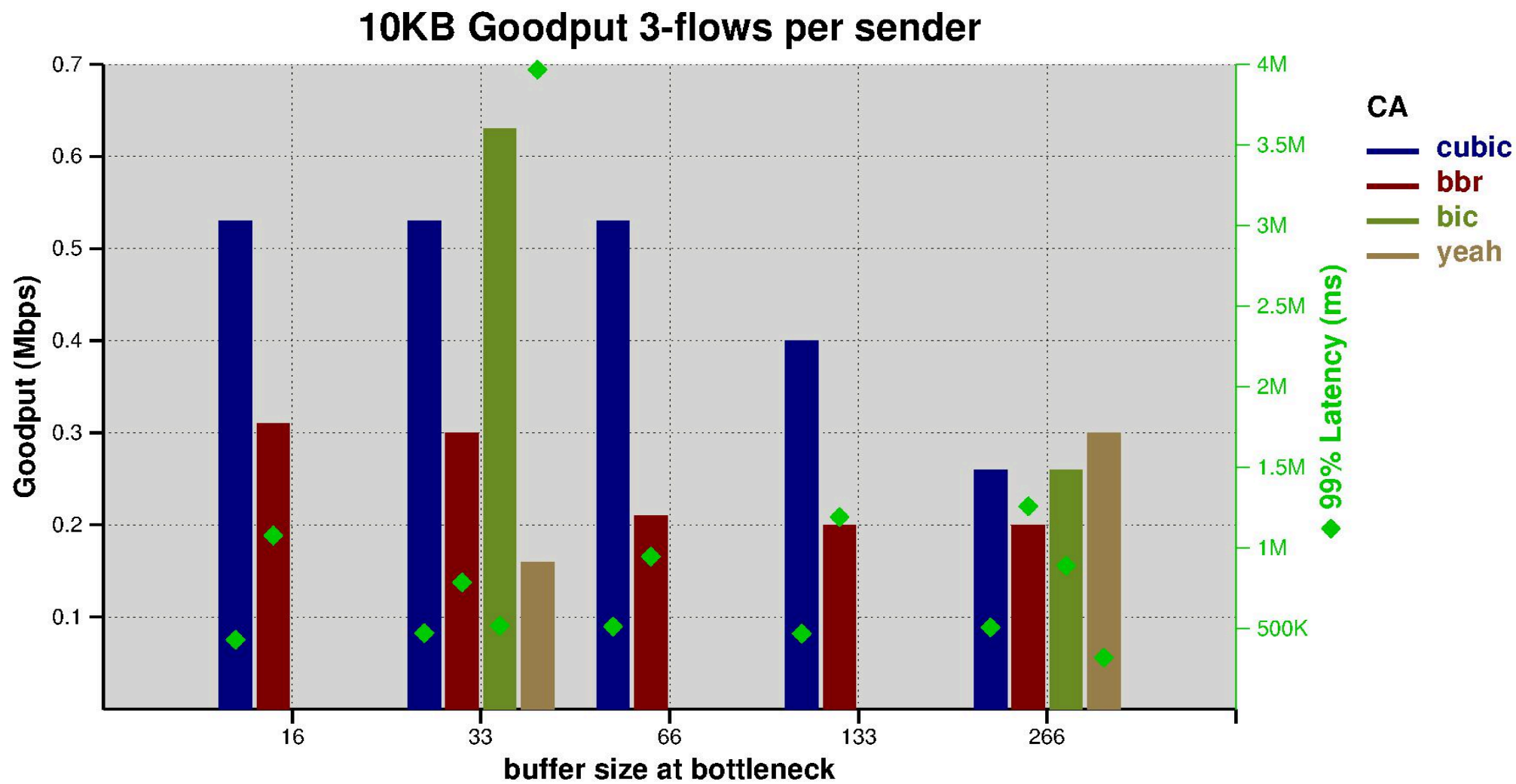
# 40MS RTT, 10 MBITS/S

Goodput 1-flow per sender

10KB Goodput 1-flow per sender

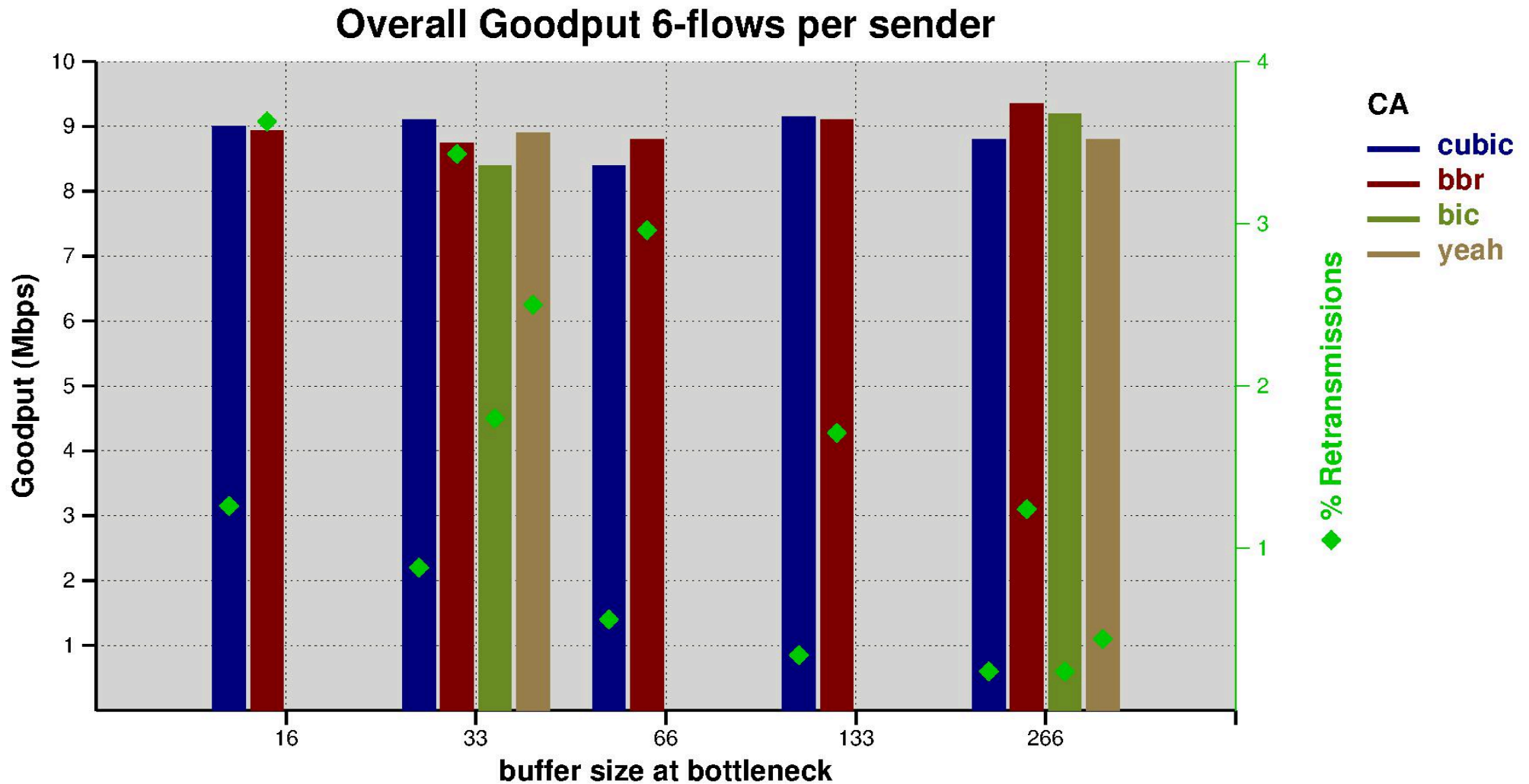Overall Goodput 3-flows per sender

10KB Goodput 3-flows per sender

Overall Goodput 6-flows per sender

10KB Goodput 6-flows per sender