

# BBR: Congestion-Based Congestion Control

by **Neal Cardwell, Yuchung Cheng, C. Stephen Gunn, Soheil Hassas Yeganeh, Van Jacobson**  
(paper review)

**Hongyu Hè**

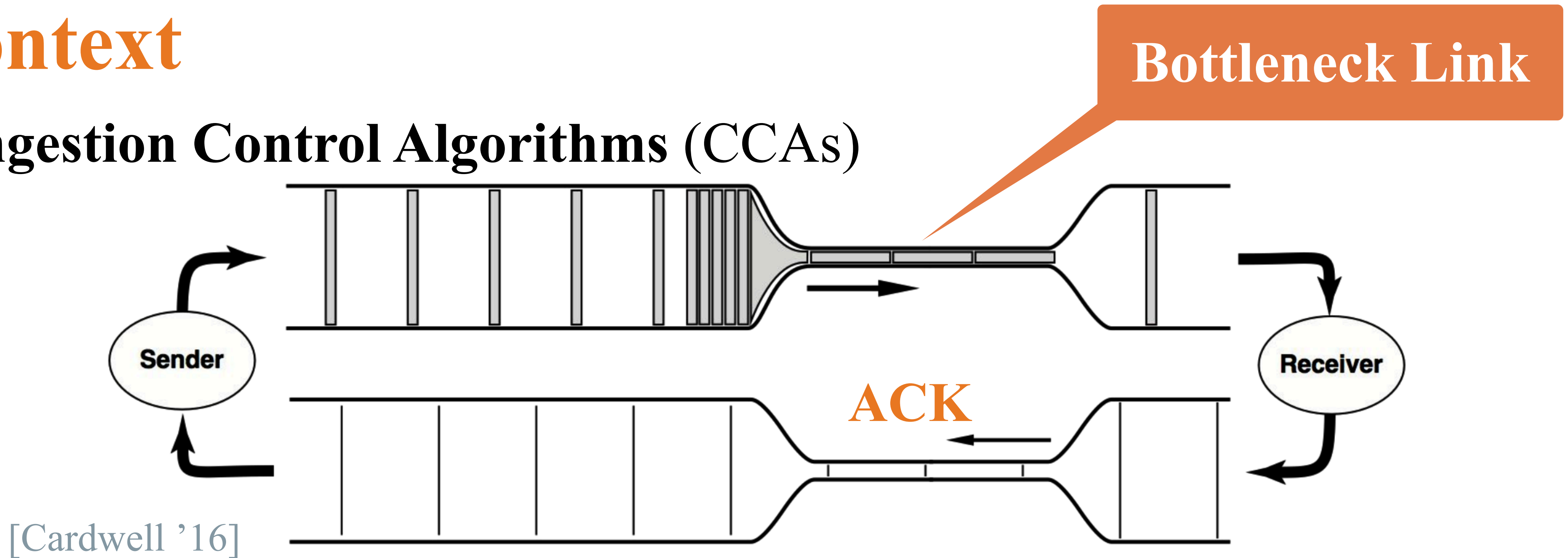
hhy@princeton.edu

September 20, 2024



# Context

## Congestion Control Algorithms (CCAs)



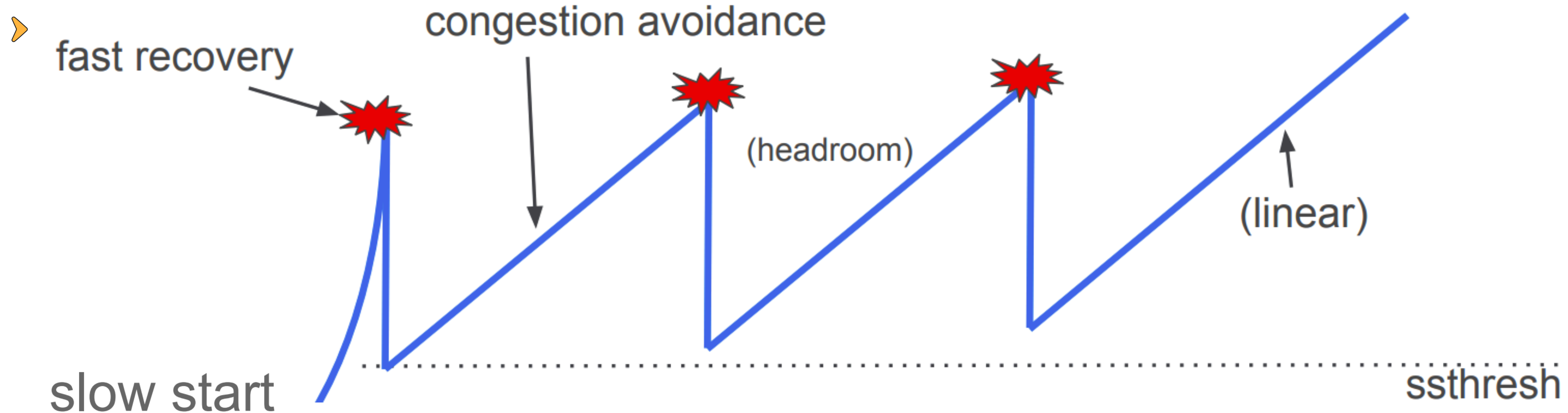
## Motivation (2011–2013)

- NICs w/ more memory → excessive buffering → TCP bufferbloat
- Single-conn HTTPv2  $\ll$  multi-conn HTTPv1
- Switches w/ shallow buffers have low TCP throughput

# Problem Statement

## Loss-based CCA is problematic

- e.g., New Reno:



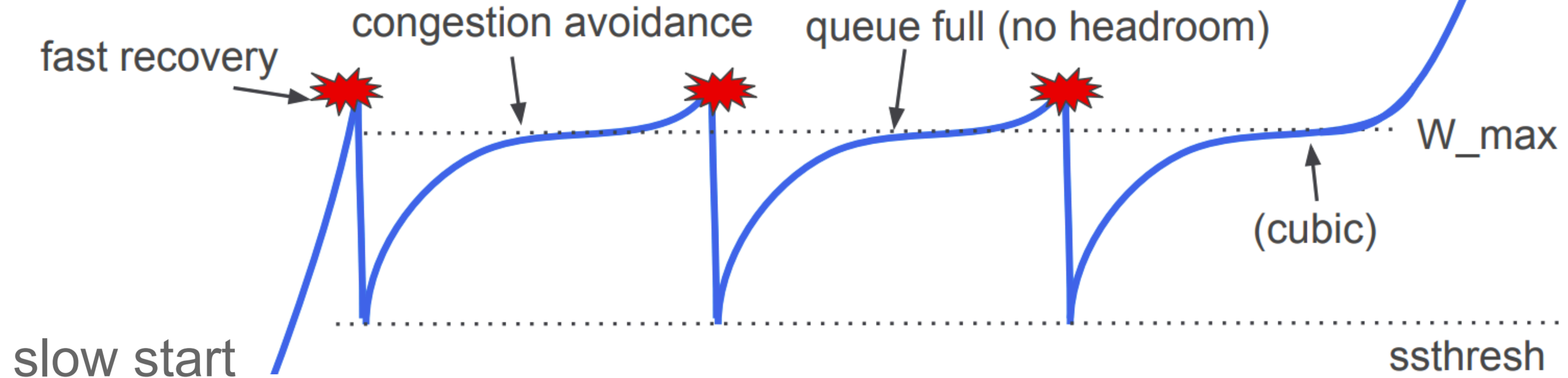
- Linear growth:  $1000\times$  more BW needs 1000 increases

⦿ saturate 10Gbps BW; 50ms RTT  $\rightarrow$   $\sim 35\text{min}$   $\Rightarrow$  loss rate  $\leq 5.7e-10$

# Problem Statement

## Loss-based CCA is problematic

- e.g., CUBIC:



- Cubic growth:  $1000\times$  more BW needs 10 increases

⦿ saturate 10Gbps BW; 50ms RTT  $\rightarrow$   $\sim 7$ min  $\Rightarrow$  loss rate  $\leq 2.86\text{-e}9$

# Problem Statement

## Loss-based CCA is problematic

- Many CCAs were loss-based (e.g., Tahoe, New Reno, and CUBIC)
- Packet loss  $\Rightarrow$  congestion?
- Loss-based CCAs + shallow or deep buffers  $\Rightarrow$  poor performance



# BBR: Bottleneck Bandwidth & Round-trip propagation time

## Key Idea:

Explicitly model bottleneck queue by probing the RTT and bottleneck BW periodically to estimate the bandwidth-delay product (BDP)

## Main Contributions:

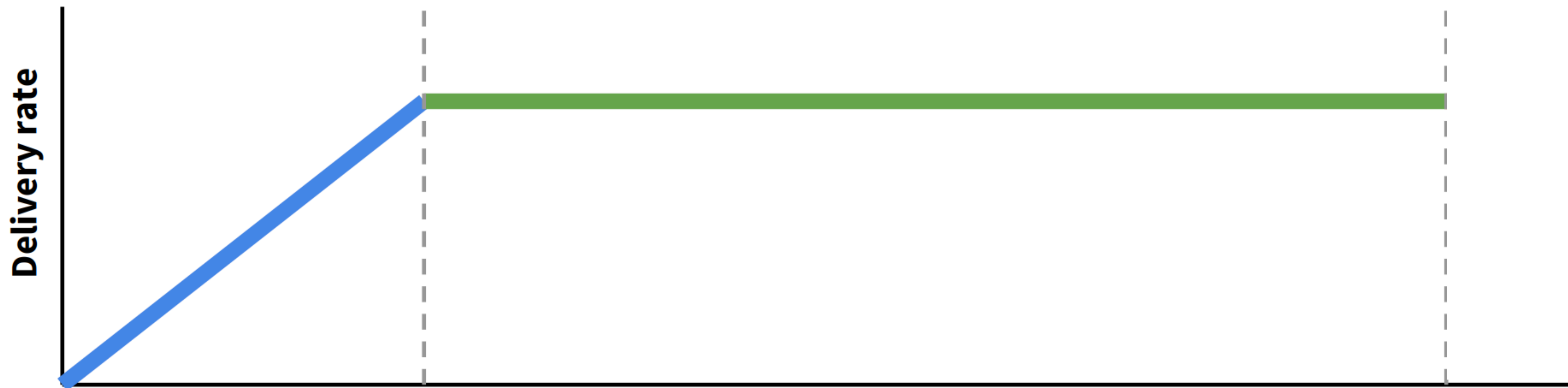
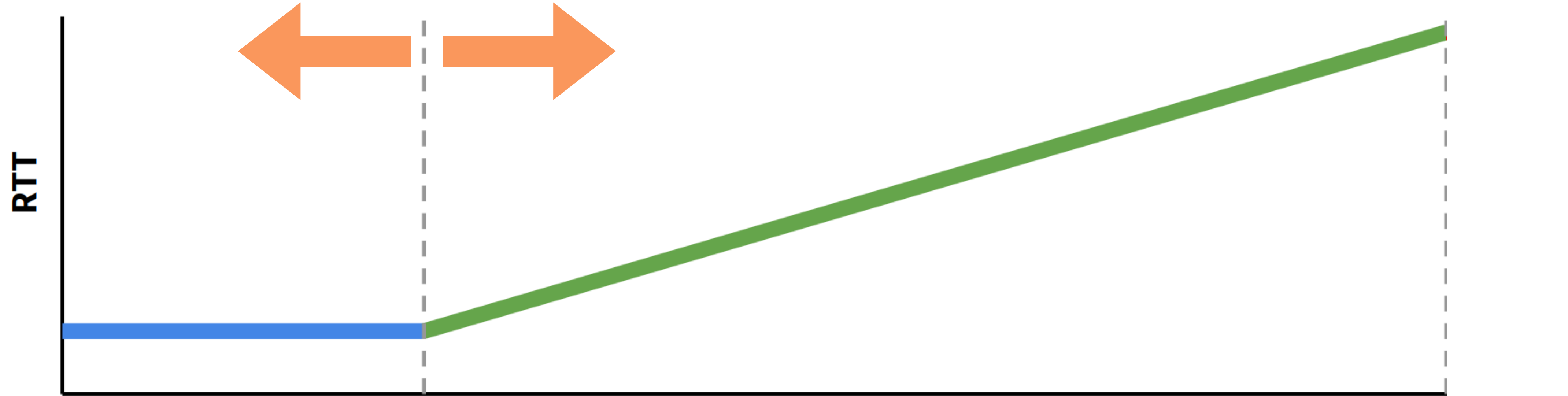
- 1) Identified and addressed a root cause for an internet-scale problem
- 2) BBR: client-side rate-based CCA, better latency and tput vs. CUBIC
  - Fast search rate  $O(\log \text{BDP})$
- 3) Production deployment, evaluation, and linux integration with LTS

# How does BBR work?

Insufficient Traffic

Queuing at Bottleneck

Bottleneck Overflow



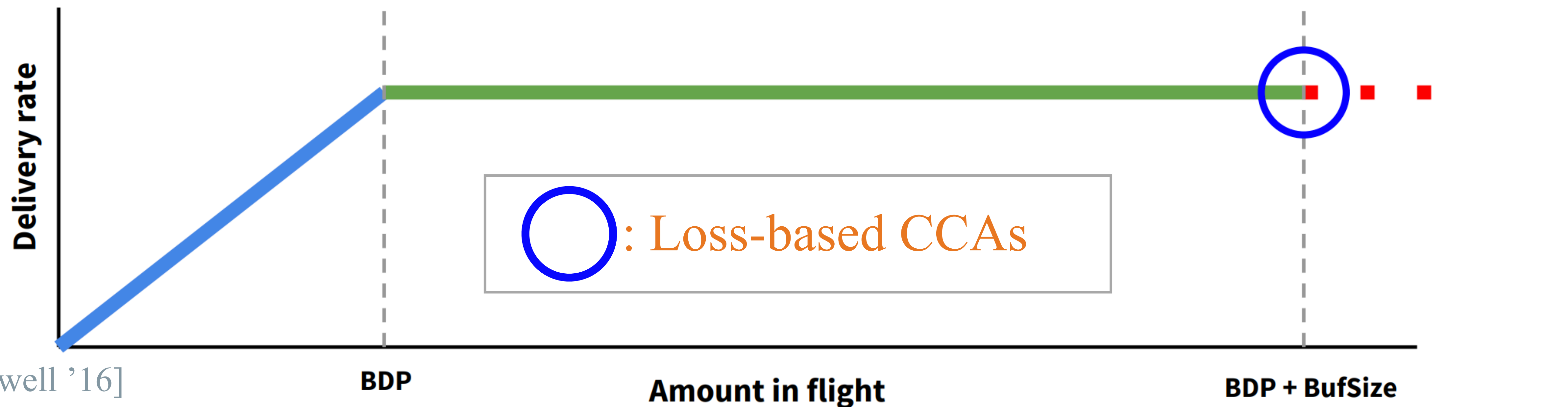
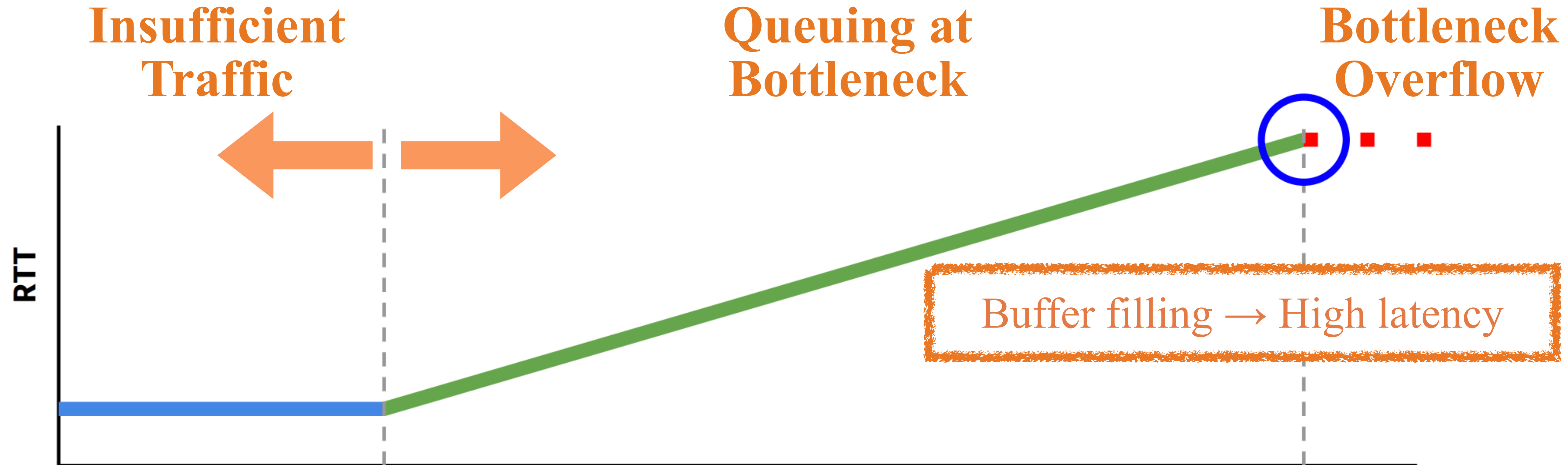
[Cardwell '16]

BDP

Amount in flight

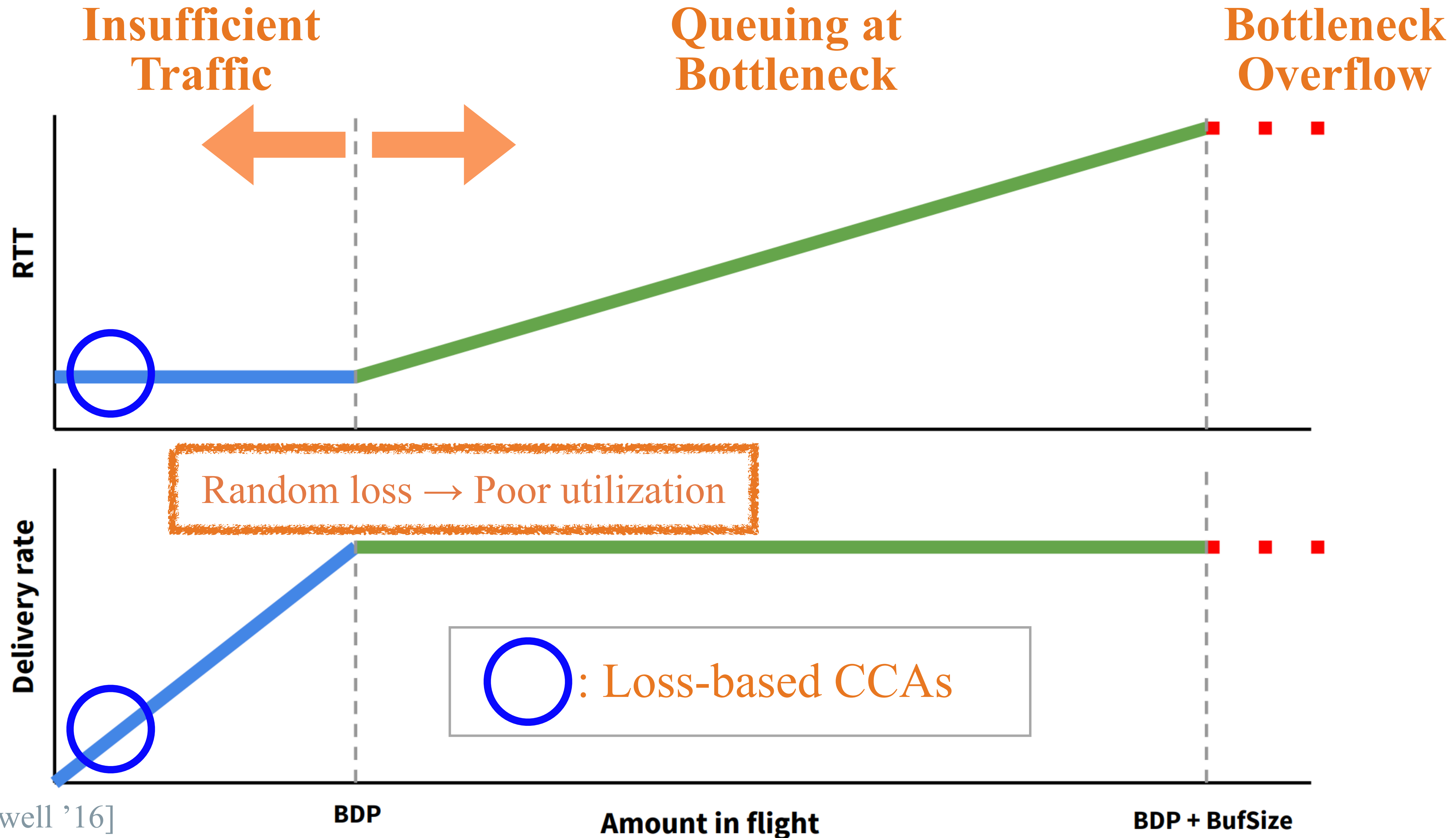
BDP + BufSize



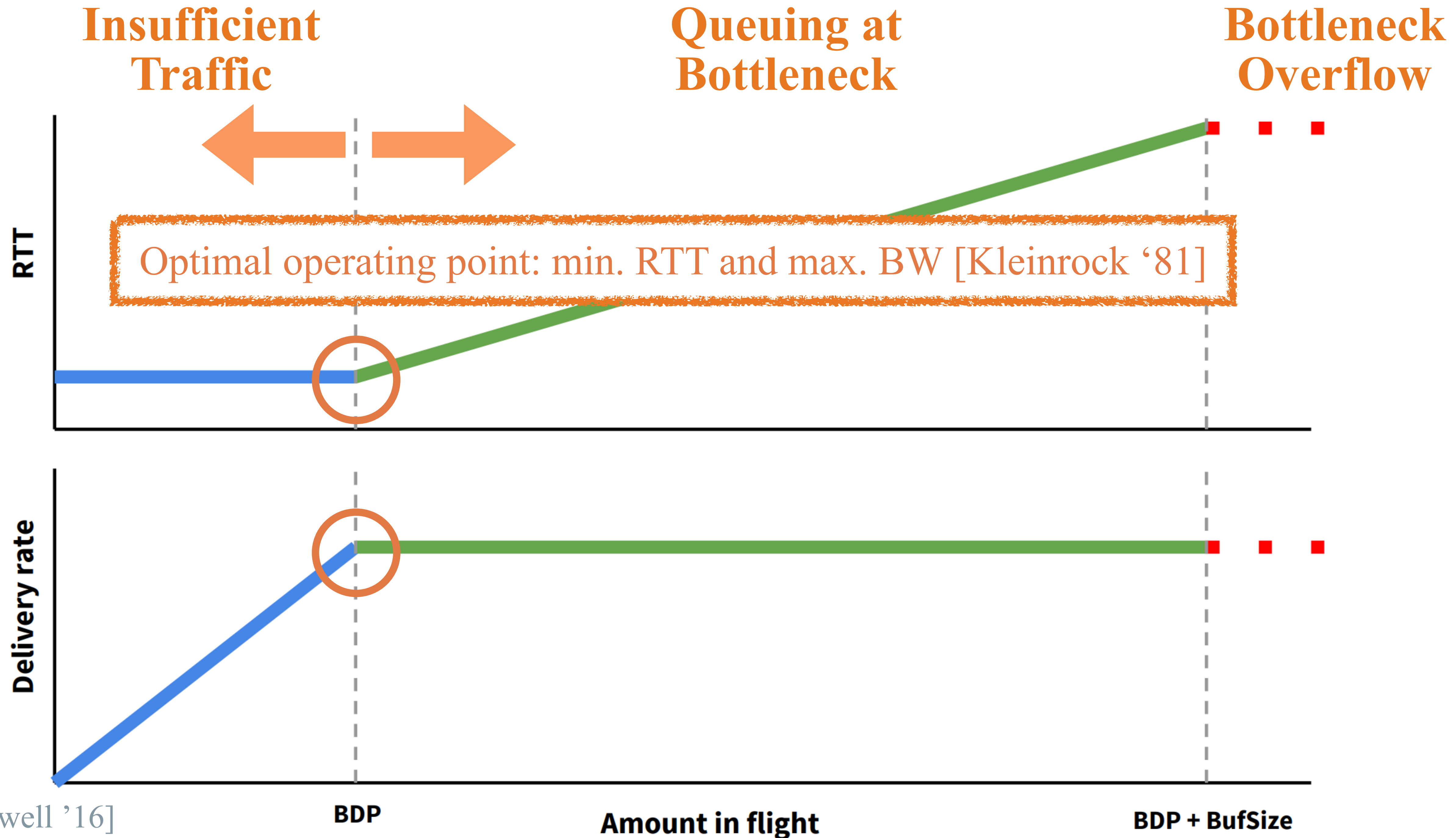


[Cardwell '16]

Amount in flight



[Cardwell '16]



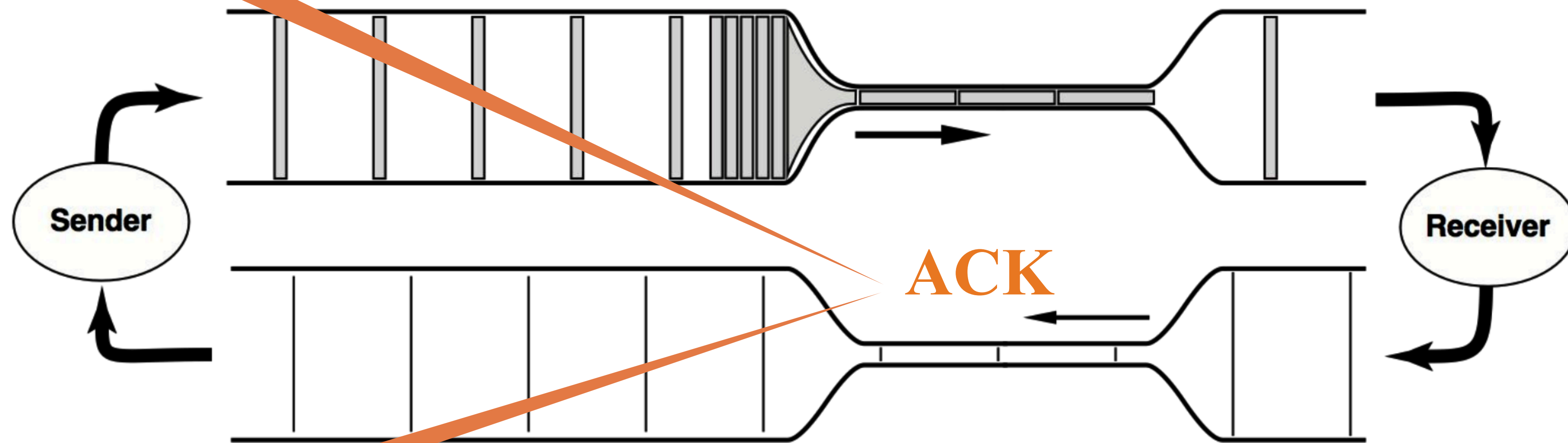
[Cardwell '16]

# BBR's Design

# Dynamically Estimate Windowed Max BW and Min RTT

$\hat{B}W$

$$\forall i \in W_{WB} : \hat{B}W \leftarrow \max \{ (\Delta_{\text{delivered}} / \Delta t)_i \}$$



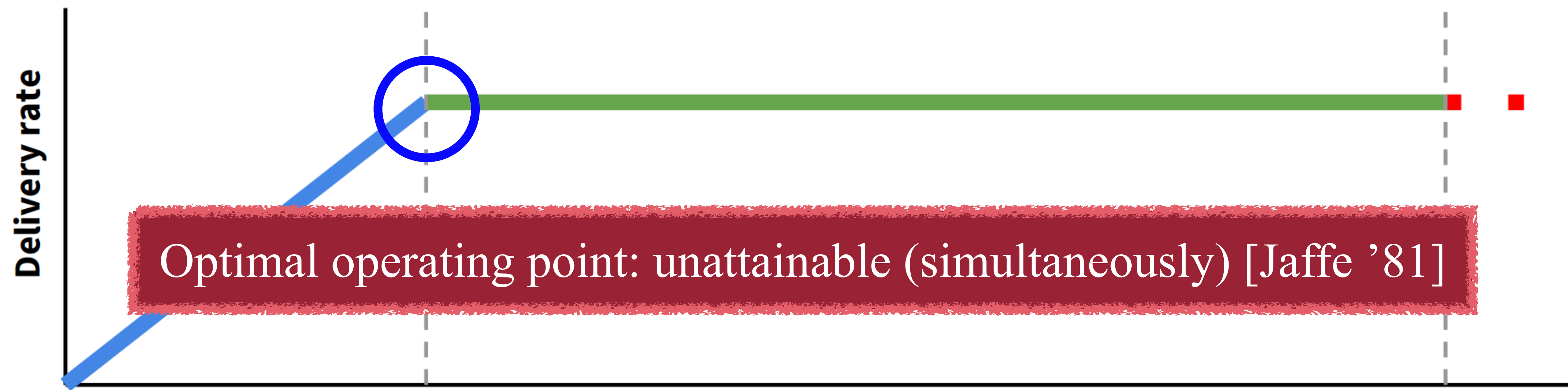
$\hat{R}TT$

$$\forall i \in W_{RTT} : \hat{R}TT \leftarrow \min \{ RTT_i \}$$

Insufficient Traffic

Queuing at Bottleneck

Bottleneck Overflow



[Cardwell '16]

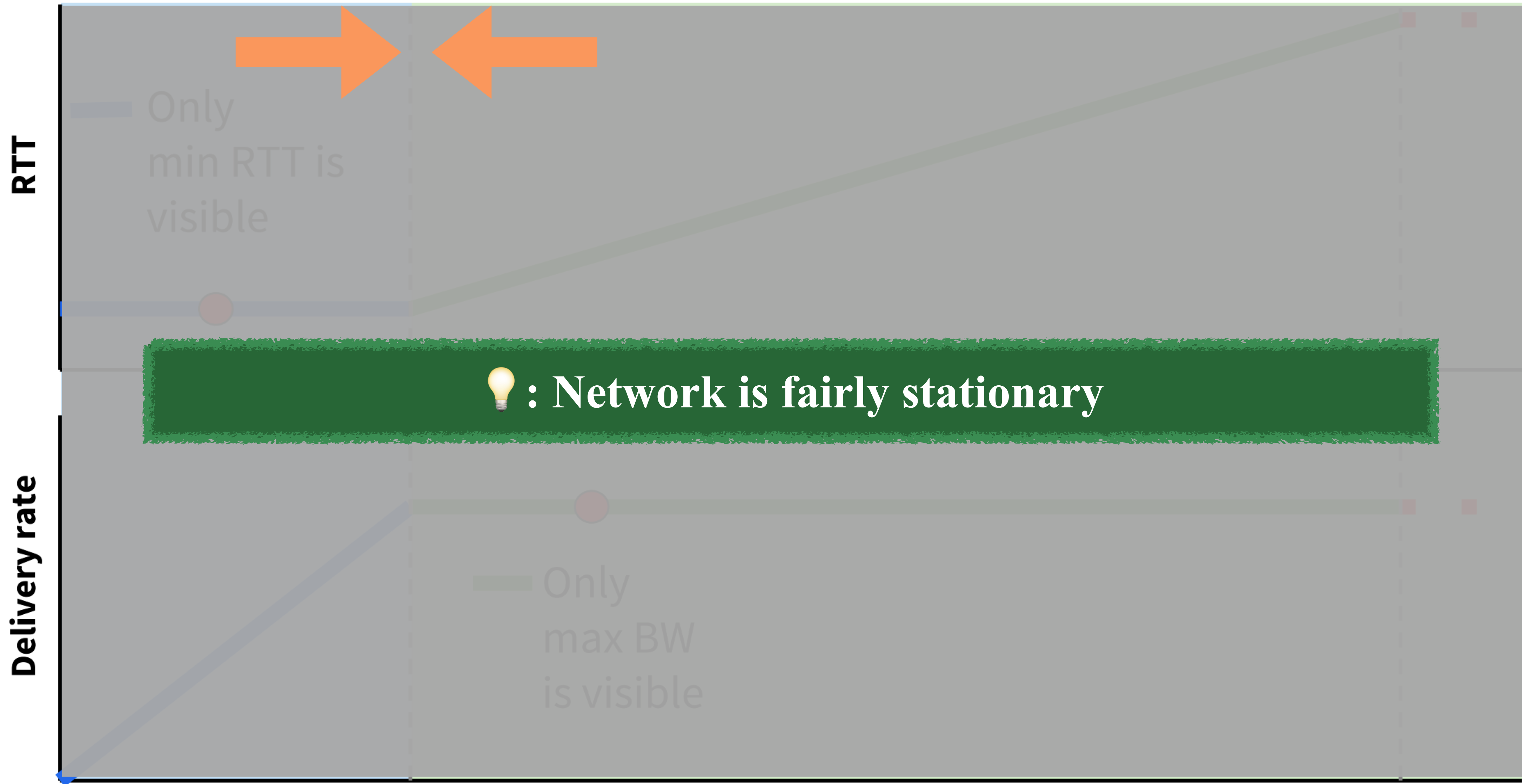
BDP

Amount in flight

BDP + BufSize



# Sequentially Probe Max BW and Min RTT



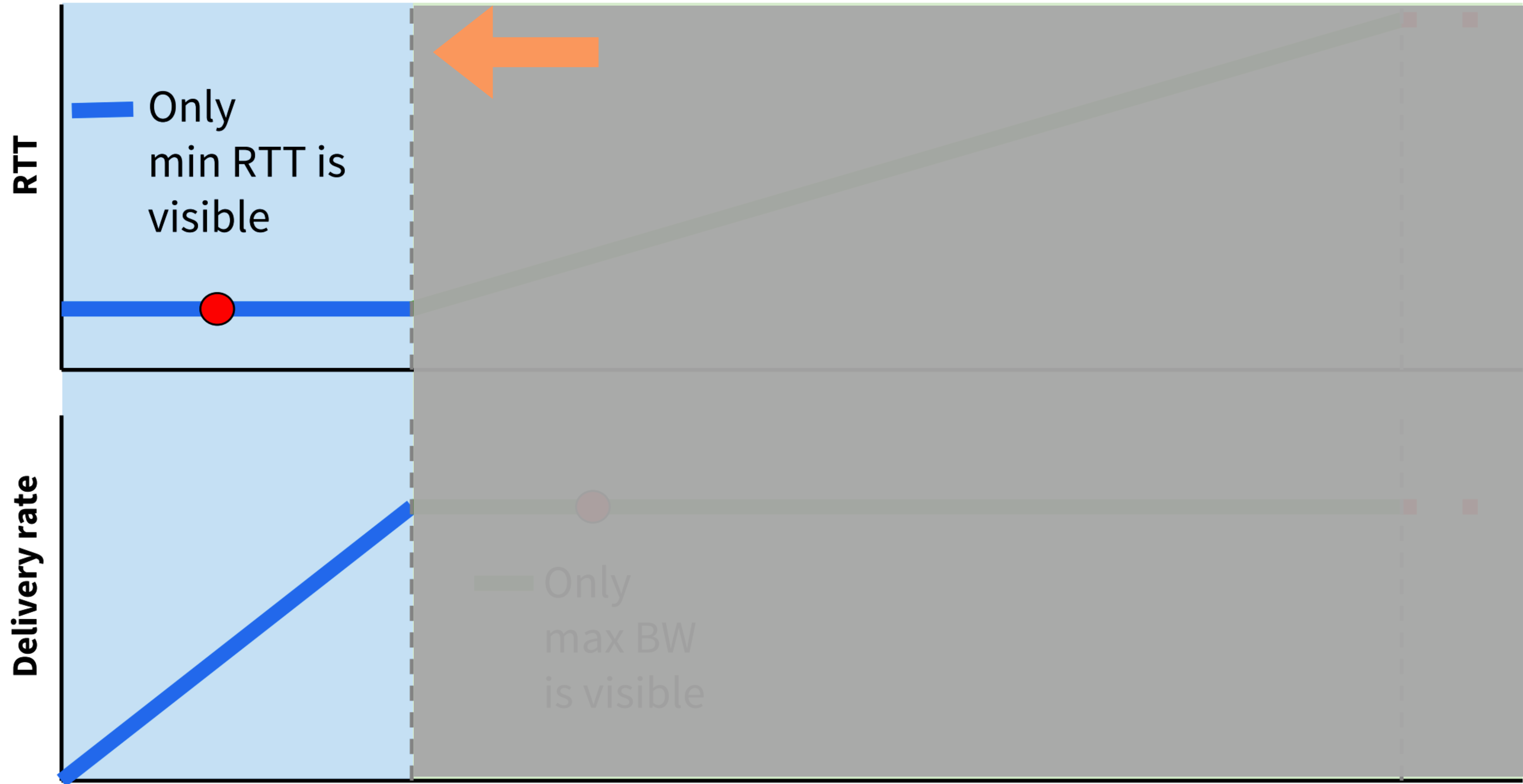
[Cardwell '16]

**BDP**

**amount in flight**

**BDP + BufSize**

# Sequentially Probe Max BW and Min RTT



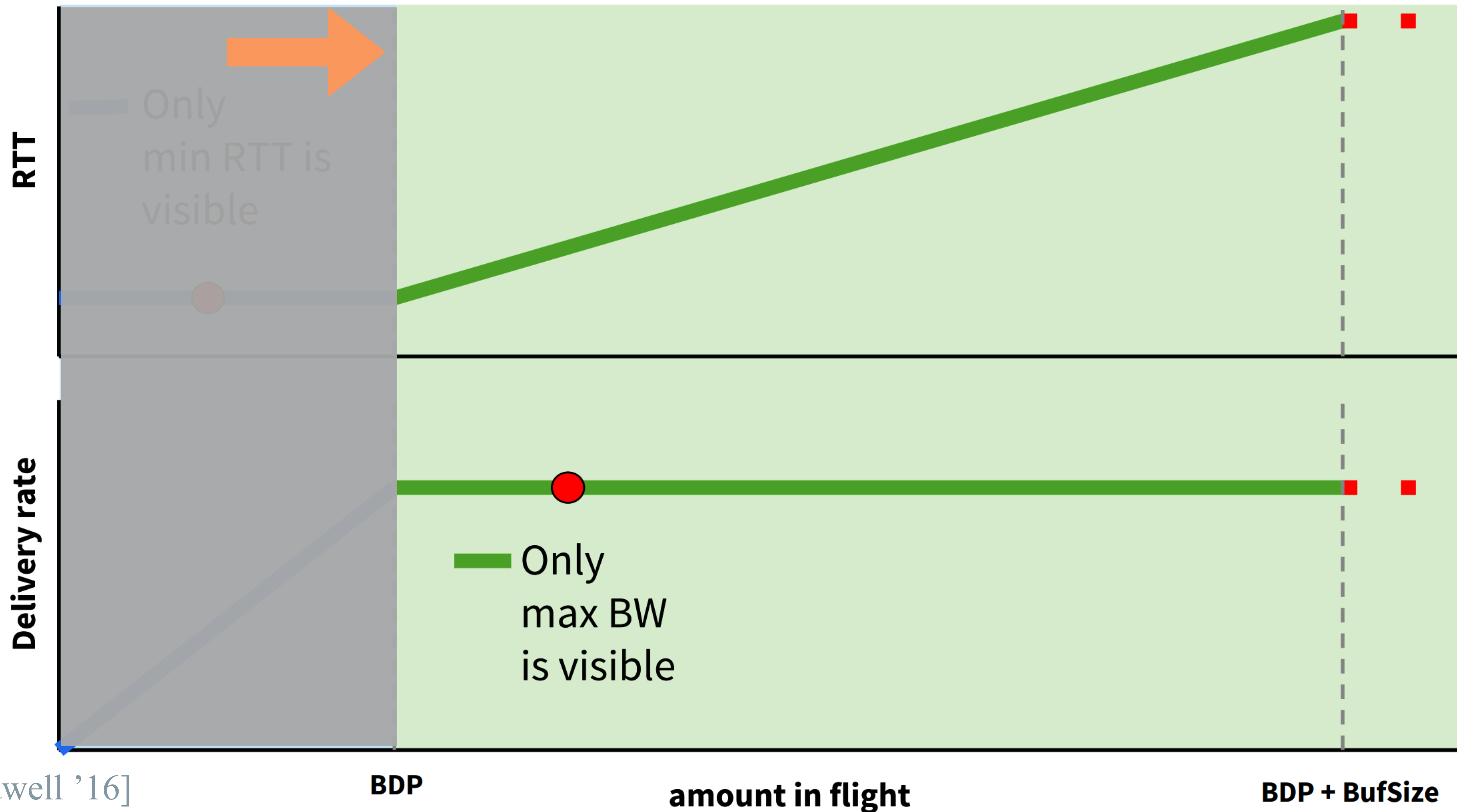
[Cardwell '16]

**BDP**

**amount in flight**

**BDP + BufSize**

# Sequentially Probe Max BW and Min RTT



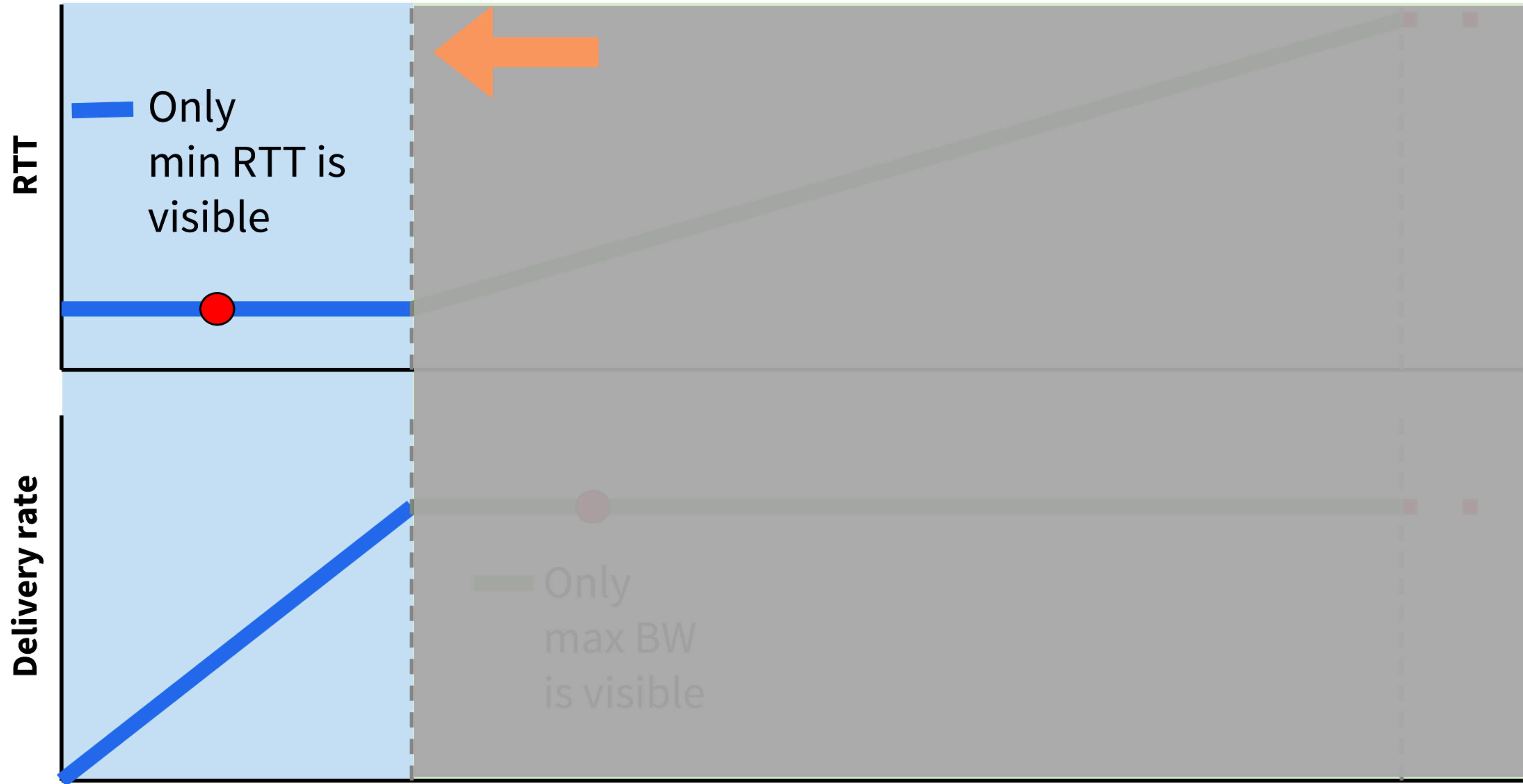
[Cardwell '16]

**BDP**

**amount in flight**

**BDP + BufSize**

# Sequentially Probe Max BW and Min RTT



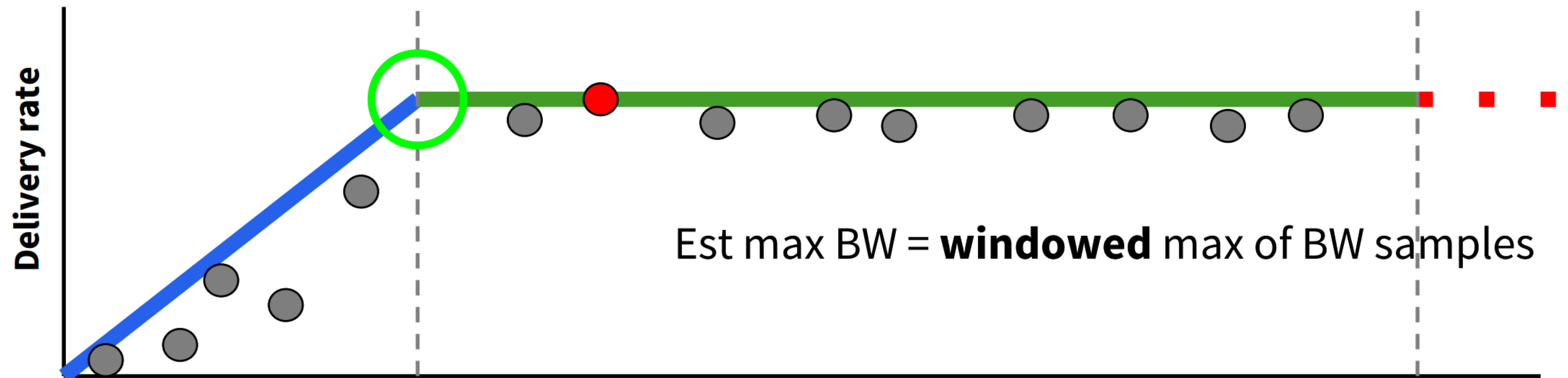
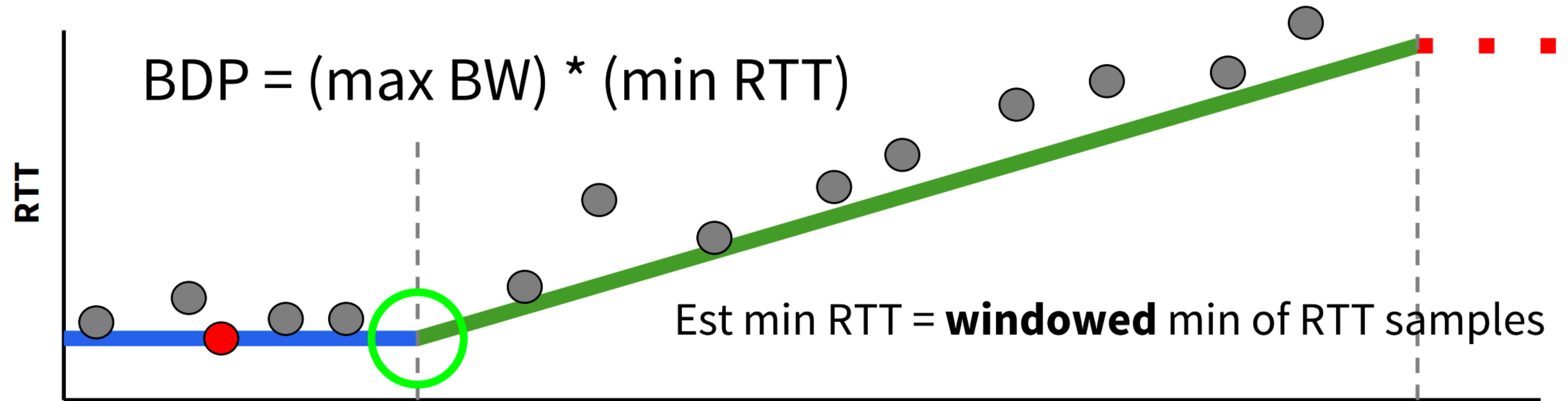
[Cardwell '16]

BDP

amount in flight

BDP + BufSize

# BDP Estimation



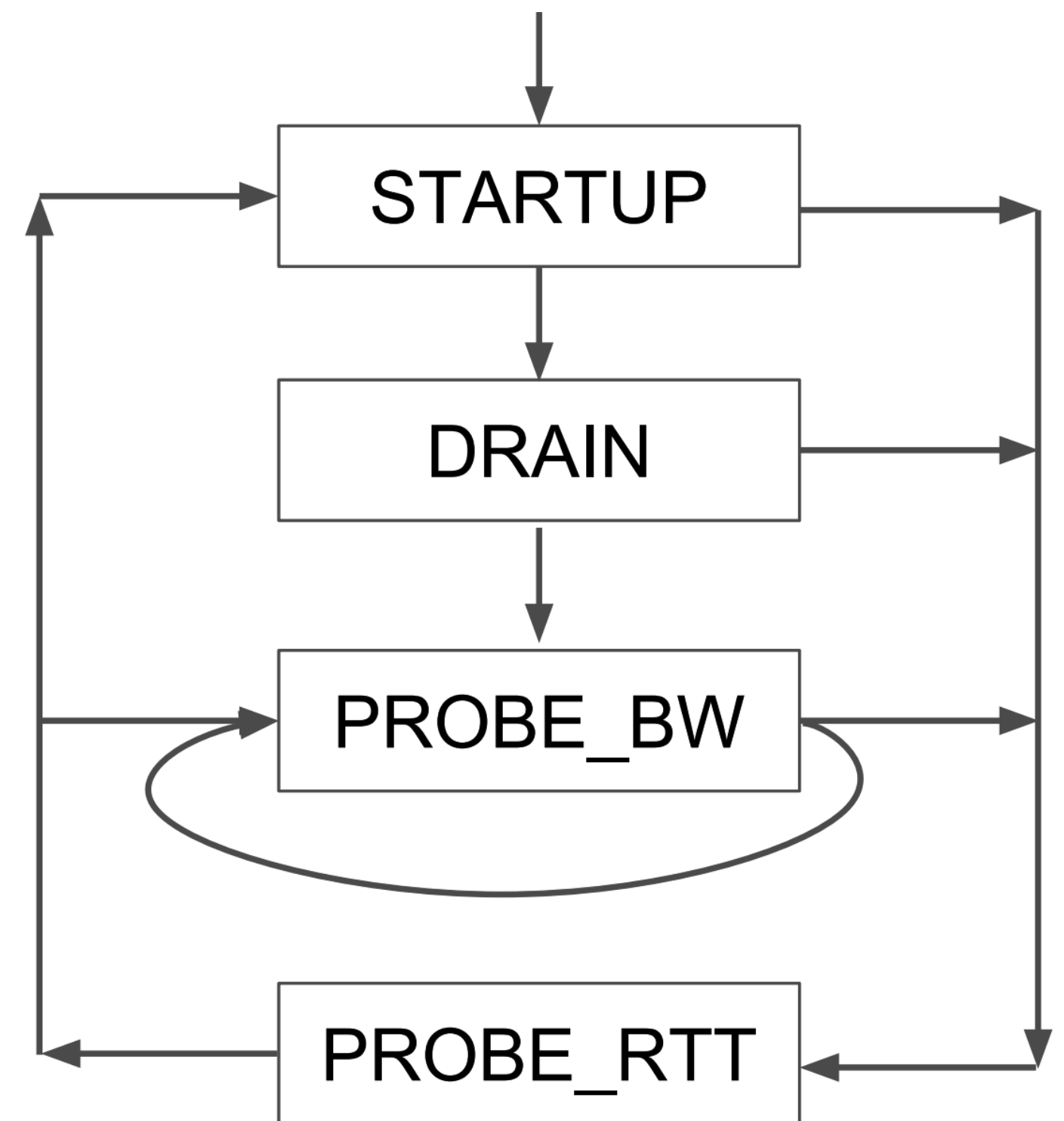
[Cardwell '16]

BDP

amount in flight

BDP + BufSize

# BBR's State Machine

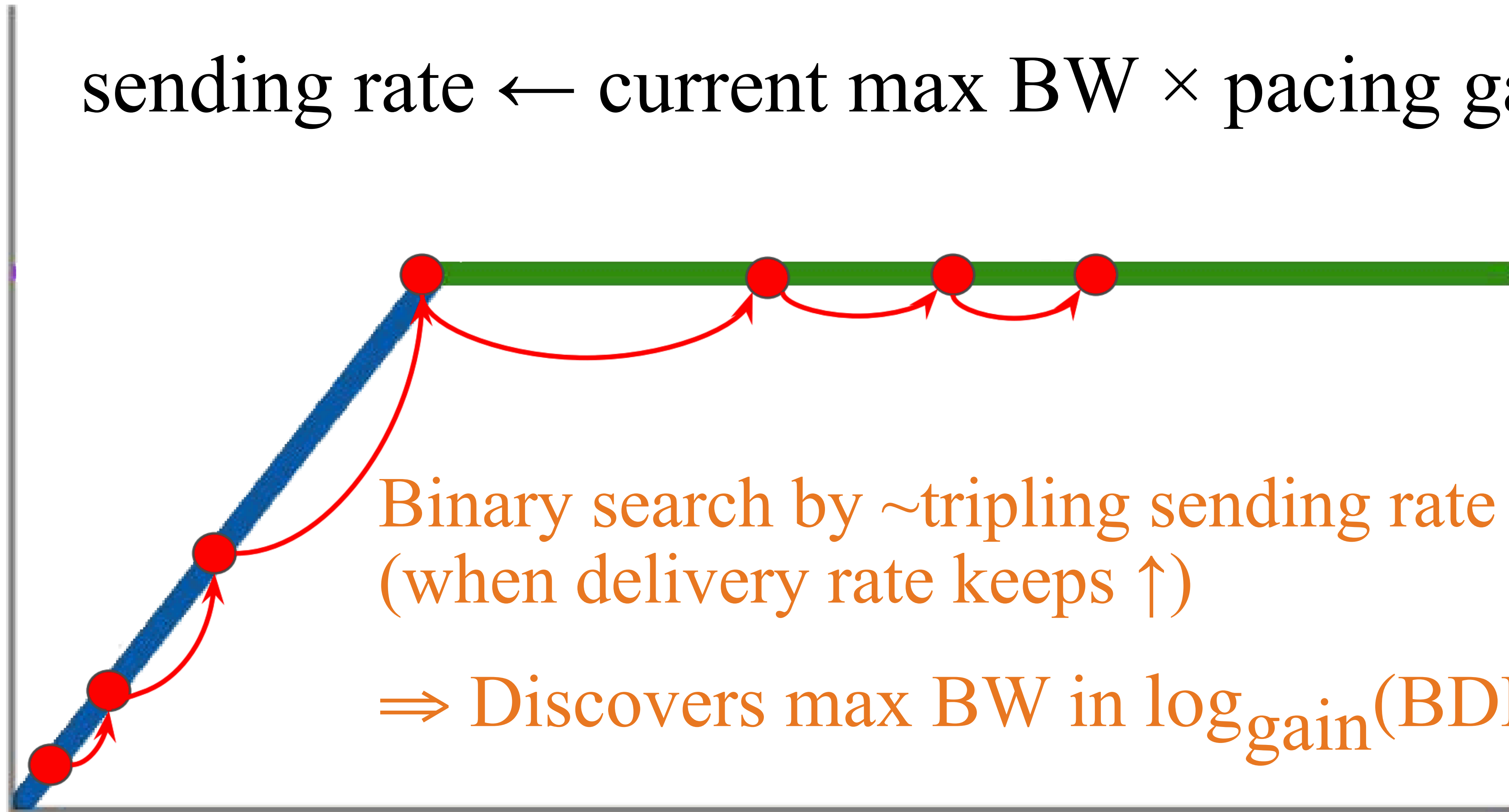




# Startup: Exponential BW search ( $\approx$ slow-start)

sending rate  $\leftarrow$  current max BW  $\times$  pacing gain

Delivery Rate

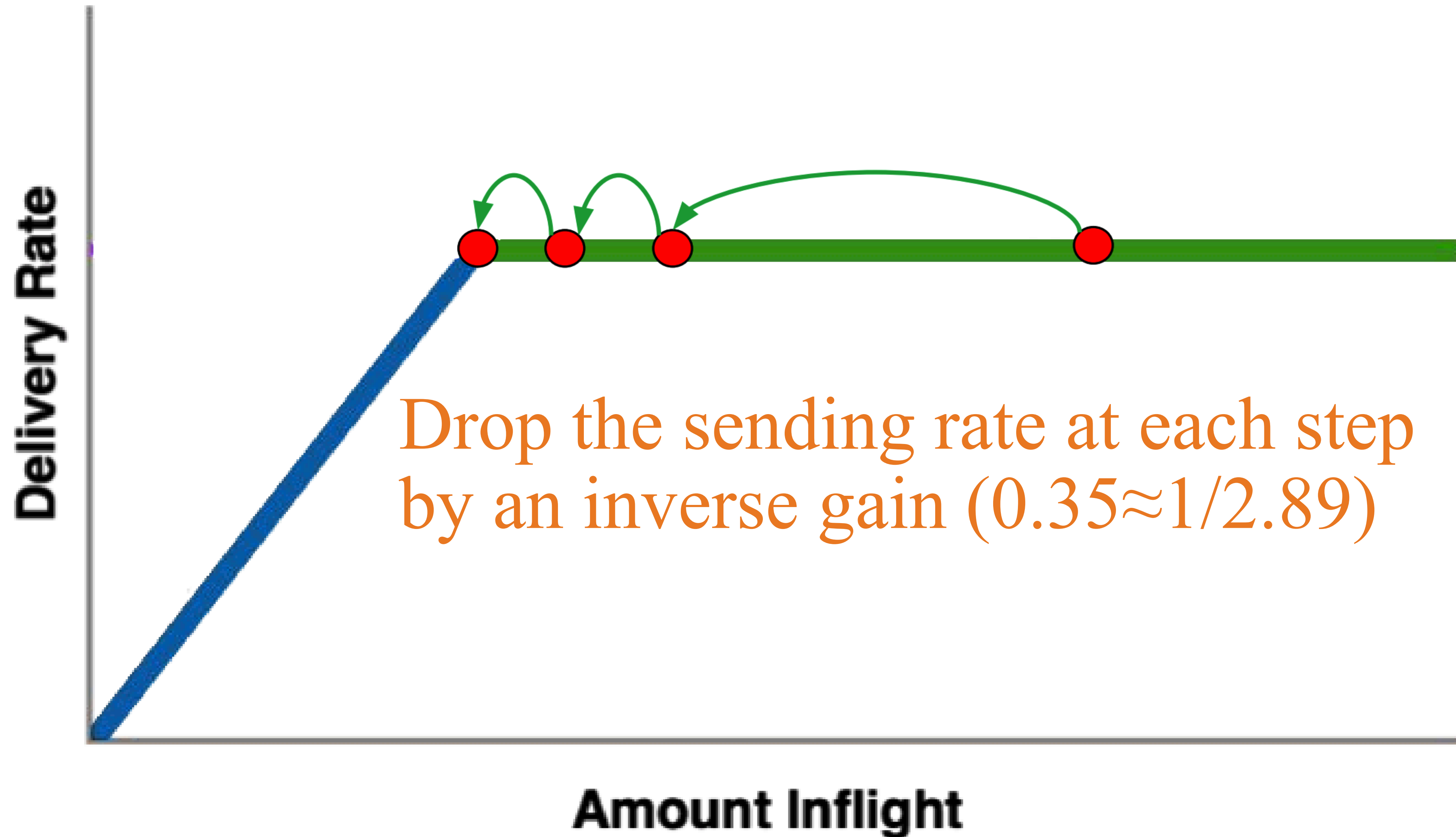


Binary search by  $\sim$ tripling sending rate (2.89 gain)  
(when delivery rate keeps  $\uparrow$ )

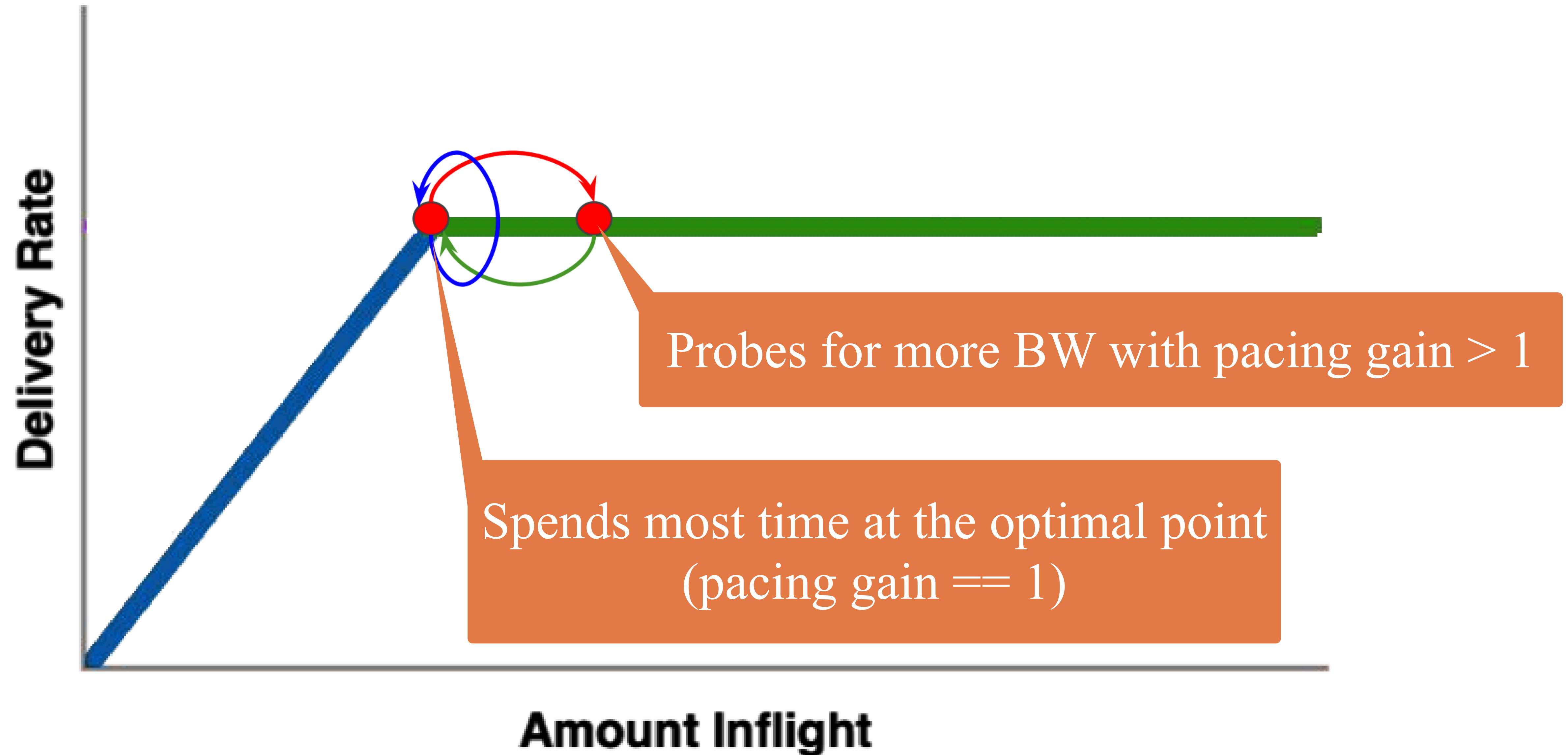
$\Rightarrow$  Discovers max BW in  $\log_{\text{gain}}(\text{BDP})$  RTTs

Amount Inflight

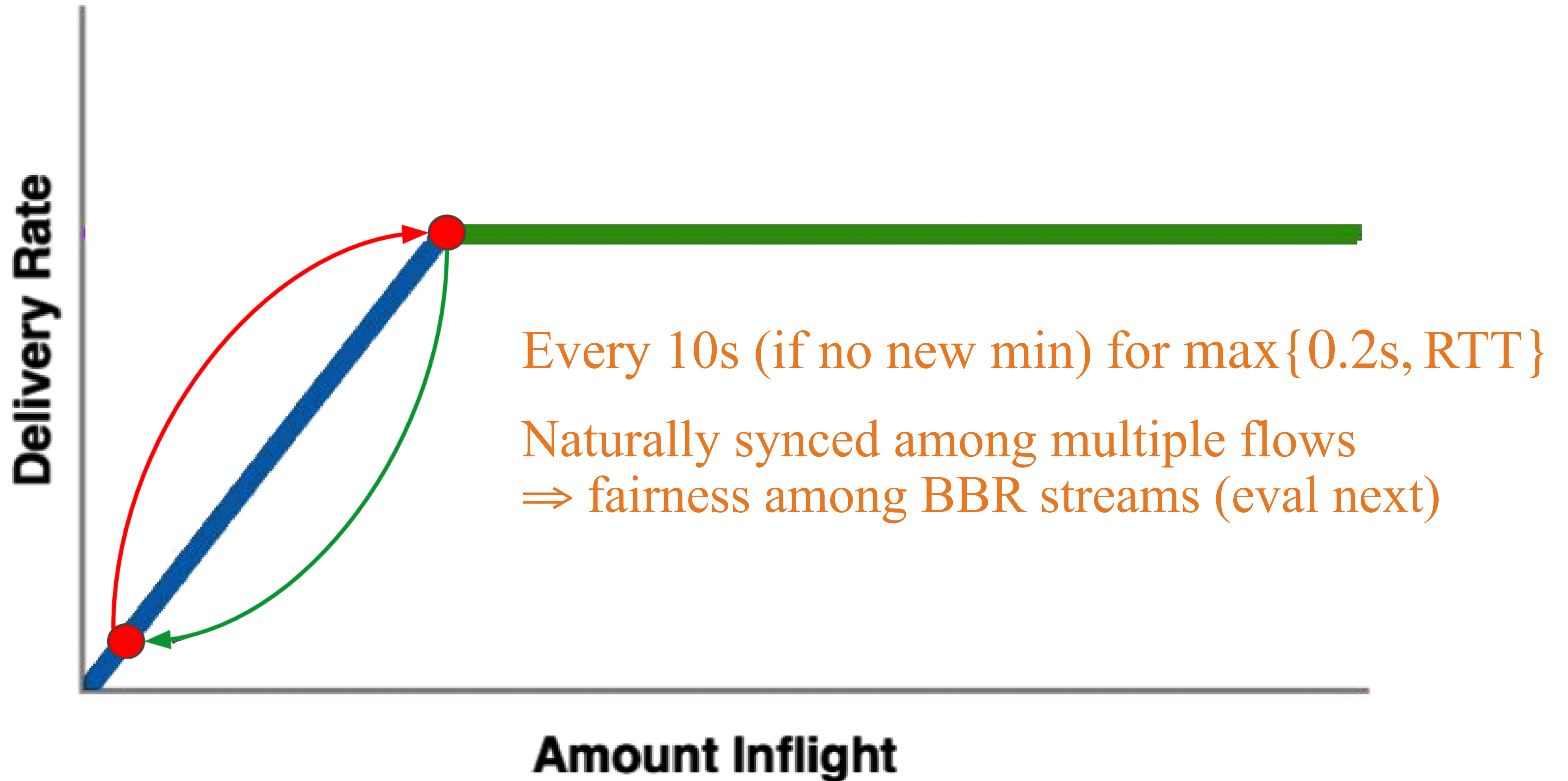
# Drain: Depleting queue (bounded by $2 \times \text{BDP}$ )



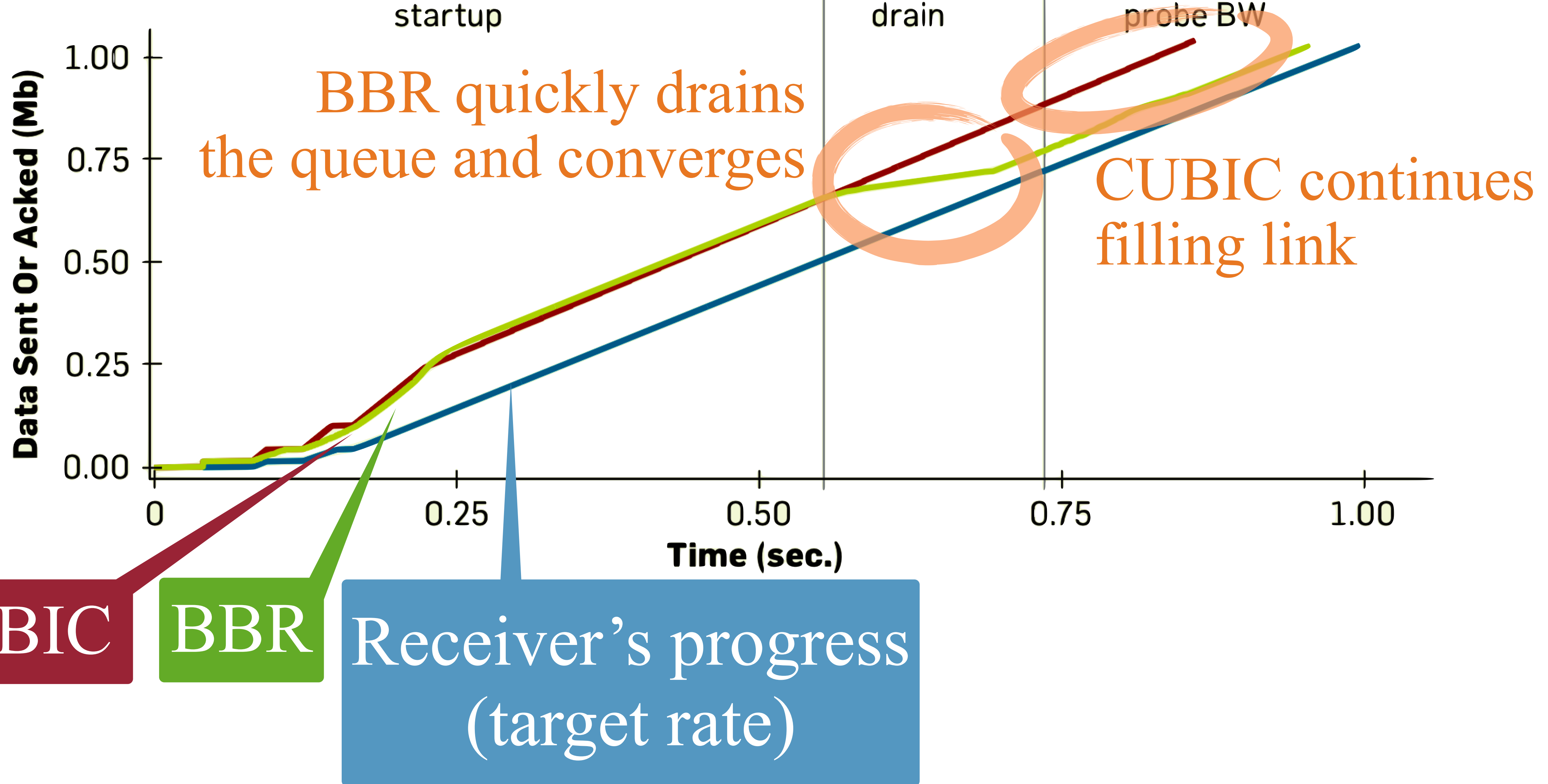
# Probe BW: Cycling Pacing Gain



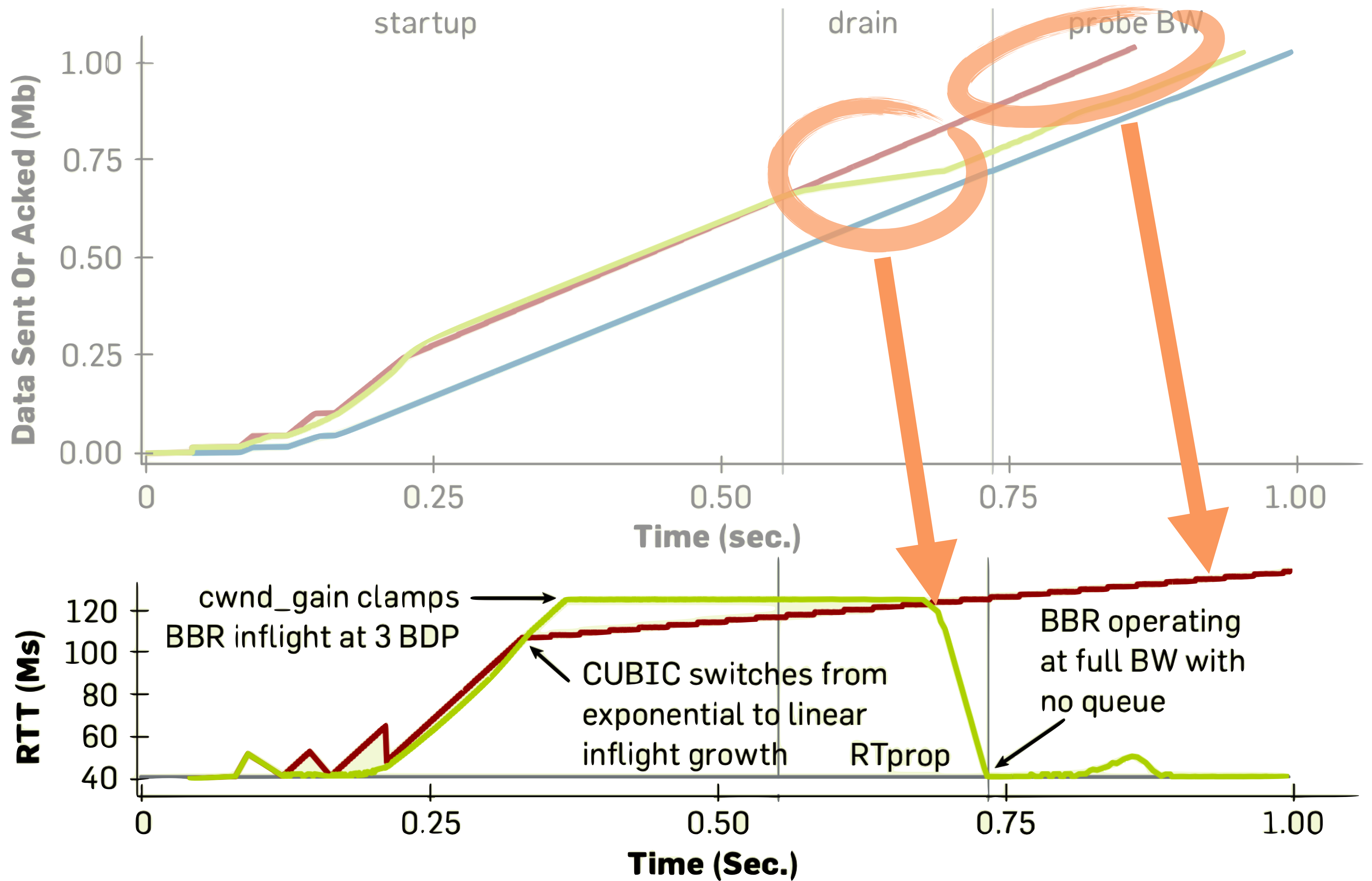
# Probe RTT: Periodic back off

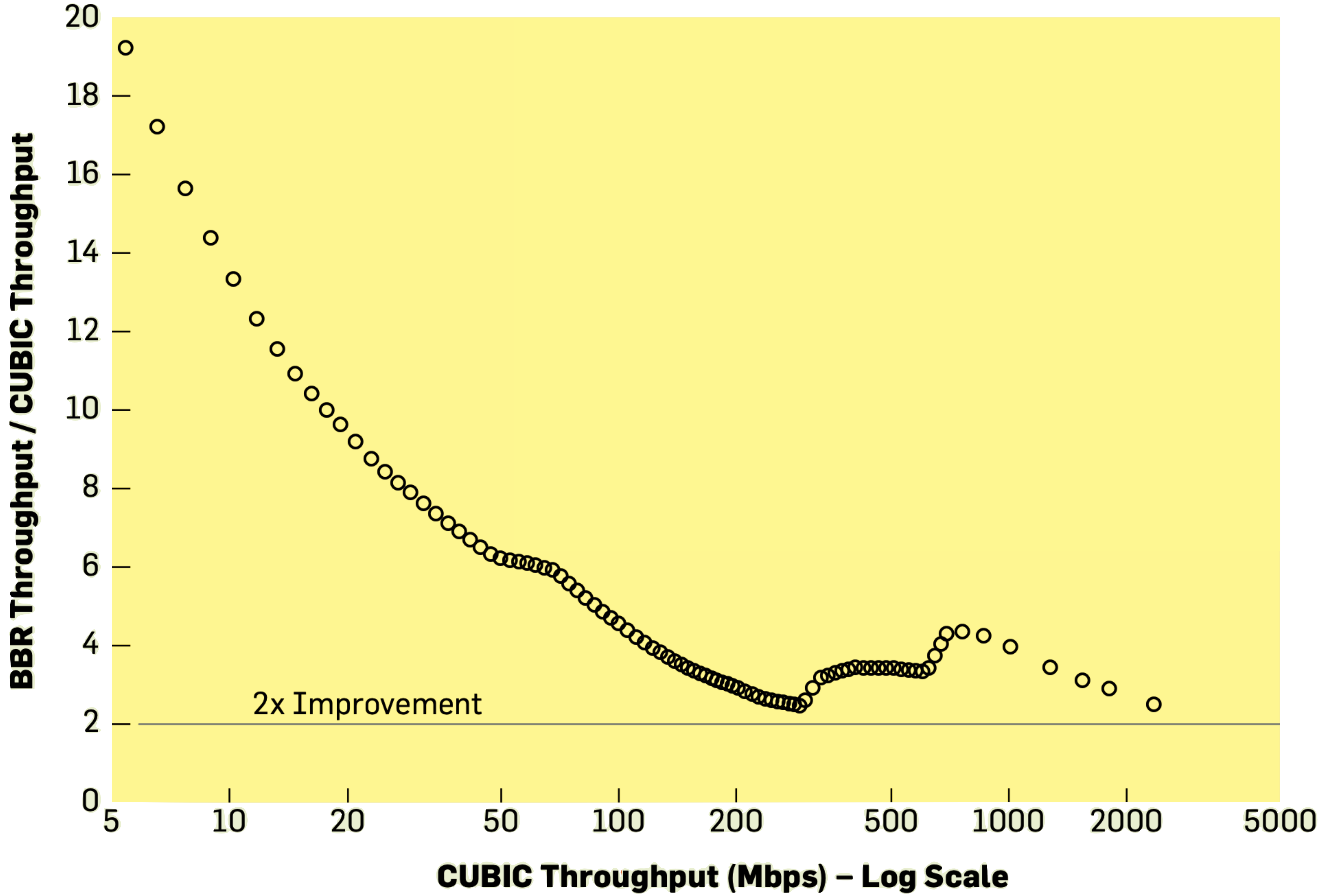


# Evaluation

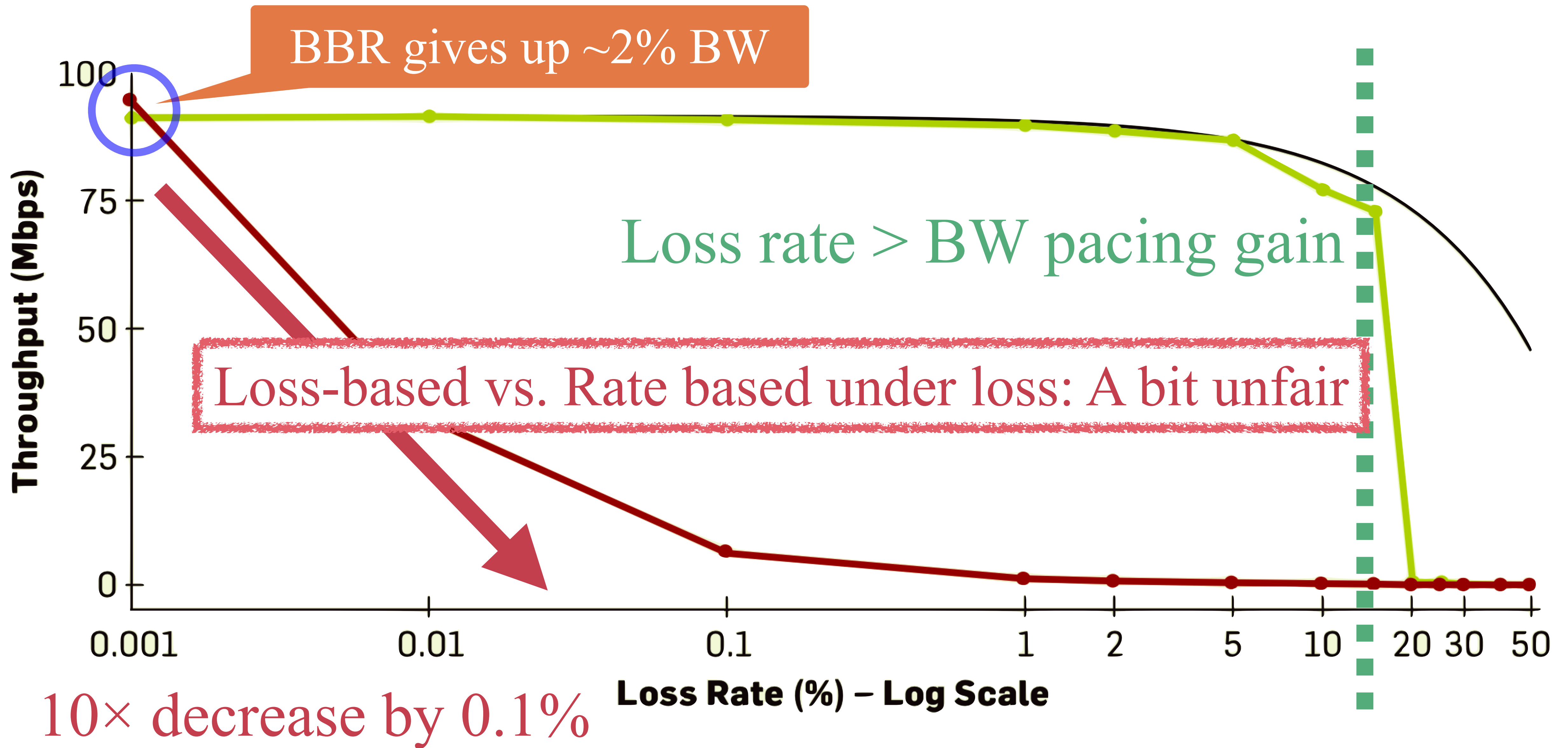








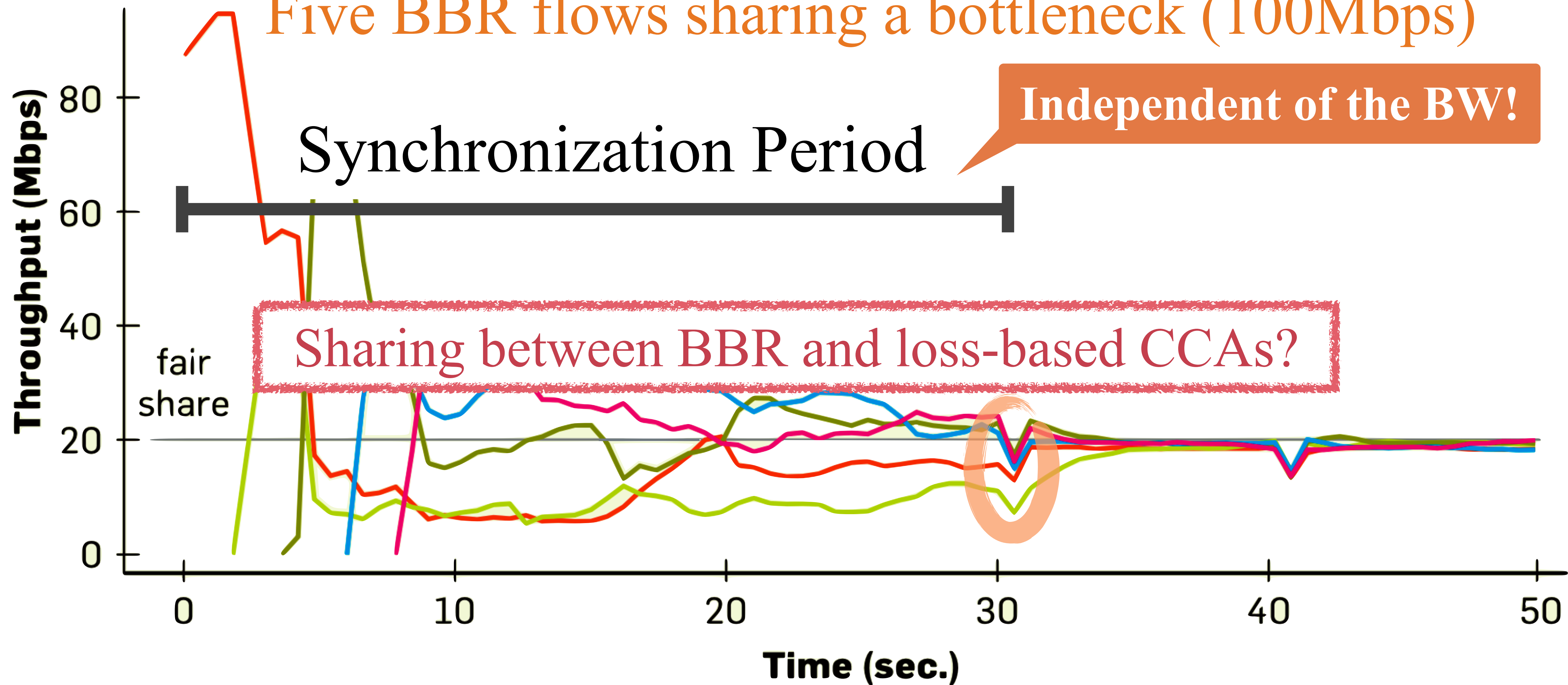
# Goodput under Loss **BBR** vs **CUBIC**





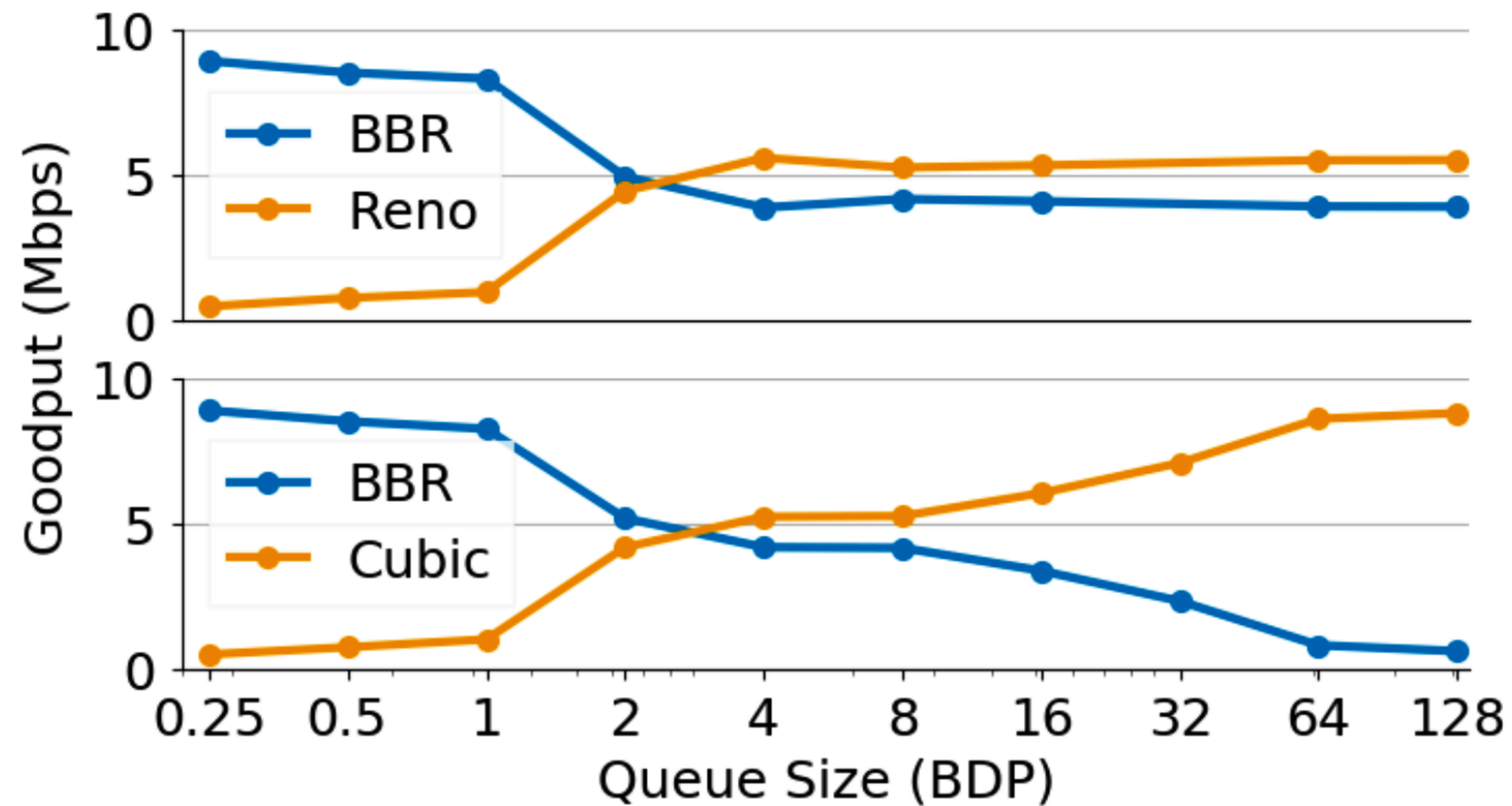
# Speaking of Fairness ...

Five BBR flows sharing a bottleneck (100Mbps)



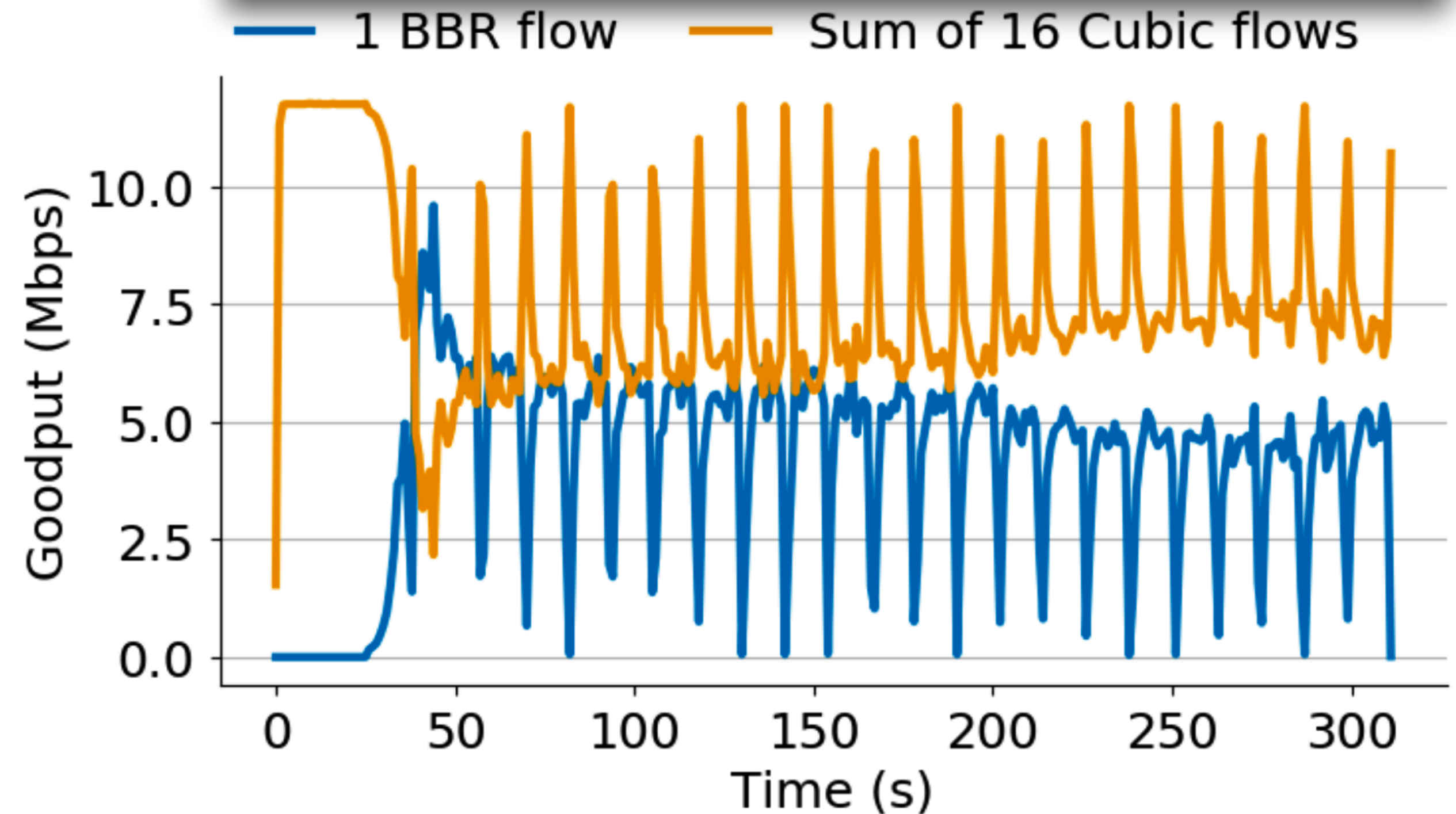
# BBR's (Un)Fairness [Ware et al., IMC '19]

BBR is generous



(b) Average goodput for two competing flows over 4 min in a  $40\text{ms} \times 10\text{Mbps}$  network with varying queue sizes.

BBR is aggressive



(c) BBR's goodput over time competing with 16 Cubic flows in a  $40\text{ms} \times 10\text{Mbps}$  network with a 32 BDP queue.

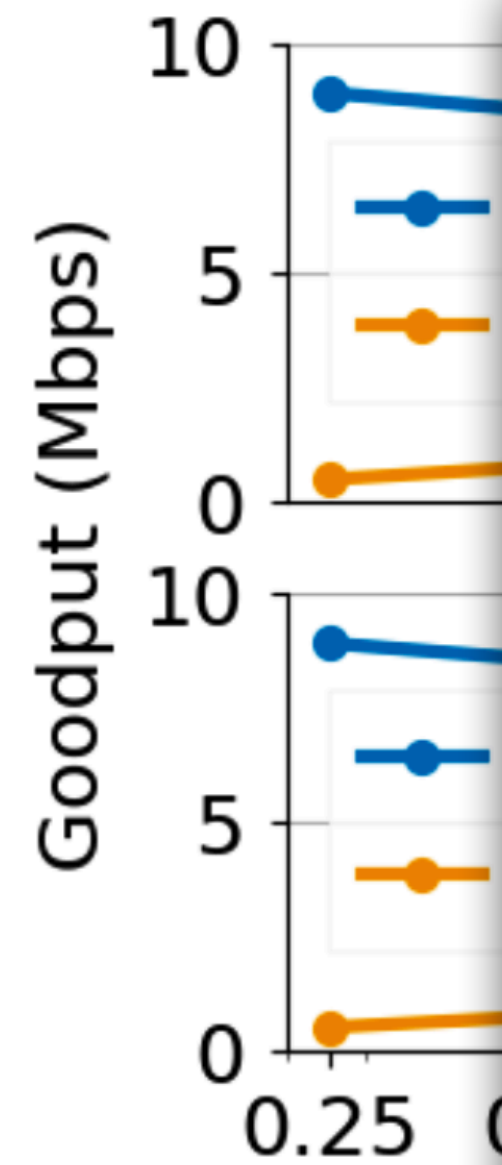


# BBR's (Un)Fairness [Ware et al., IMC '19]

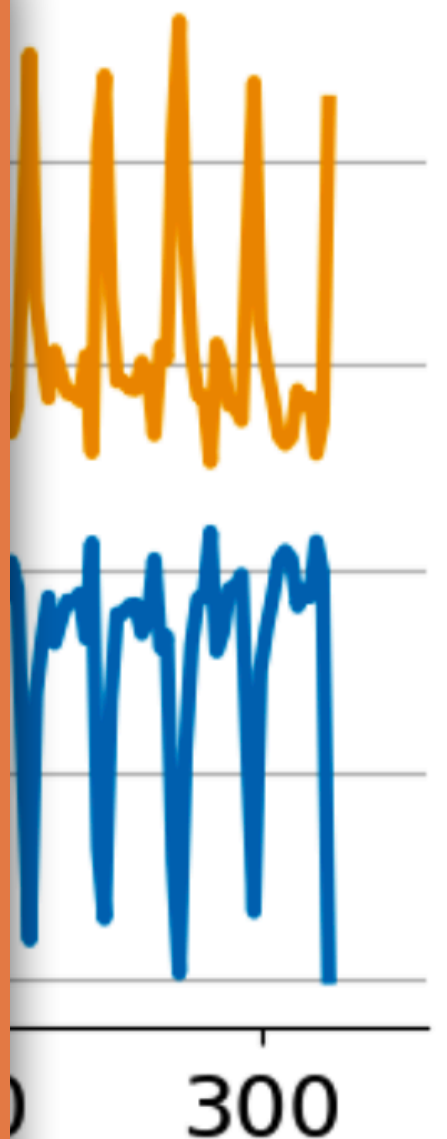
BBR is generous

BBR is aggressive

— 1 BBR flow — Sum of 16 Cubic flows



- BBR has an inflight cap ( $2 \times \text{BDP}$ ), a safety cap  $\rightarrow$  delayed/aggregated ACKs
- It dictates BBR's behaviors
- The # of flows doesn't change the cap (More in [Ware et al.] )



(b) Average goodput of competing flows over 4 min in a  $40\text{ms} \times 10\text{Mbps}$  network with varying queue sizes.

with 16 Cubic flows in a  $40\text{ms} \times 10\text{Mbps}$  network with a 32 BDP queue.



# BBR's Deployment at Google (~2017)

- YouTube: deployed for a small percentage of users
- Internal: test programs for Google data-centers
  - deployed as default TCP congestion control for internal Google traffic

Compared to CUBIC:

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

# Related and Future Work

- New versions of BBR (v2, v3 draft) → backup slides
  - v2 explicitly bounds loss rate
- Other rate based CCAs, e.g., Copa [NSDI '18] (next session)
- BBR's (un)fairness [NSDI '18; NotNets '19; IMC '19]
- Optimizing BBR's retransmission [Bi et al, ATC '23]

# My Review

## Strengths

S1. A performant, scalable solution to a fundamental challenge

S2. Solid evaluation, long-term development, and high impact

## Weaknesses

W1. Didn't push BBR over the edge in the paper (e.g., scaling/unfairness)

W2. Handling token bucket policers [BBRv2], TSO, and middle boxes

W3. (minor) None of the figures have legends :/

## Future directions

D1. Security vulnerabilities?

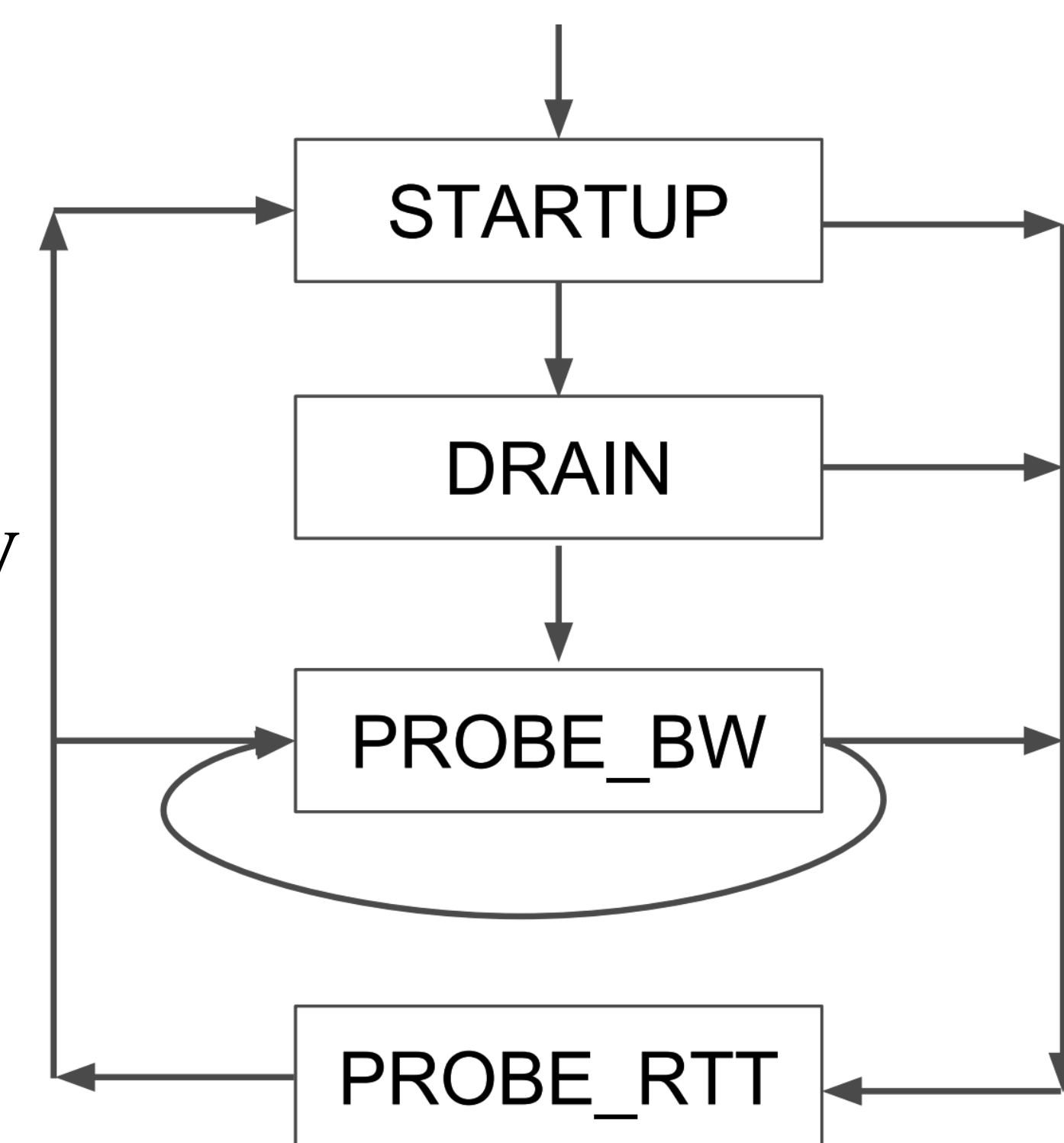
D2. Wireless environments (e.g., fairness and work with MPTCP)?

# Class Discussion

**Backup slides ...**

# BBR's FSM

- **STARTUP**: exponential growth to quickly fill pipe
  - stop growth when BW estimate plateaus, not on loss or delay
  - $\text{pacing\_gain} = 2.89$
- **DRAIN**: drain the queue created in STARTUP
  - $\text{pacing\_gain} = 0.35 = 1/2.89$
- **PROBE\_BW**: cycle pacing gain to explore and fairly share bandwidth
  - $[1.25, 0.75, 1, 1, 1, 1, 1, 1]$  (1 phase per min RTT)
  - $\text{pacing\_gain} = 1.25 \geq$  probe for more BW
  - $\text{pacing\_gain} = 1.0 \geq$  cruise with full utilization and low, bounded queue
- **PROBE\_RTT**: if needed, occasionally send slower to probe min RTT
  - $\text{pacing\_gain} = 0.75 \Rightarrow$  drain queue and yield BW to other flows

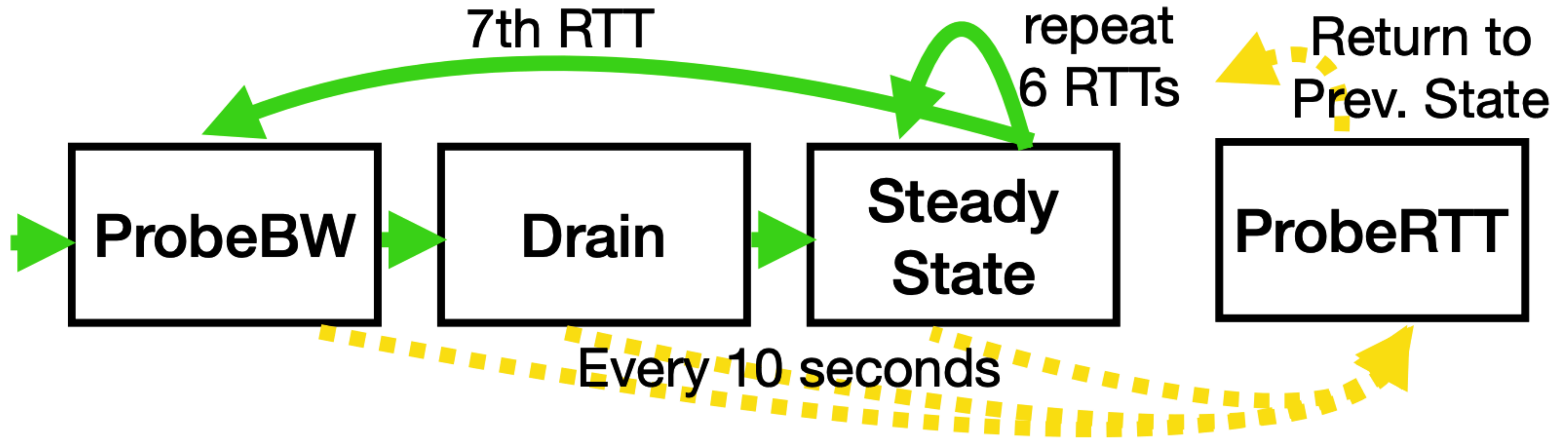




# BBR v2

	<b>CUBIC</b>	<b>BBR v1</b>	<b>BBR v2</b>
Model parameters to the state machine	N/A	Throughput, RTT	Throughput, RTT, max aggregation, max inflight
Loss	Reduce cwnd by 30% on window with any loss	N/A	Explicit loss rate target
ECN	<a href="#">RFC3168</a> (Classic ECN)	N/A	DCTCP-inspired ECN
Startup	Slow-start until RTT rises (Hystart) or any loss	Slow-start until tput plateaus	Slow-start until tput plateaus or ECN/loss rate > target

# BBR's FSM (2)



**Figure 4: BBR's steady-state operation.**

[Ware et al., IMC '19]