

BBR Congestion Control: An Update

Neal Cardwell, Yuchung Cheng,
C. Stephen Gunn, Soheil Hassas Yeganeh,
Van Jacobson

<https://groups.google.com/d/forum/bbr-dev>

Outline

- Review of BBR [also see [IETF 97 ICCRG BBR slides](#)]:
 - Model
 - Algorithm
 - Behavior
- Deployment at Google: fully deployed for Google TCP WAN traffic
 - Google.com
 - YouTube
 - Internal WANs
- Active and upcoming work

Problems with loss-based congestion control

2011: many reported excessive buffering and delays on the Internet (a.k.a. bufferbloat)

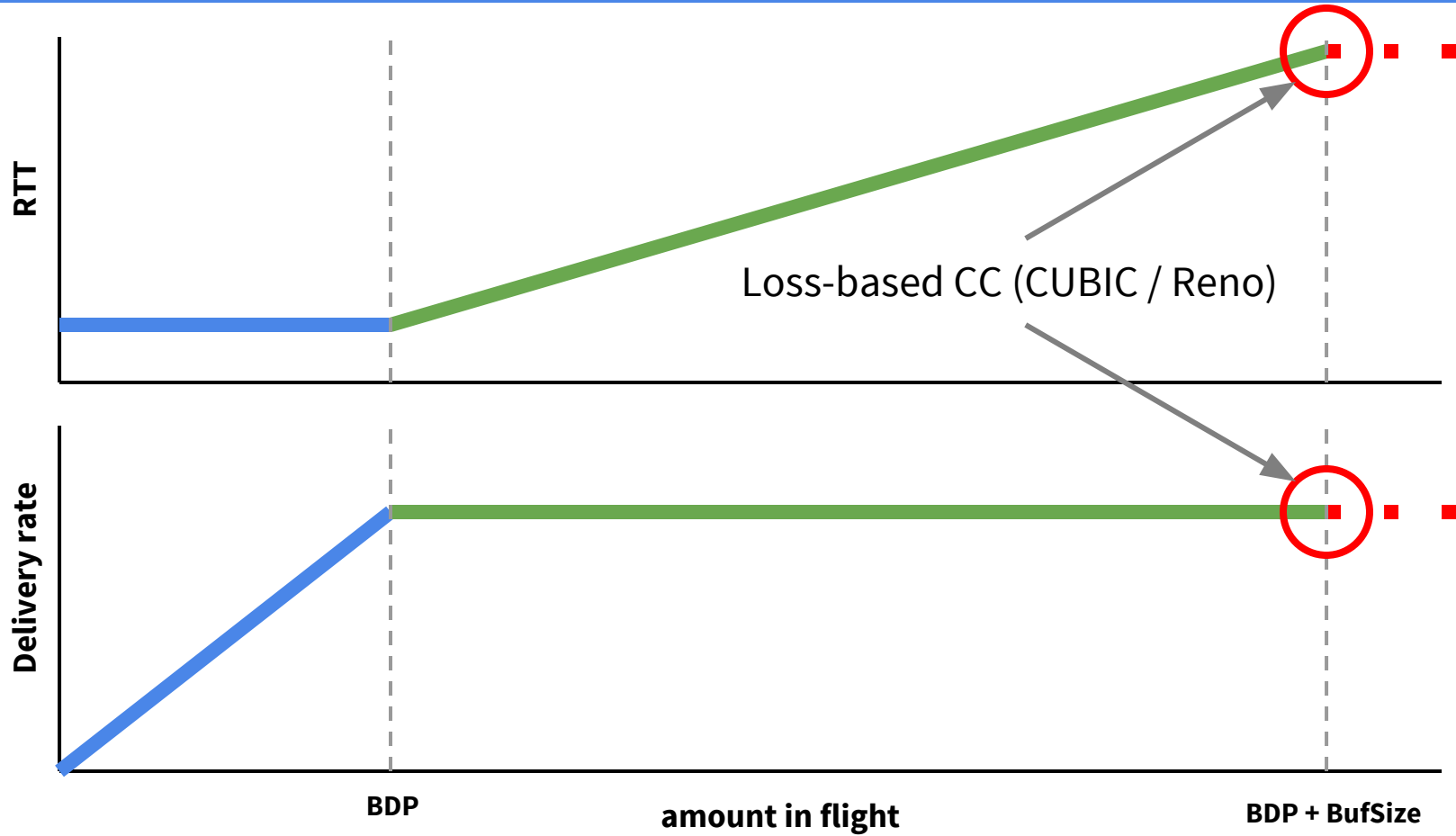
2012: single-connection HTTP/2 was much slower than multi-conn HTTP/1 on lossy links

2013: poor TCP throughput on WANs w/ commodity shallow-buffer switches

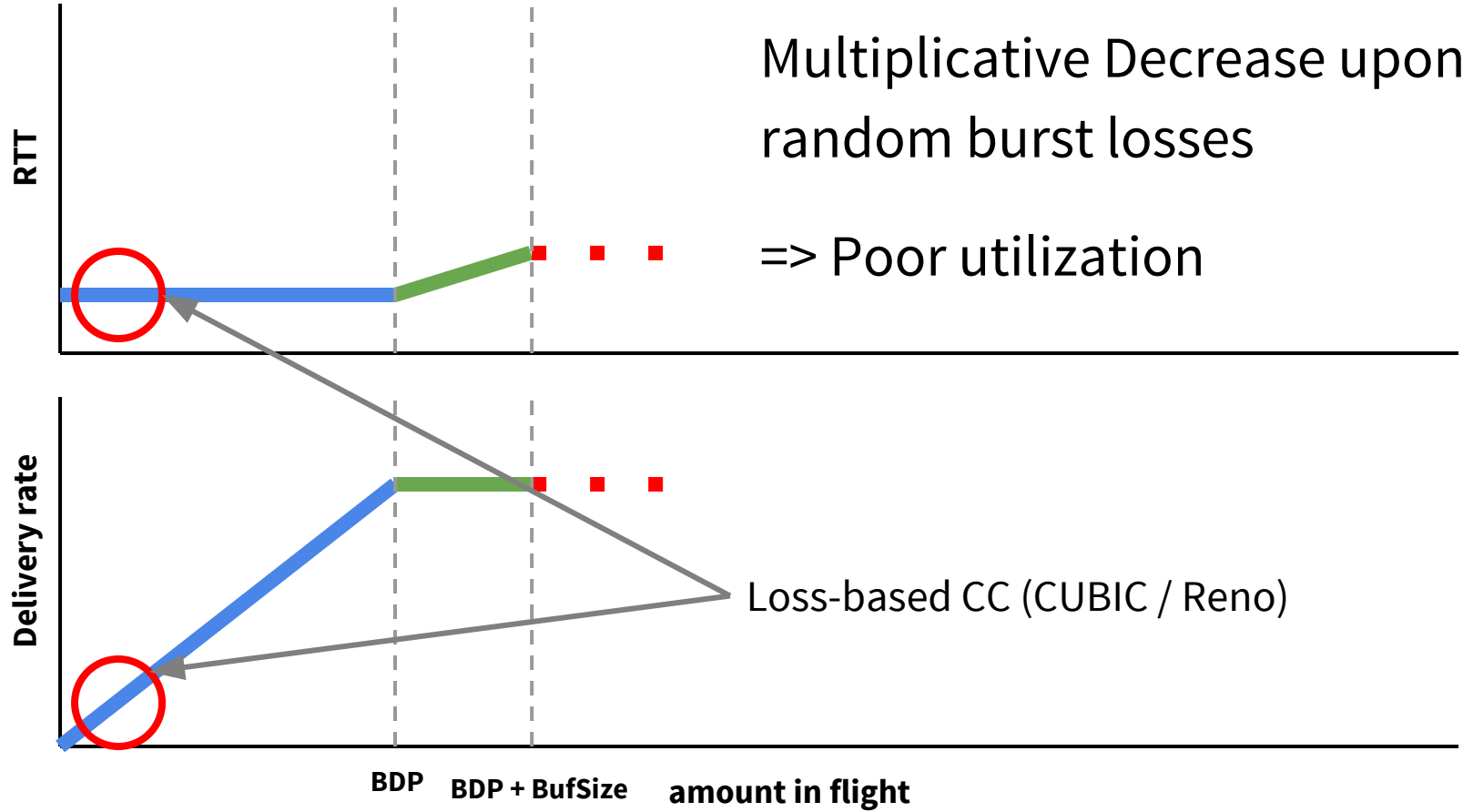
Culprit: loss-based congestion control (Reno, then CUBIC)

- Packet loss alone is **not** a good proxy to detect congestion
- Loss-based CC is overly sensitive to losses that come **before** congestion
 - 10Gbps over 100ms RTT needs $<0.000003\%$ packet loss
 - 1% loss over 100ms RTT gets 3Mbps
- Loss-based CC bloats buffers if loss comes **after** congestion

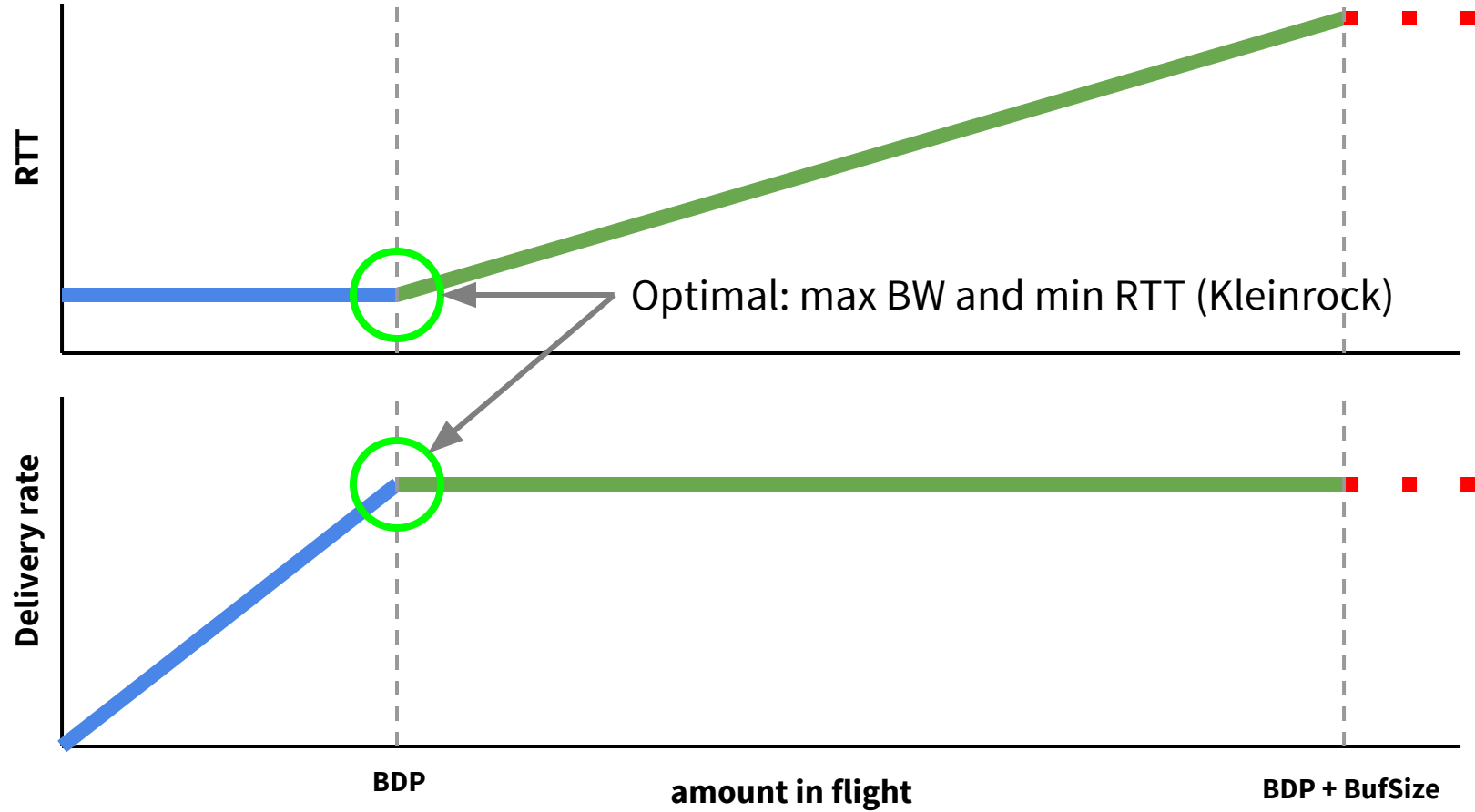
Loss-based congestion control in deep buffers



Loss-based congestion control in shallow buffers



Optimal operating point

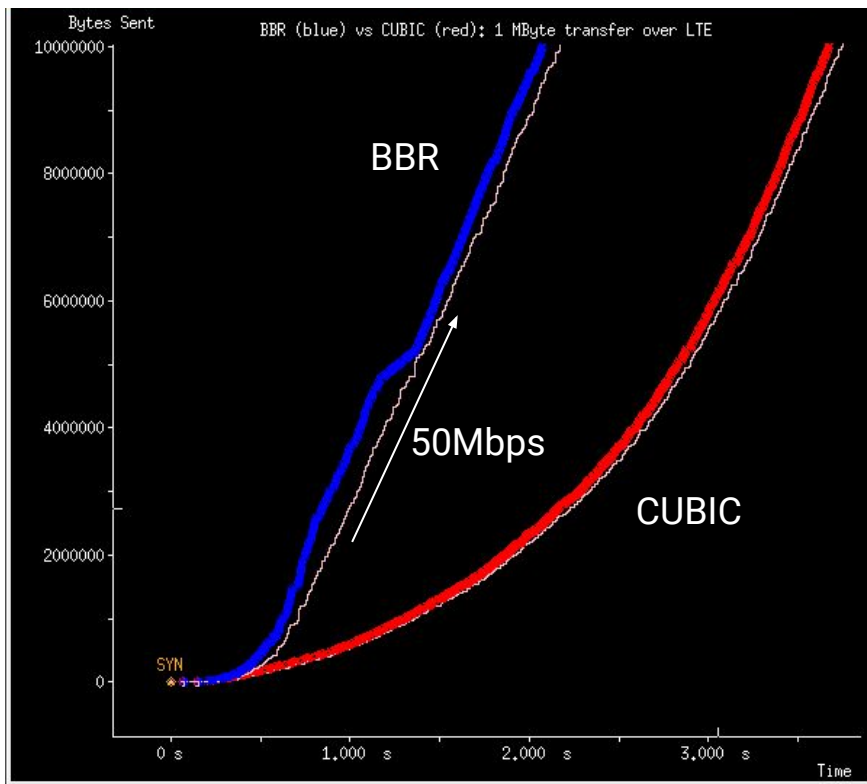


BBR = **B**ottleneck **B**andwidth and **R**ound-trip propagation time

- Model network path
 - **Dynamically** estimate windowed max BW and min RTT on each ACK
- Control sending based on the model, to...
 - **Sequentially** probe max BW and min RTT, to feed the model samples
 - Pace near estimated BW
 - Vary pacing rate to keep inflight near BDP
- Seek high throughput with a small queue
 - Approaches maximum available throughput for random losses up to 15%
 - Maintains small queue independent of buffer depth

BBR: faster for short flows, too

	Cubic (Hystart)	BBR
Initial rate	10 packets / RTT	
Acceleration	2x per round trip	
Exit acceleration	A packet loss or significant RTT increase	Delivery rate plateaus



BBR and Cubic time series overlaid. BBR downloads 1MB 44% faster than Cubic. Trials produced over LTE on Neal's phone in New York

BBR is deployed for WAN TCP traffic at Google

vs CUBIC, BBR yields:

- 2% lower search latency on google.com
- 13% larger Mean Time Between Rebuffers on YouTube
- 32% lower RTT on YouTube
- Loss rate increased from 1% to 2%

Lessons from deploying a new C.C. in the real Internet

- Cellular or Wi-Fi gateways adjust link rate based on the backlog
 - BBR needs to deliberately keep some queue
- Delay alone is an extremely noisy signal
 - e.g. waiting for a slot on a shared medium (radio or cable)
- NICs or middleboxes hold up or decimate ACK packets
 - e.g. one TCP ACK for up to +200 packets
- Token-bucket traffic policers allow bursts, then drop
 - Common in developing regions and mobile carriers
- Middle-boxes modify TCP receive window to throttle sender
- ...

Lessons from implementing a new C.C.

- RFC 793 is not a sufficient roadmap for how to implement a CC protocol
- Need a framework providing much more detailed accounting of what's going on
 - Per-connection totals: SACKed, lost, retransmitted, in-flight
 - Attaching detailed per-TSO-chunk accounting info
 - (Re)transmit/SACK state => enables SACK scoreboard => RFC6675 recovery
 - Transmission times => enables time-based loss reasoning => RACK
 - Delivery rate state => enables delivery-based CC => BBR

An opportunity: leveraging BBR's path model

- Linux TCP rate sampling code is available to Linux apps using any C.C.
 - `getsockopt(TCP_INFO)`
- BBR exports its bandwidth (bw) and two-way propagation delay (min_rtt) estimates
 - `getsockopt(TCP_CC_INFO)`
 - Better than cwnd/rtt
- Possible applications:
 - Exporting BW to video apps to pick the best video format
 - ...

Improving BBR

BBR can be even better:

- Smaller queues: lower delays, less loss, more fair with Reno/CUBIC
 - Potential: cut RTT and loss rate in half for bulk flows
- Higher throughput with wifi/cellular/DOCSIS
 - Potential: 10-20% higher throughput for some paths
- Lower tail latency by adapting magnitude of PROBE_RTT
 - Potential: usually PROBE_RTT with $cwnd = 0.75 * BDP$ instead of $cwnd=4$

End goal: improve BBR to enable it to be the default congestion control for the Internet

We have some ideas for tackling these challenges

We also encourage the research community to dive in and improve BBR!

Following are some open research areas, places where BBR can be improved...

Open research challenges and opportunities with BBR

Some of the areas with work (experiments) planned or in progress:

- Reducing queuing/losses on shallow-buffered networks and/or with cross-traffic:
 - Quicker detection of full pipes at startup
 - Gentler [PRR](#)-inspired packet scheduling during loss recovery
 - Refining the bandwidth estimator for competition, app-limited traffic
 - Refining cwnd provisioning for TSO quantization
 - More frequent pacing at sub-unity gain to keep inflight closer to available BDP
 - Explicit modeling of buffer space available for bandwidth probing
- Improving fairness vs. other congestion controls
- Reducing the latency impact of PROBE_RTT by adaptively scaling probing
- Explicitly modeling ACK timing, to better handle wifi/cellular/cable ACK aggregation

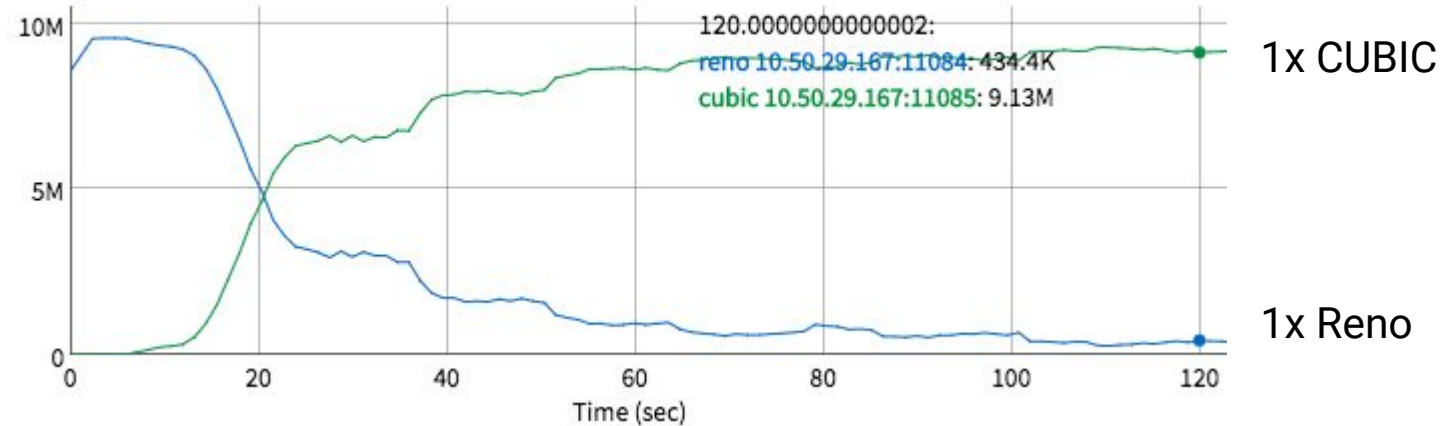
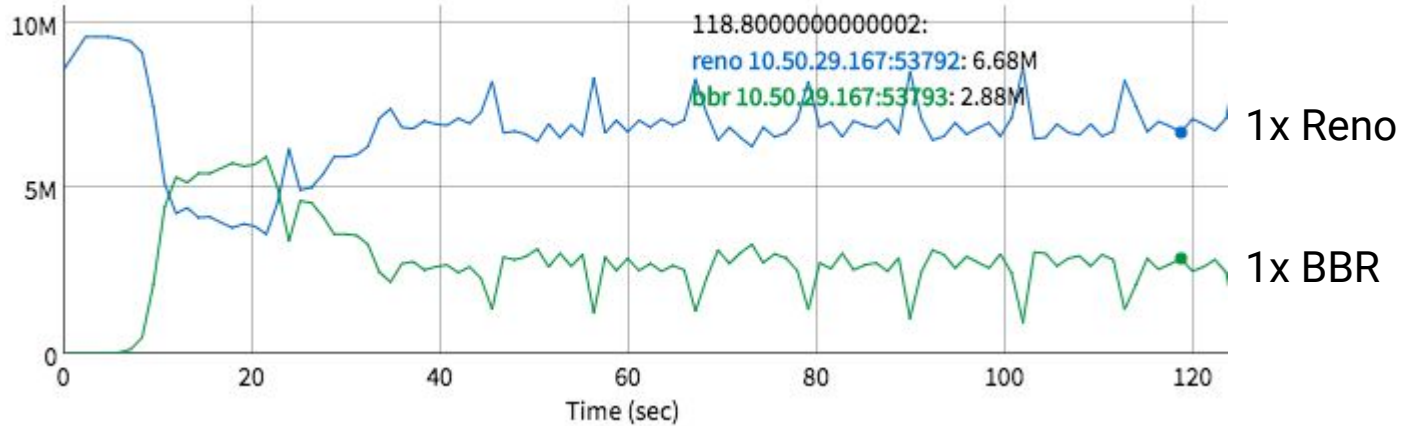
Most recent experiments

- Reducing queuing/losses on shallow-buffered networks and/or with cross-traffic:
 - Quicker detection of full pipes at startup
 - Gentler [PRR](#)-inspired packet scheduling during loss recovery
 - More frequent lower-rate pacing to keep inflight closer to available BDP

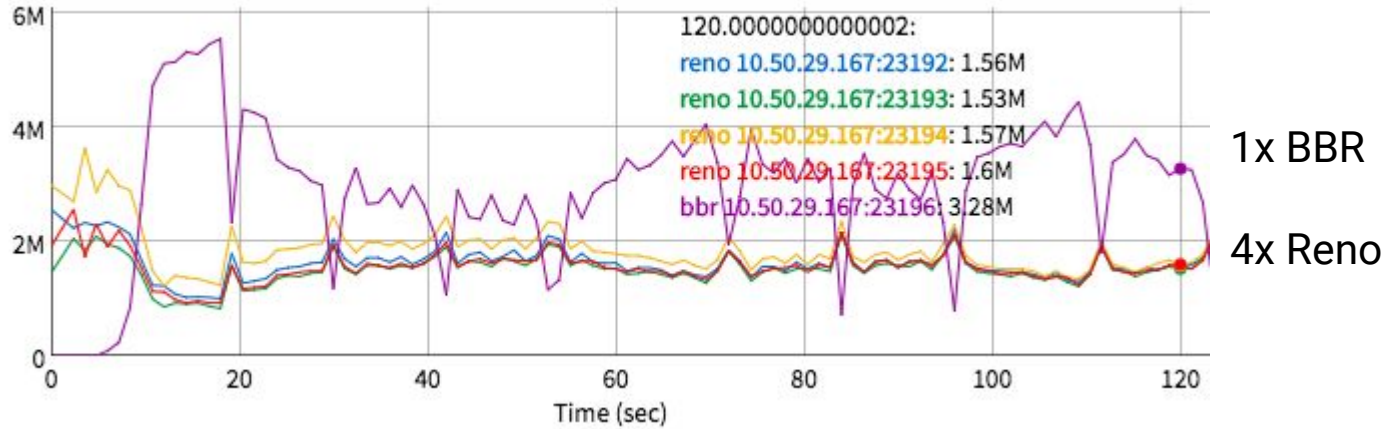
.... resulting fairness?

In **deep** buffers, BBR's fairness to Reno matches or exceeds CUBIC's fairness to Reno...

In deep buffers: BBR, CUBIC friendliness to 1x Reno



In deep buffers: BBR, CUBIC friendliness to 4x Reno

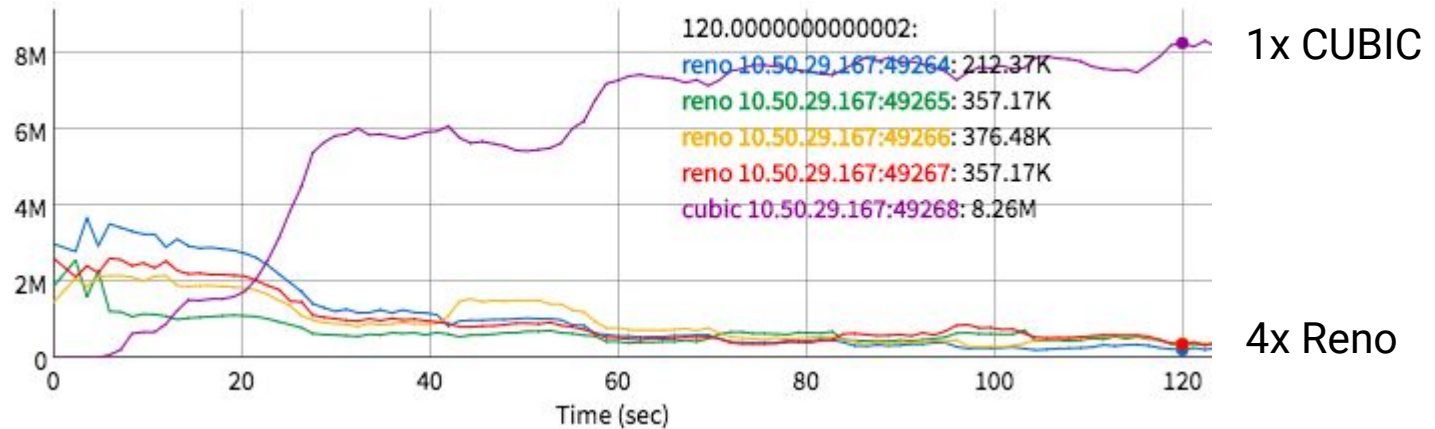


10 Mbps bw

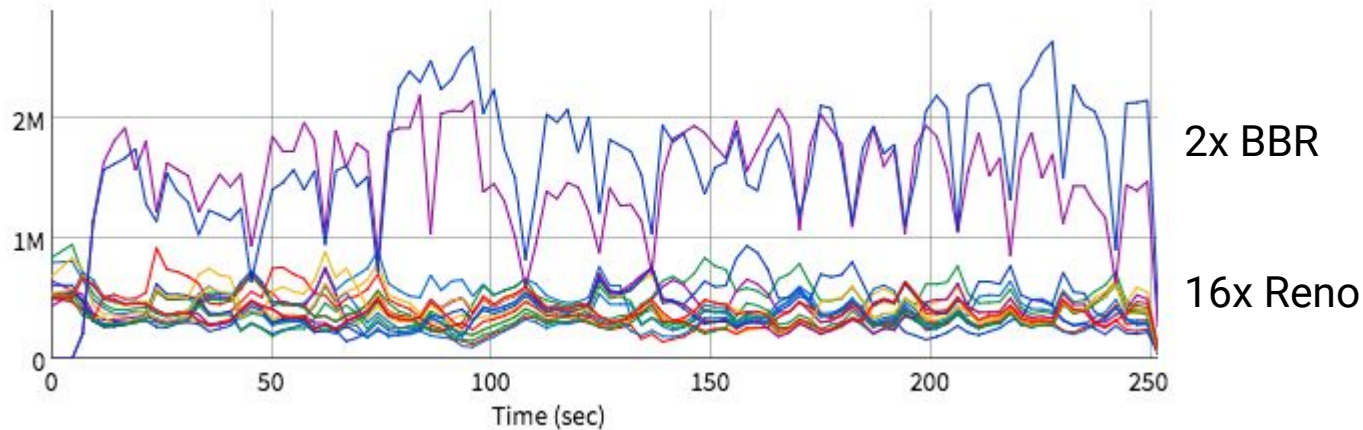
40ms RTT

1 MByte buffer

120 sec test



In deep buffers: BBR, CUBIC friendliness to 16x Reno

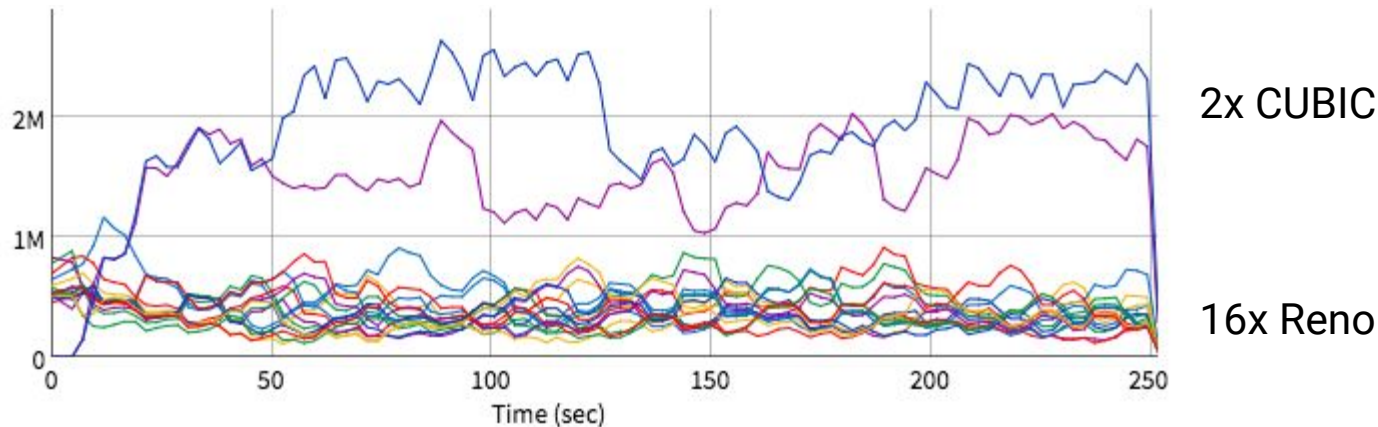


10 Mbps bw

40ms RTT

1 MByte buffer

240 sec test



Latest experiment: modeling available buffer space

Goal: How to reduce buffer pressure and improve fairness in **shallow** buffers?

What if: we try to use no more than half of flow's estimated share of the bottleneck buffer?

full_rtt: average of RTT samples in first round of loss recovery phases in last N secs

if (full_rtt)

my_buffer_target = (full_rtt - min_rtt) * bw / 2

my_max_cwnd = bw * min_rtt + my_buffer_target

Next: how to probe gently but scalably when there are no recent losses?

e.g.: my_buffer_target *= 1.25 for each second of active sending?

Looking for ways to contribute?

- More data is always useful
- And Google uses many kinds of networking traffic, but not all
- So... testing, testing, testing!
 - wifi LANs
 - AQM: experiences with Codel, PIE, DOCSIS AQM, ...
 - ...
- We love pcaps (or recipes to reproduce behaviors)
- If you're trying to use BBR in a realistic scenario, and it's not working, let us know!

Conclusion

- Loss-based C.C., a 30-year-old approach, is failing on modern fast networks
 - Packet loss signal alone is too late (big queues) or too noisy (underutilization)
- BBR uses BW & RTT (instead of a window) to model the network
 - Goal: maximize bandwidth, then minimize queue
- Deployed on Google.com, YouTube, B4/B2 (for TCP)
 - Better performance for web, video, RPC traffic
 - Open sourced in [Linux TCP](#), [QUIC](#)
- Applying BBR to [QUIC](#), FreeBSD TCP @ Netflix
- Actively working on improving the algorithm

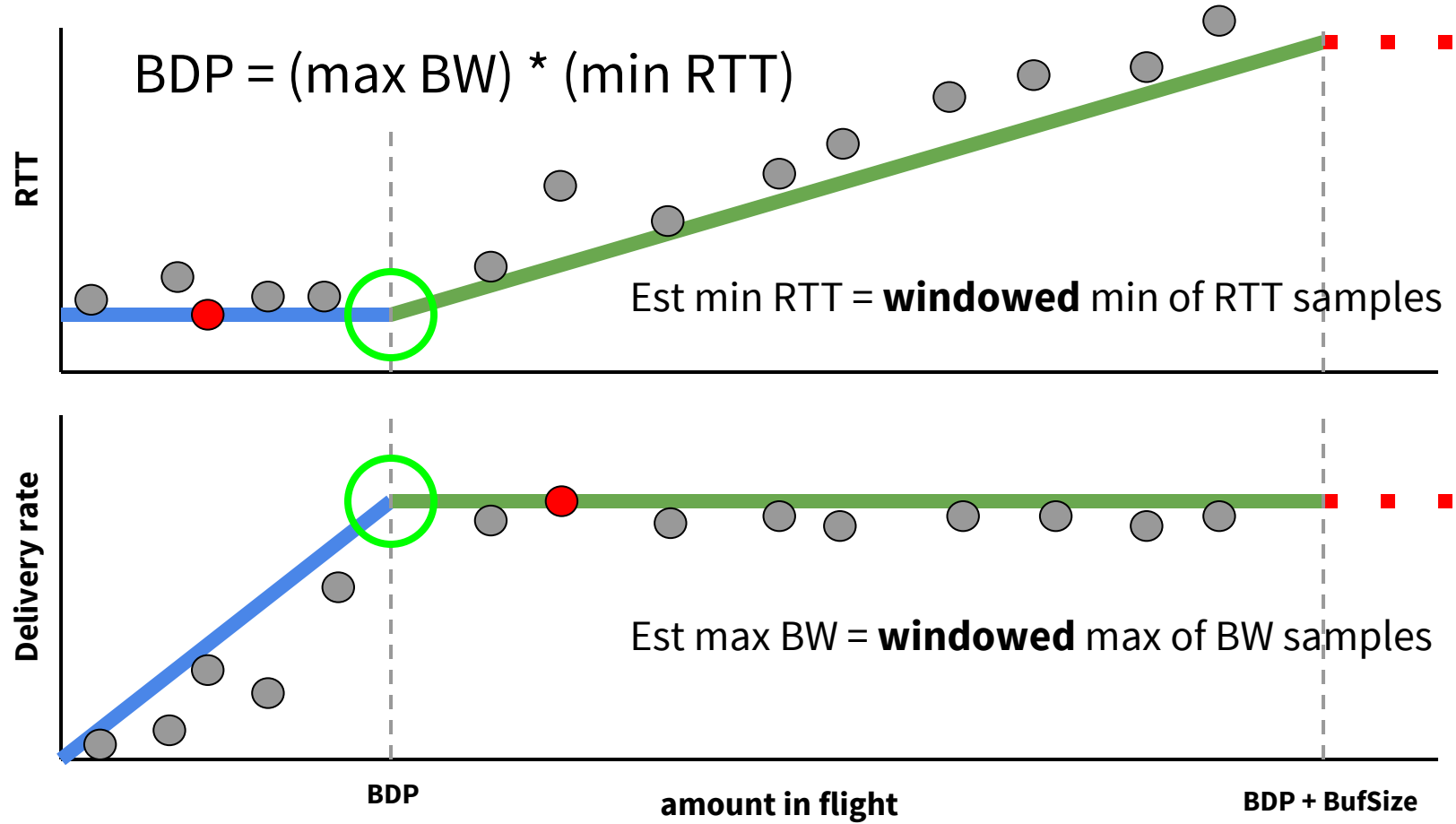
<https://groups.google.com/d/forum/bbr-dev>

research paper, code, mailing list, talks, etc.

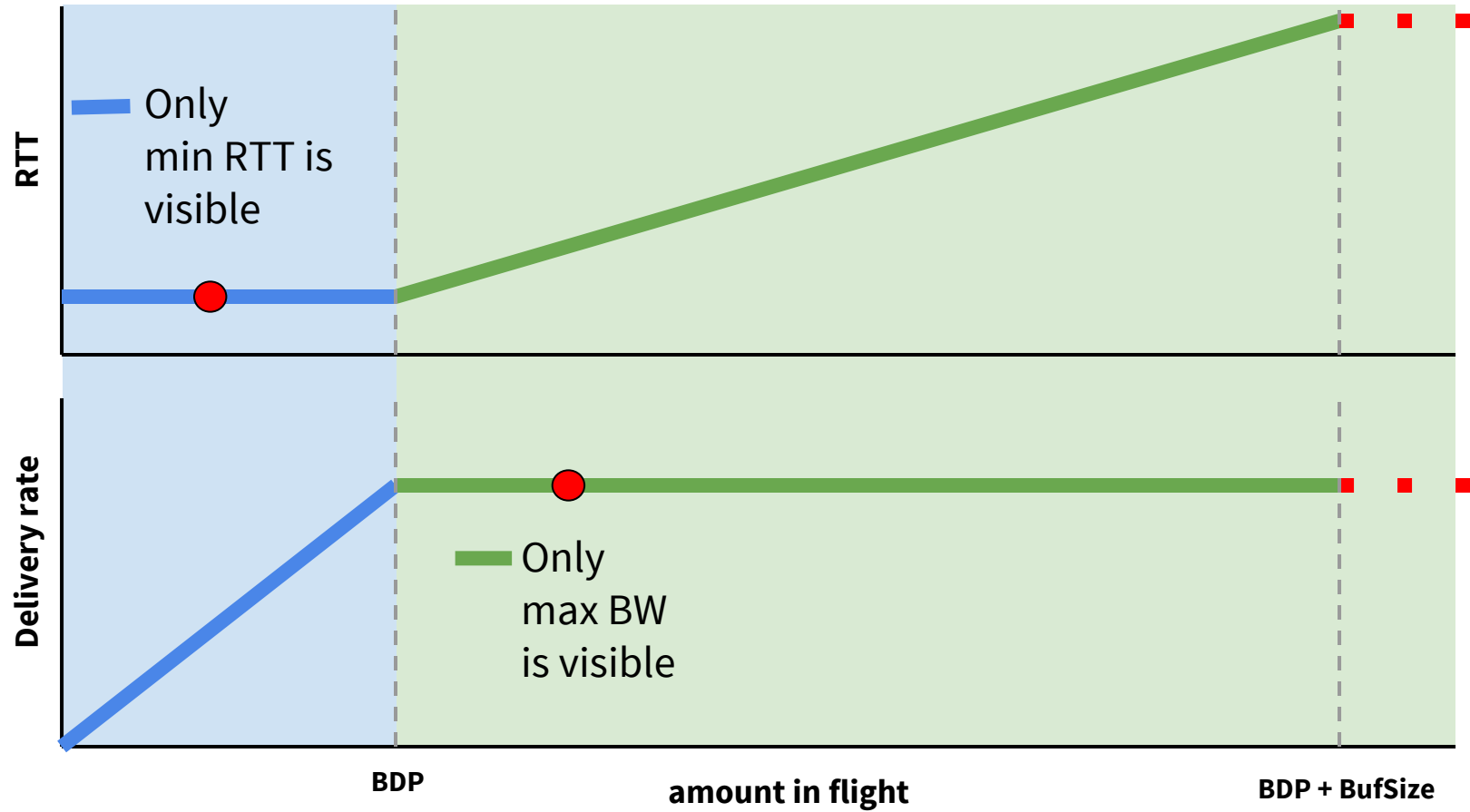
Special thanks to Eric Dumazet, Ian Swett, Jana Iyengar, Victor Vasiliev, Nandita Dukkipati, Pawel Jurczyk, Biren Roy, David Wetherall, Amin Vahdat, Leonidas Kontothanassis, and {YouTube, google.com, SRE, BWE} teams.

Backup slides...

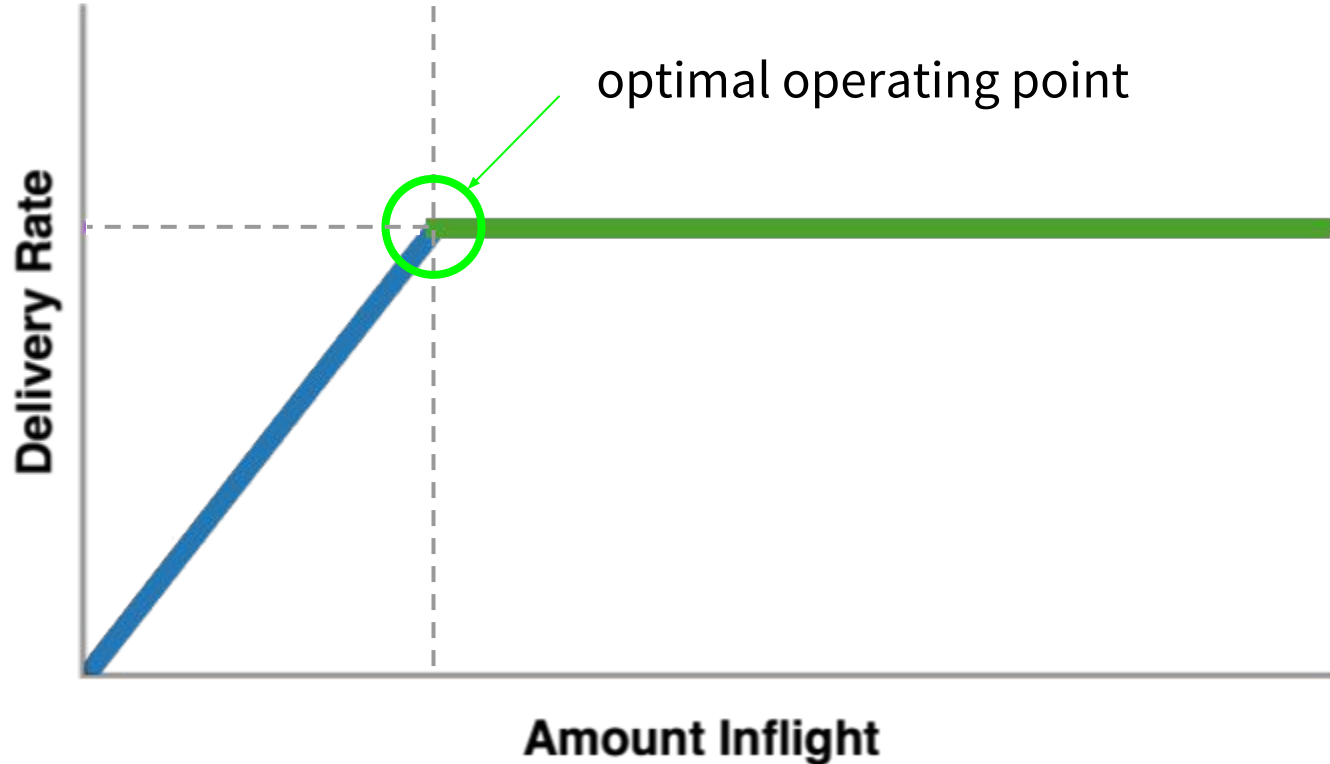
Estimating optimal point (max BW, min RTT)



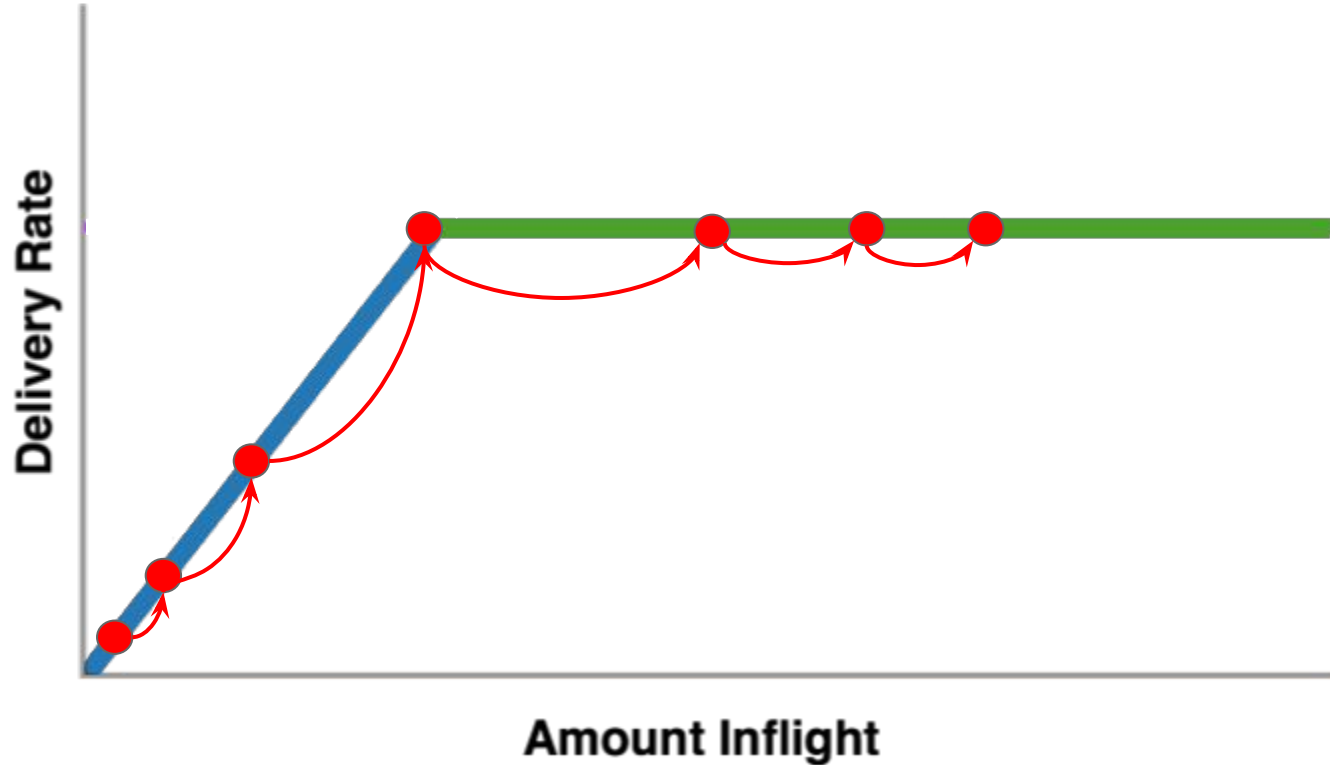
To see max BW, min RTT: probe both sides of BDP



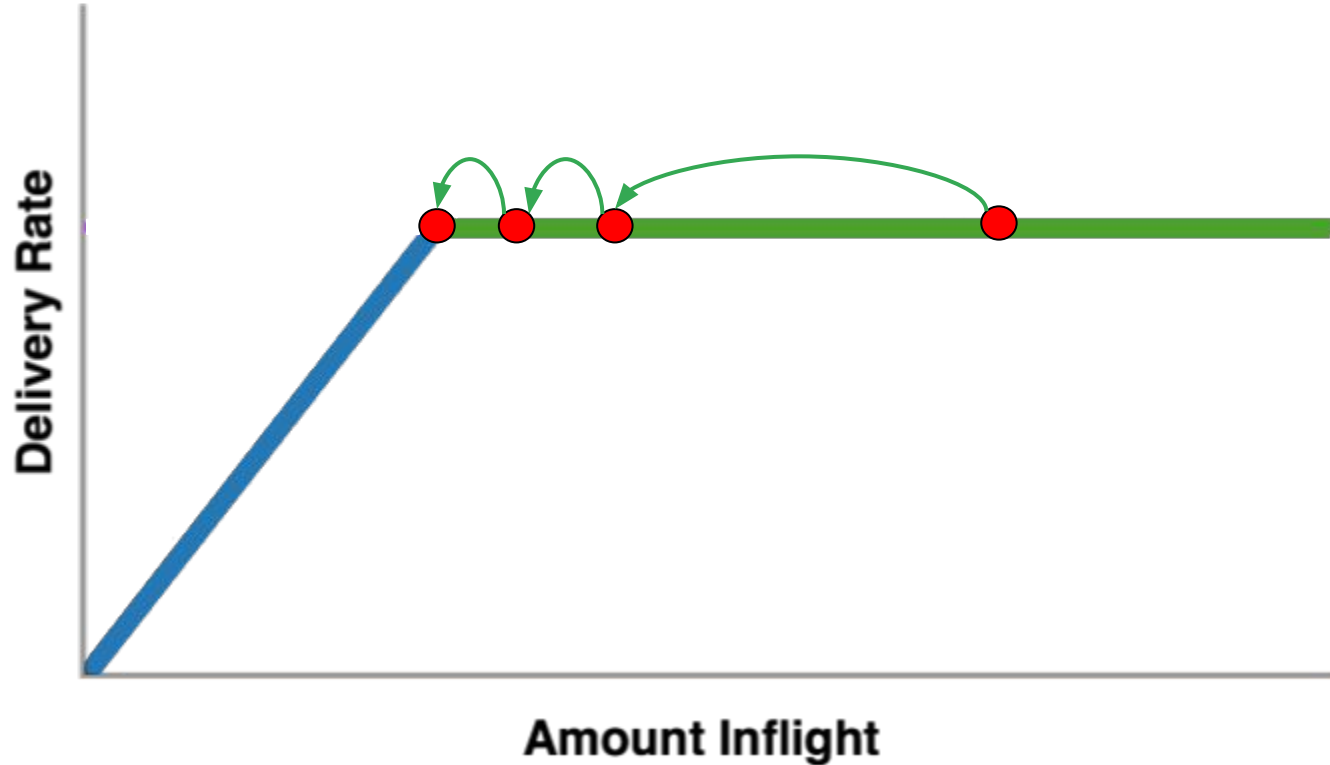
BBR: model-based walk toward max BW, min RTT



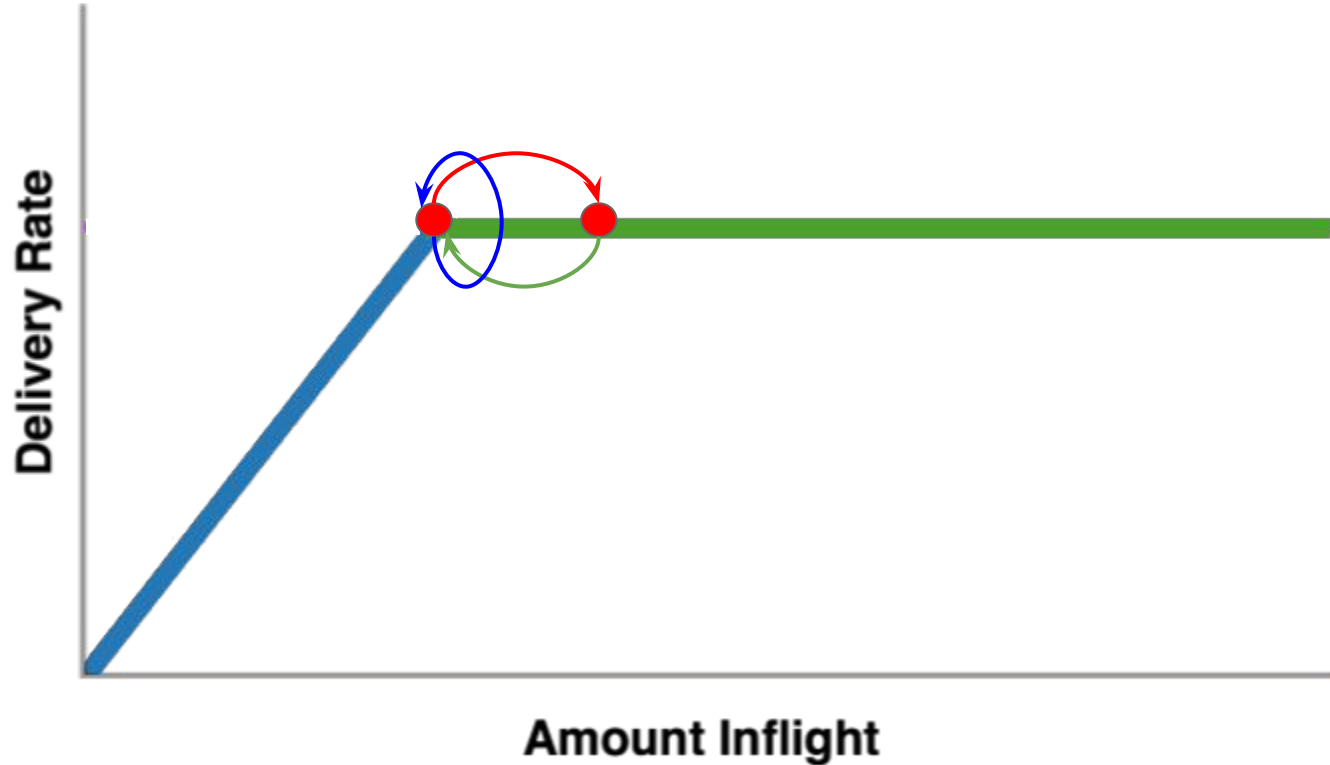
STARTUP: exponential BW search



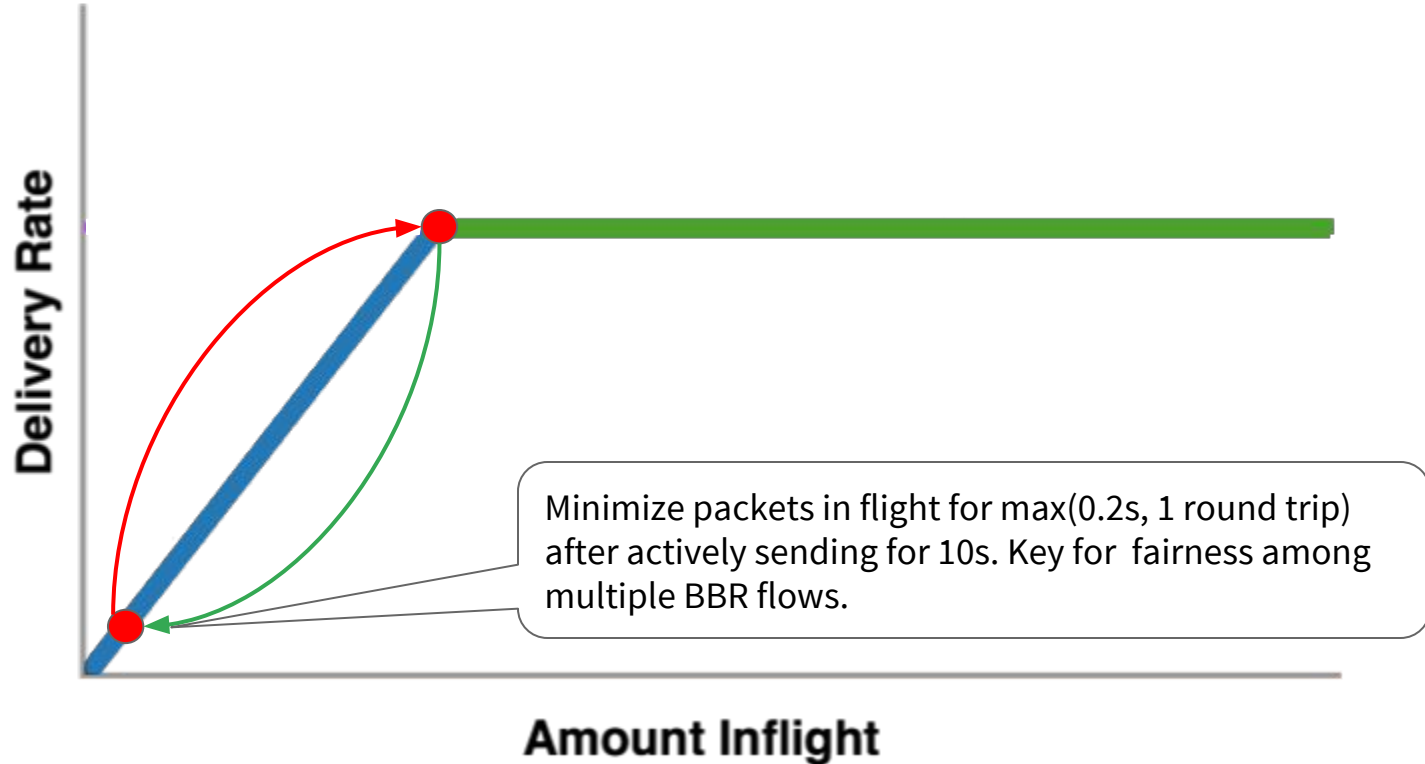
DRAIN: drain the queue created during STARTUP

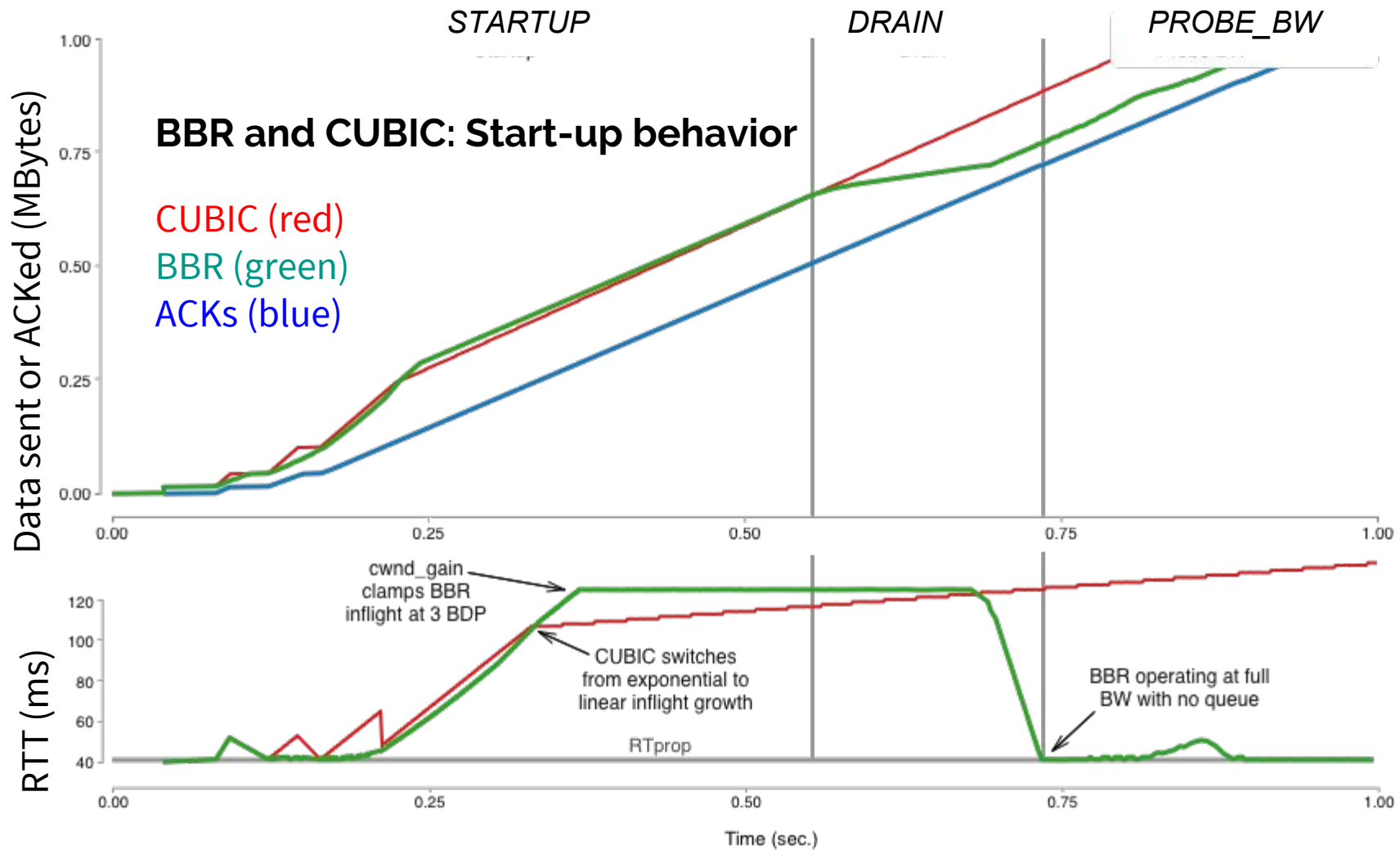


PROBE_BW: explore max BW, drain queue, cruise

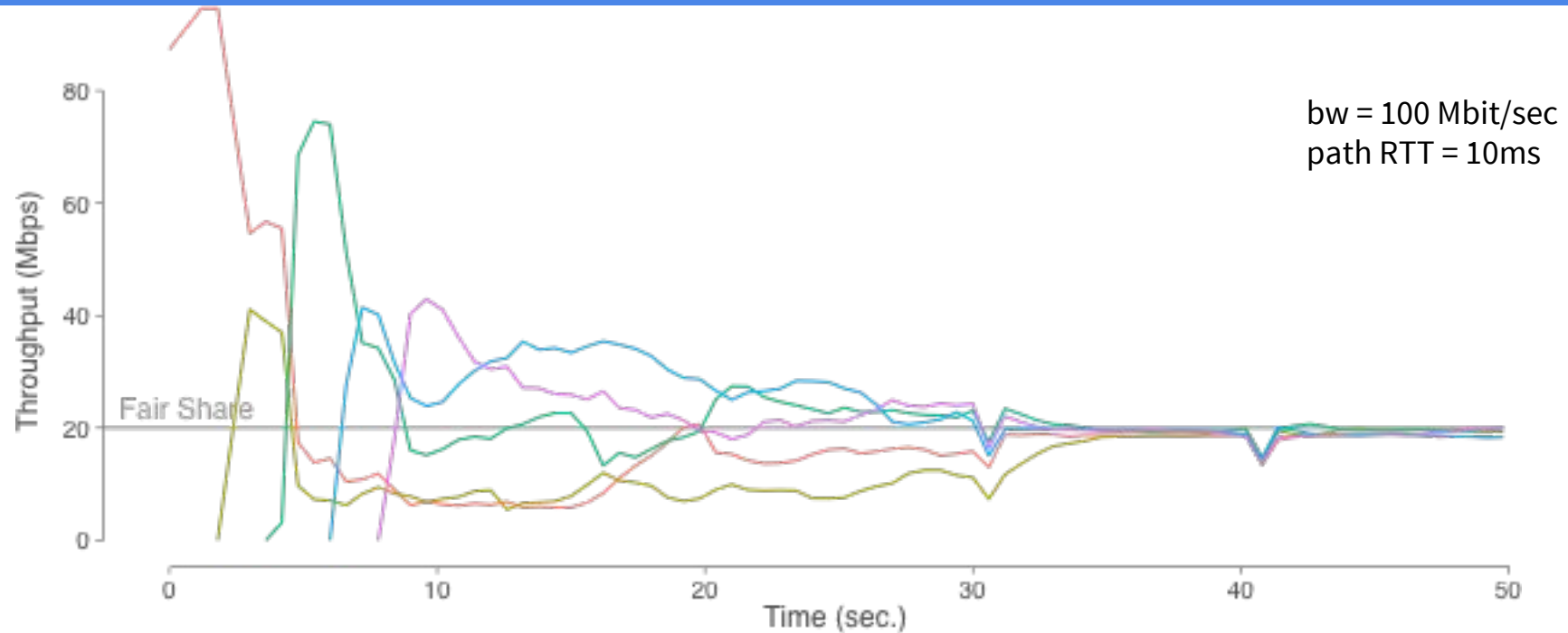


PROBE_RTT: drains queue to refresh min RTT



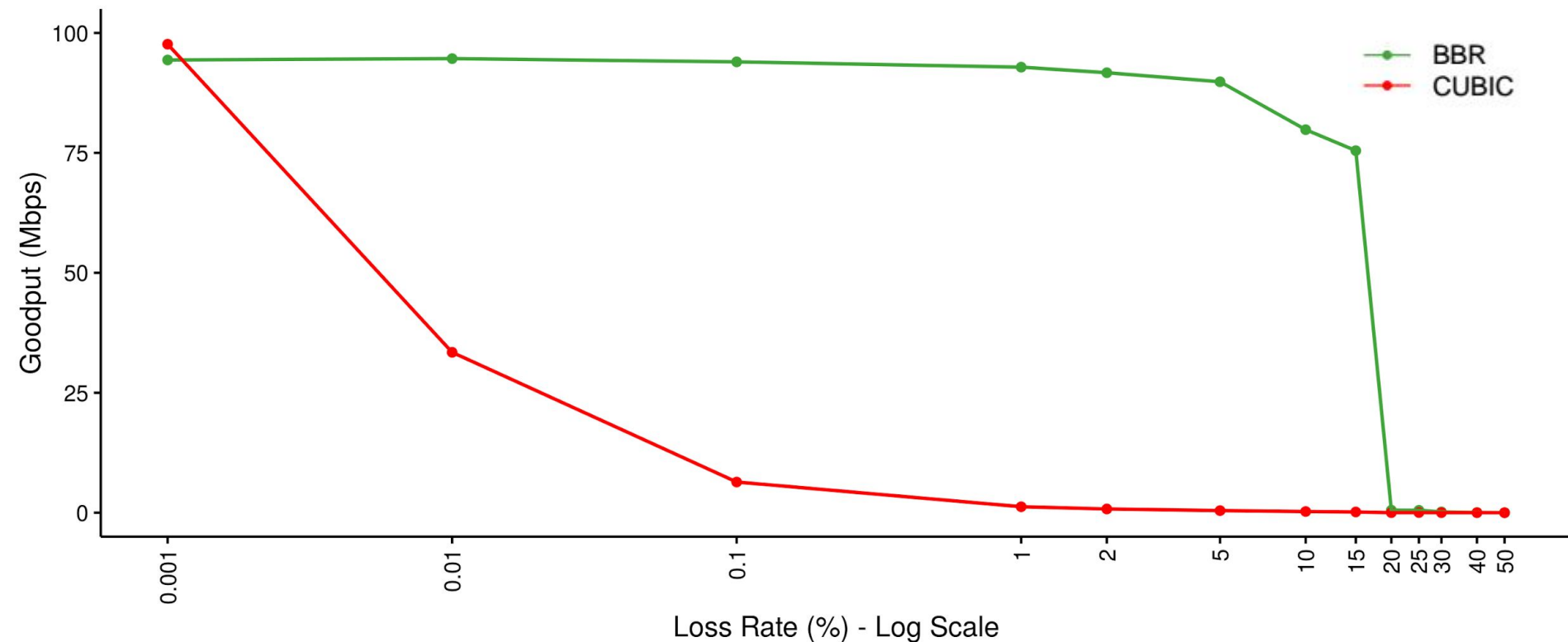


BBR multi-flow convergence dynamics



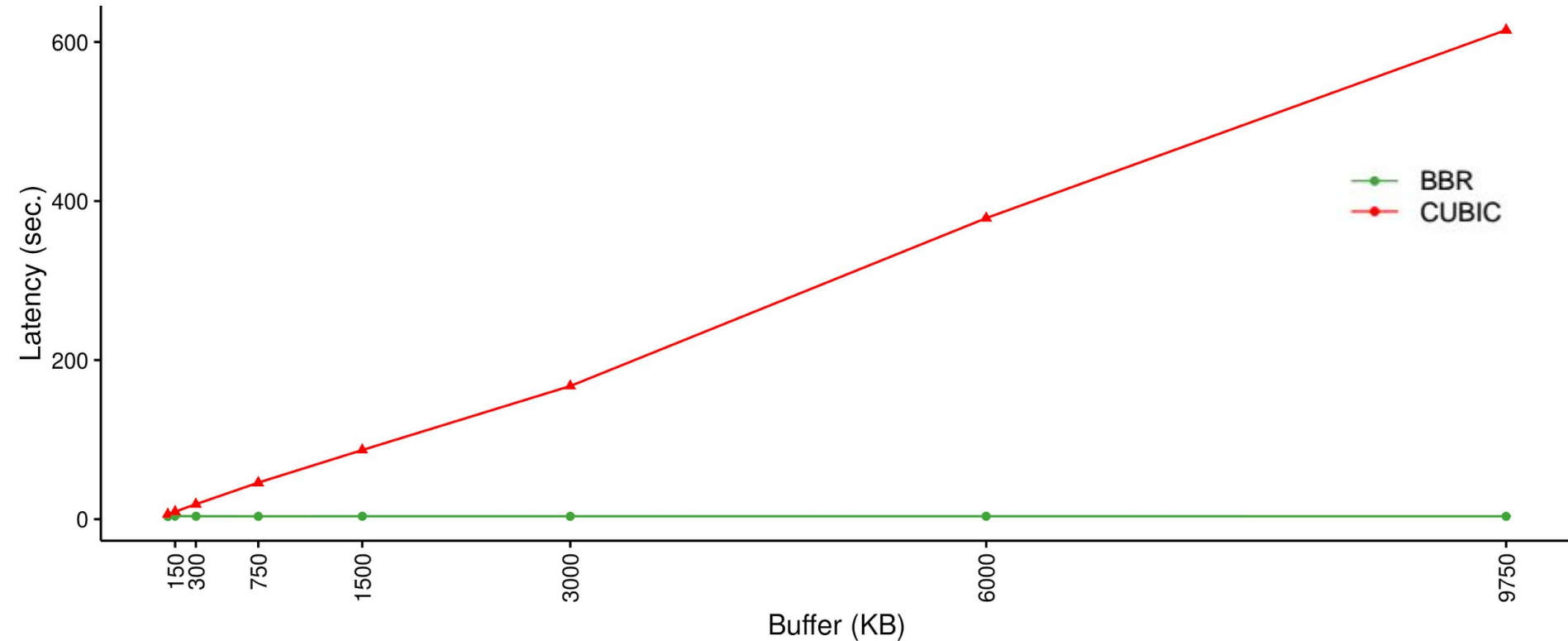
1. Flow 1 briefly slows down to reduce its queue every 10s (PROBE_RTT mode)
2. Flow 2 notices the queue reduction via its RTT measurements
3. Flow 2 schedules to enter slow down 10 secs later (PROBE_RTT mode)
4. Flow 1 and Flow 2 gradually converge to share BW fairly

BBR: fully use bandwidth, despite high packet loss



BBR vs CUBIC: synthetic bulk TCP test with 1 flow, bottleneck_bw 100Mbps, RTT 100ms

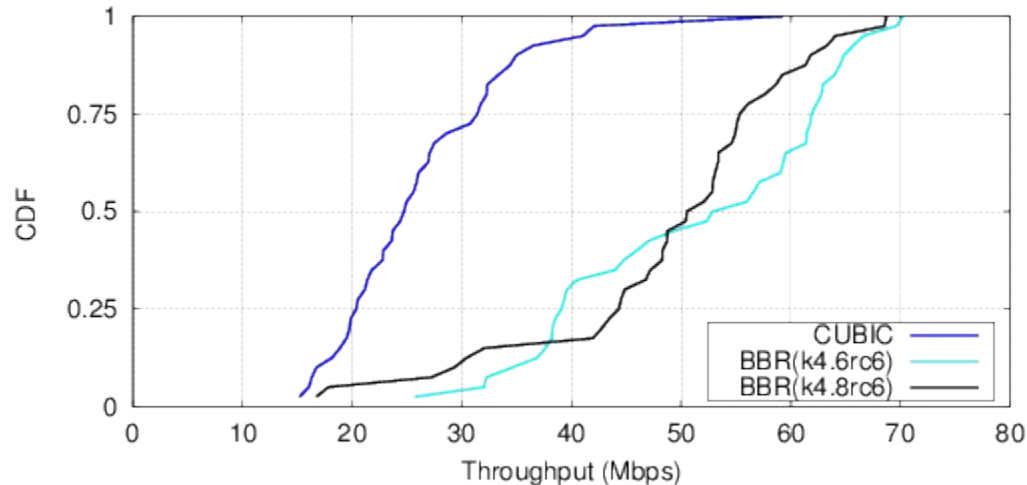
BBR: low queue delay, despite bloated buffers



BBR vs CUBIC: synthetic bulk TCP test with 8 flows, bottleneck_bw=128kbps, RTT=40ms

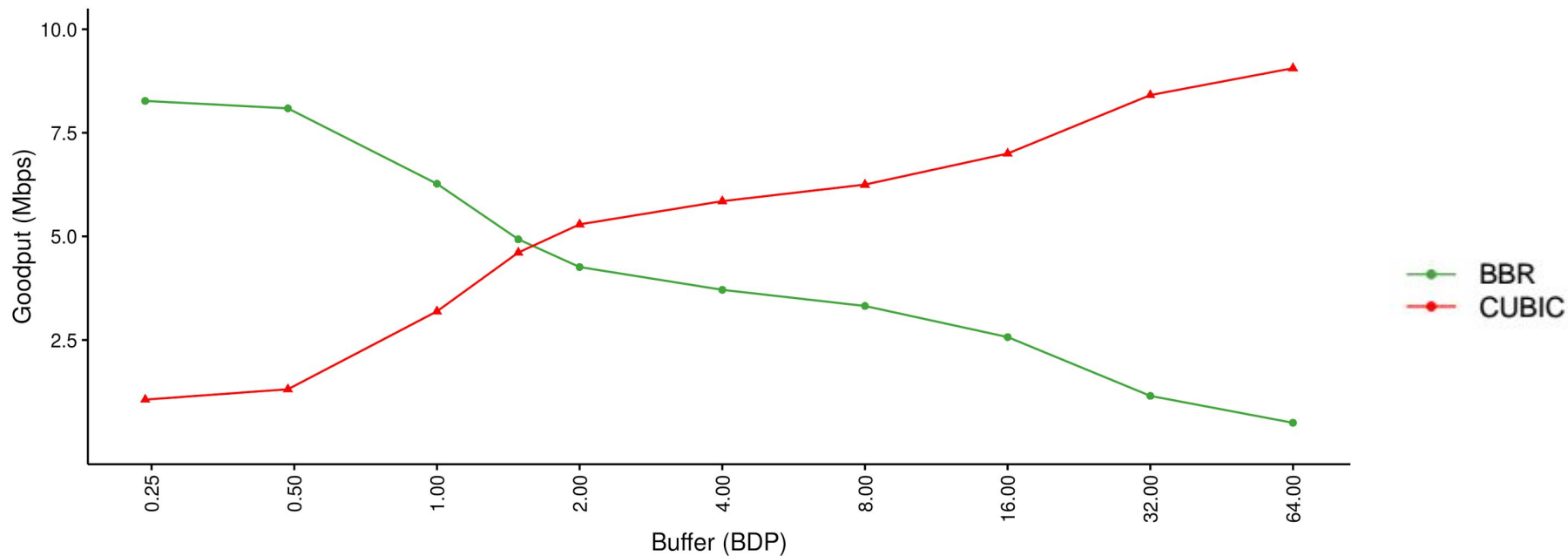
BBR: robust detection of full pipes -> faster start-up

- **BBR STARTUP**: estimate reached full BW if **BW** stops increasing significantly
- **CUBIC Hystart**: estimate reached full BW if **RTT** increases significantly
- But delay (RTT) can increase significantly well before full BW is reached!
 - Shared media links (cellular, wifi, cable modem) use slotting, aggregation
- e.g.: 20 MByte transfers over LTE (source: [post by Fung Lee on bbr-dev list, 2016/9/22](#)):

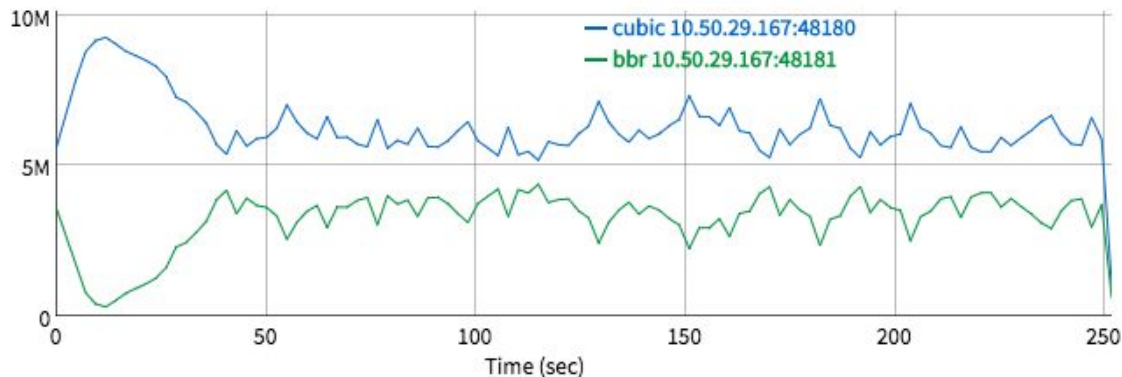


Improving dynamics w/ with loss-based CC

1xCUBIC v 1xBBR goodput: bw=10Mbps, RTT=40ms, 4min transfer, varying buffer sizes



BBR and loss-based CC in deep buffers: an example



At first CUBIC/Reno gains an advantage by filling deep buffers

But BBR does not collapse; it adapts: BBR's bw and RTT probing tends to drive system toward fairness

Deep buffer data point: $8 \cdot \text{BDP}$ case: $\text{bw} = 10\text{Mbps}$, $\text{RTT} = 40\text{ms}$, $\text{buffer} = 8 \cdot \text{BDP}$

-> CUBIC: 6.31 Mbps vs BBR: 3.26 Mbps