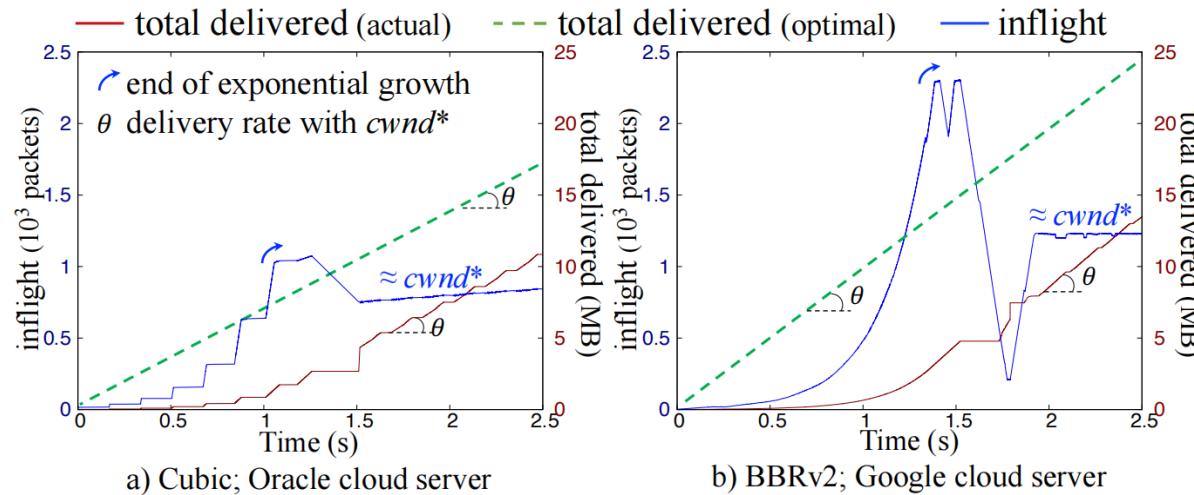


SUSS: Improving TCP Performance by Speeding Up Slow-Start

Mahdi Arghavani, Haibo Zhang, David Eyers, and Abbas
Arghvani *SIGCOMM 2024*

TCP slow-start drawbacks

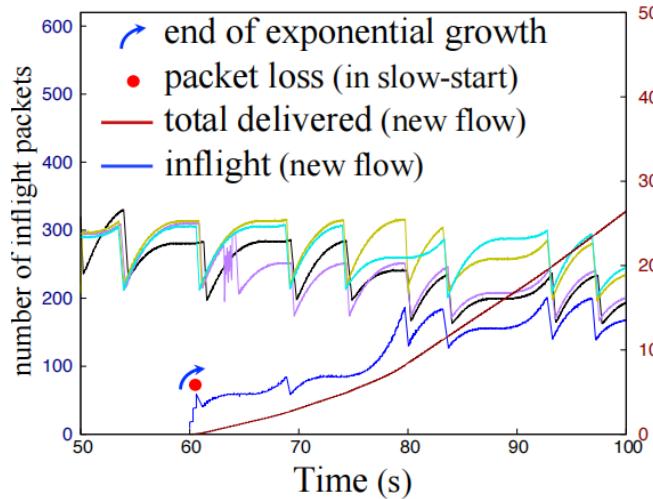
- slow ramping-up of the data delivery rate
 - inefficient bandwidth utilization
- prolonged completion time for small-size flows (flows less than a few megabytes in size)
 - significantly diminish the Quality of Experience (QoE)



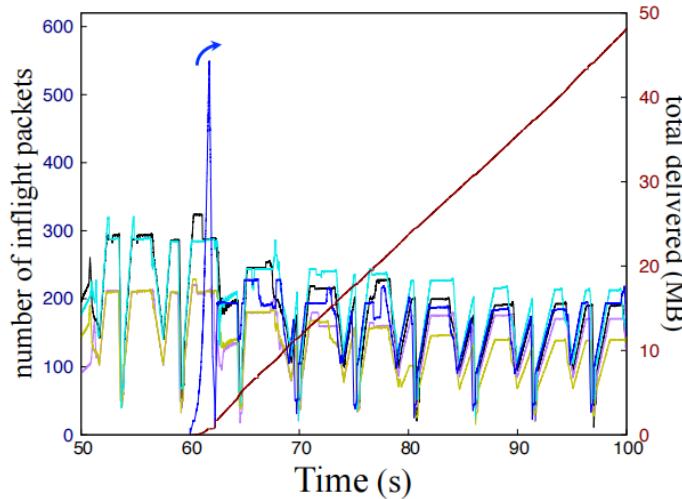
- Two plots highlight the capability for transferring more data during the initial phase

TCP slow-start drawbacks

- Inefficiency of CUBIC's slow-start in a congested network path



a) Cubic (loss-sensitive slow-start)



b) BBRv2 (loss-tolerant slow-start)

- CUBIC flow cannot gain a fair share of network resources rapidly, resulting in prolonged unfairness
 - reason: sensitivity in packet loss → exit slow-start before reaching $cwnd^*$

Related Works

- network-assisted approaches
 - require feedback from routers and switches → challenging for deployment and testing
- End-to-end solutions
 - designed to estimate $cwnd^*$, based on bandwidth measurement and optimal rate estimation
 - accurate bandwidth measurement and rate estimation are very challenging, especially in early rounds
- intuitive thoughts: use a larger initial window size (iw) → now 15KB (10 segments, RFC6928)
 - unaware of the flow length, lack prior knowledge of the network path conditions at the beginning
 - cannot select iw appropriately
 - the consequences of releasing a large burst of data packets in early rounds

A lightweight, sender-side add-on to CUBIC's slow-start mechanism

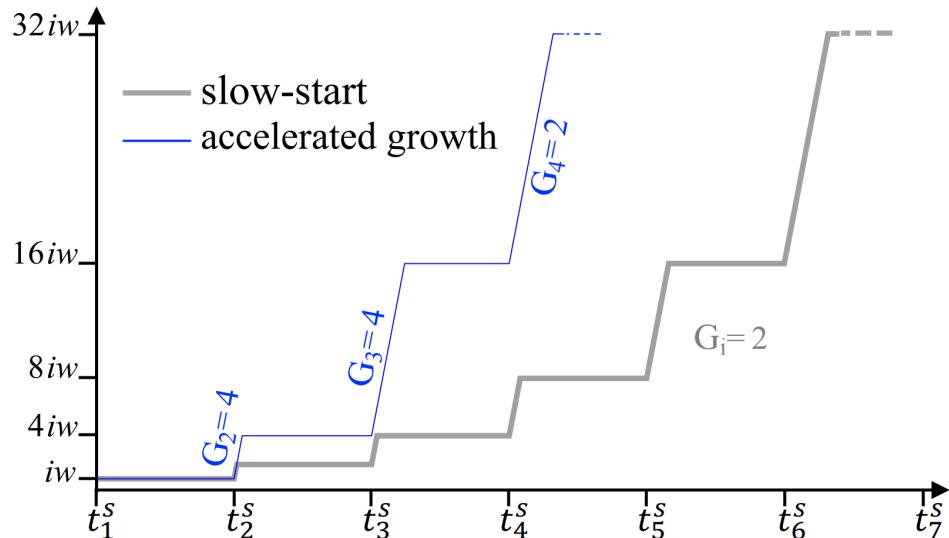
- Aims
 - safely accelerate cwnd growth in large-BDP networks
 - enhance the delivery rate during the early RTTs of slow-start
- Key idea: predict the continuation of exponential growth of cwnd in the subsequent RTT
 - if *cwnd* is expected to grow in the next RTT, then quadruple *cwnd*
 - allow *cwnd* to quickly ramp up in the early RTTs of slow-start
- primary challenges
 - How to predict the likelihood of exponential *cwnd* growth in the next RTT
 - based on HyStart: $cwnd \approx cwnd^*$
 - How to control additional data packet transmission to avoid bursts and packet loss
 - propose a novel combination of ACK clocking and packet pacing
 - ACK clocking: protect the functionality of HyStart and its estimation of the growth factor
 - packet pacing: avoid the bursts of traffic that might be caused by surges in *cwnd*

Key parameters

- $\text{round}(i)$: the time interval $[t_i^s, t_{i+1}^s)$, where t_i^s is the time when the ACK of the first data packet sent in $\text{round}(i - 1)$ is received for $i > 1$. For $i = 1$, $\text{round}(1)$ initiates at the beginning of the data transfer.
- $cwnd_i$: maximum congestion window size in $\text{round}(i)$.
- iw : volume of data sent in $\text{round}(1)$, i.e., $cwnd_1 = iw$.
- *Data train*: sequence of data packets sent in a round.
- *ACK train*: sequence of ACKs for a data train.
- Δt_i^{at} : time taken to receive the ACK train in $\text{round}(i)$.
- R_i : the transmission rate of the data train in $\text{round}(i)$.
- minRTT : measured minimum round-trip time since the initiation of the TCP connection.
- moRTT_i : the minimum observed RTT in $\text{round}(i)$.
- $BtlBw$: bandwidth of the bottleneck link.
- G_i : growth factor for $cwnd$ in $\text{round}(i)$.
- $cwnd^*$: the optimal congestion window size when a fair share of $BtlBw$ is achieved.

Theory Behind SUSS

- After receive the ACK train:
 - evaluate whether the exponential growth of $cwnd$ is extrapolated to persist for the next round
 - if true, then send another $2 * cwnd(i - 1)$ this round $\rightarrow G_i = 4$
 - Why not consider more rounds? the network condition may vary over several upcoming rounds



Key Challenge 1

How to determine whether the exponential growth of cwnd is expected to continue in the next round?

→ implement the approach from HyStart

- HyStart
 - Upon receiving an ACK, the amount of time elapsed from the start of the current round does not exceed half of the minRTT
 - The minimum observed RTT in the current round ($moRTT_i$) must not be greater than $1.125 * \text{minRTT}$.
- Condition 1 for SUSS
 - $\Delta t_{i+1}^{at} \leq \text{minRTT}/2$
 - need to estimate Δt_{i+1}^{at}
- Condition 2 for SUSS
 - $moRTT_{i+1} \leq 1.125 * \text{minRTT}$
 - need to estimate $moRTT_{i+1}$

Formulating Condition 1 for SUSS

$$\Delta t_{i+1}^{at} \leq minRTT/2$$

- Δt_{i+1}^{at} is not known in round(i) → estimate it by its relationship with Δt_i^{at}
 - $\Delta t_i^{at} = \frac{cwnd_{i-1}}{BtlBw}$
- As the same as slow-start, $2 * cwnd_{i-1}$ is sent during Δt_i^{at}
 - $\Delta t_{i+1}^{at} = \frac{2*cwnd_{i-1}}{BtlBw} = 2 * \Delta t_i^{at}$
- **Condition 1:** $\Delta t_i^{at} \leq minRTT/4$

Formulating Condition 2 for SUSS

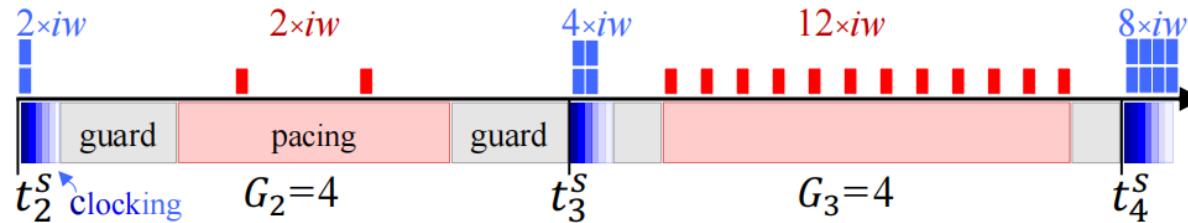
$$moRTT_{i+1} \leq 1.125 * minRTT$$

- $1.125 * minRTT \rightarrow$ initial signs of queueing delay \rightarrow beginning of congestion
- suppose that $minRTT$ was last updated r rounds ago
 - the average increase in queuing delay since then is approximately $\frac{moRTT_i - minRTT}{r}$
 - $moRTT_{i+1} \approx moRTT_i + \frac{moRTT_i - minRTT}{r}$
- **Condition 2:** $moRTT_i + \frac{moRTT_i - minRTT}{r} \leq 1.125 * minRTT$

Key Challenge 2

How to control additional data packet transmission to avoid bursts and packet loss?

- Combination of ACK clocking and packet pacing when $G_i > 2$
 - only ACK clocking → a burst of packets upon receiving each ACK train
 - only packet pacing → the measurement of Δt_i^{at} will no longer be accurate



- In *clocking period*, SUSS always sends twice the amount of data acknowledged by that ACK
 - facilitate the measurement of Δt_i^{at}
- In *pacing period*, the remaining data will be sent at a predetermined pace to mitigate burstiness
- guard period* → avoid interference with the clocking period in both the current and the next round

3.1 Clocking Period

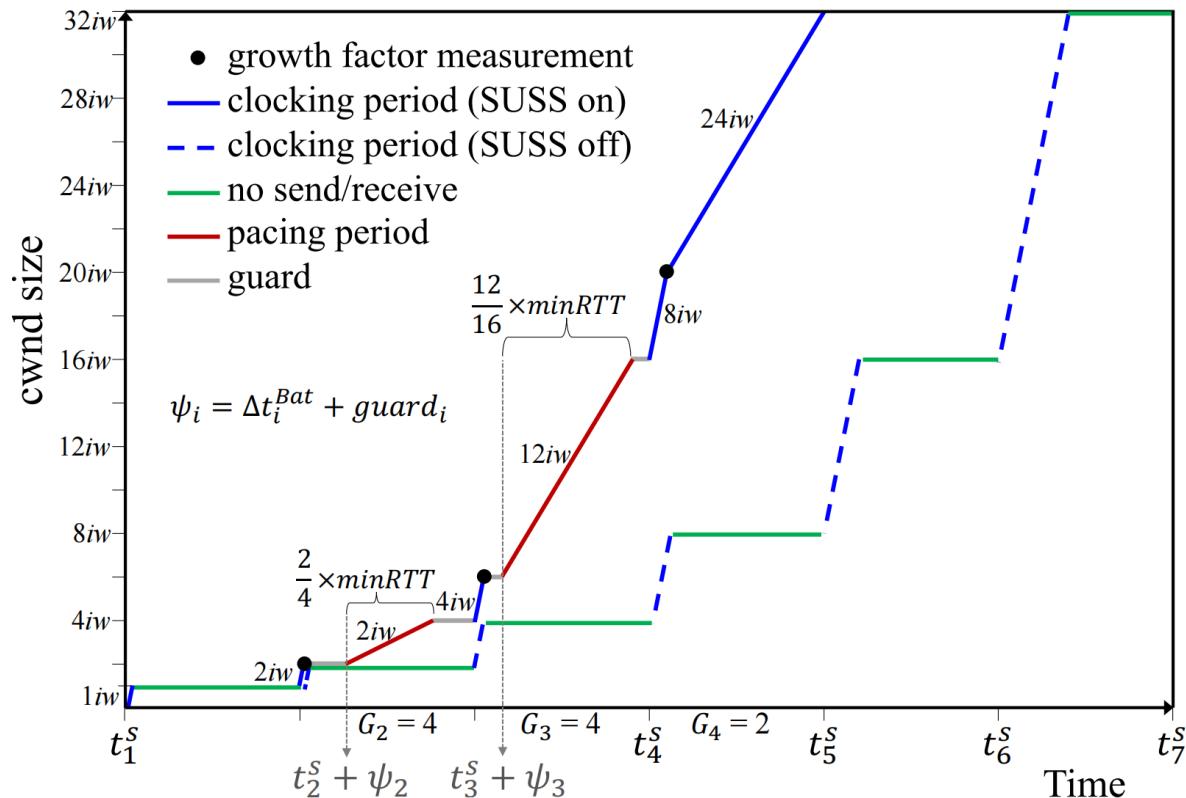
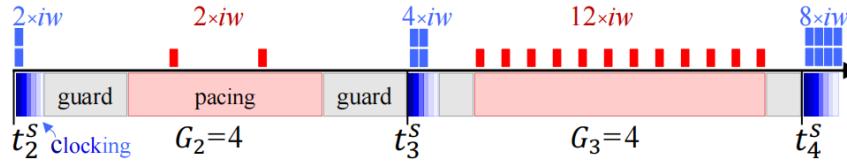
- $S_i^{Bdt} = iw * 2^{i-1}$: the amount of data sent in the clocking period for round(i)
- Δt_i^{Bat} : the amount of time taken to receive the ACKs for the blue packets send in round(i-1)
- $\Delta t_i^{at} = cwnd_{i-1} * \frac{\Delta t_i^{Bat}}{S_{i-1}^{Bdt}}$

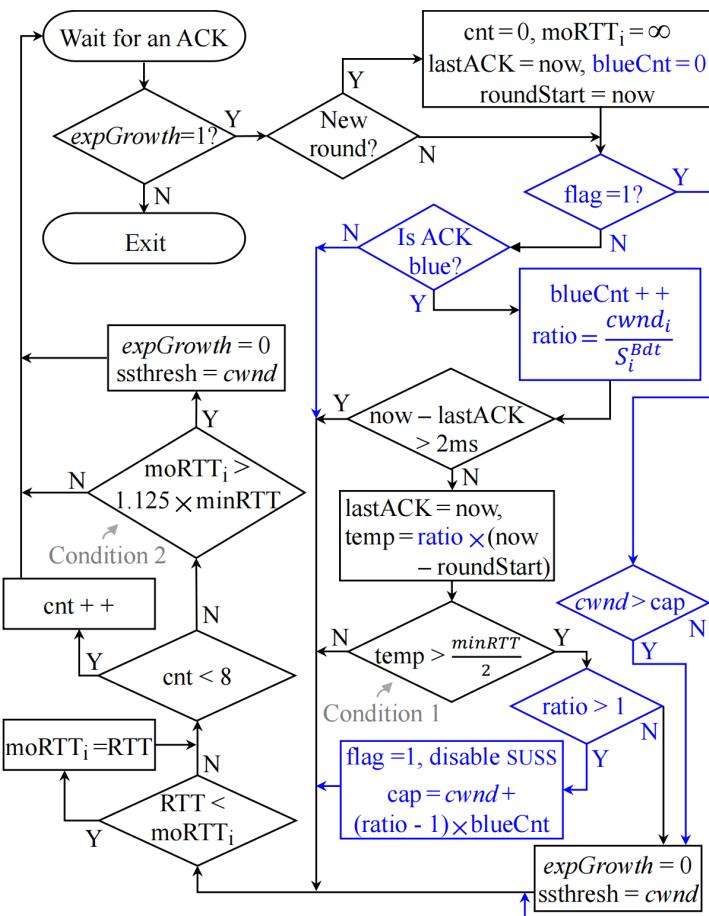
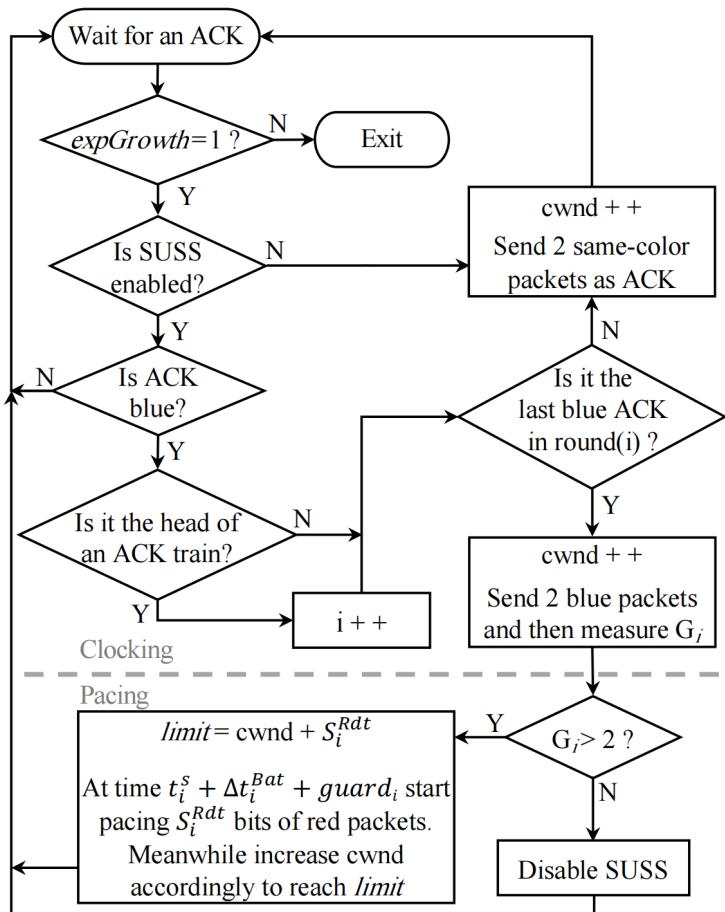
3.2 Pacing Period

- S_i^{Rdt} : the amount of data to be sent during the pacing period of round(i)
 - $cwnd_i = S_i^{Bdt} + S_i^{Rdt}$
 - $S_i^{Rdt} = G_i * (S_{i-1}^{Bdt} + S_{i-1}^{Rdt}) - S_i^{Bdt} = G_i * S_{i-1}^{Rdt} + (G_i - 2) * 2^{i-2} * iw$
 - the growth of S_i^{Rdt} is higher than S_i^{Bdt} when $G_i > 2$
 - the duration and starting time of the pacing period in each round need to be dynamic
- time of pacing period: $\frac{S_i^{Rdt}}{cwnd_i} * minRTT \rightarrow$ sending rate: $rate_i = \frac{S_i^{Rdt}}{\frac{S_i^{Rdt}}{cwnd_i} * minRTT} = \frac{cwnd_i}{minRTT}$

3.3 Guard Period

$$\begin{aligned} guard_i &= \frac{minRTT - \left(\Delta t_i^{Bat} + \frac{S_i^{Rdt}}{cwnd_i} * minRTT \right)}{2} \\ &= \frac{S_i^{Bdt}}{2 * cwnd_i} * minRTT - \frac{\Delta t_i^{Bat}}{2} \end{aligned}$$

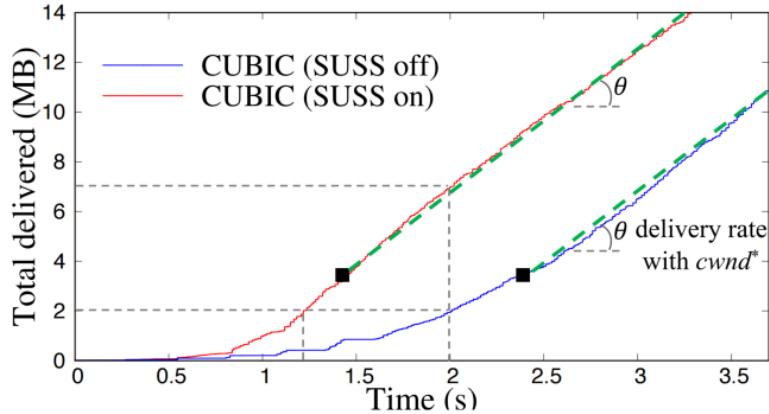
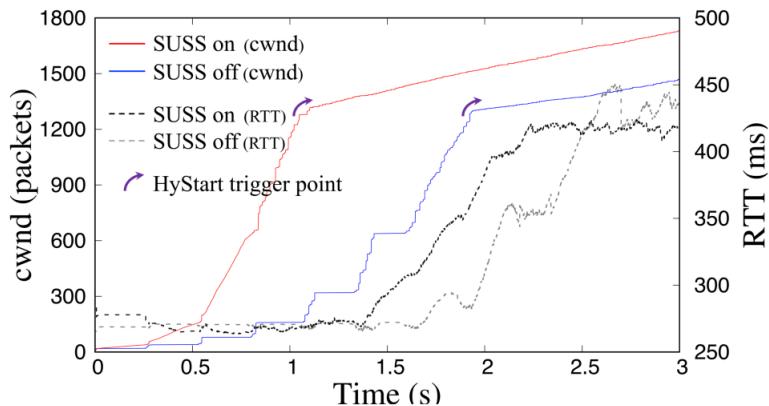




Performance Evaluation

- server deployment: Google server(BBR), Oracle server(CUBIC), both 3 servers located in different places
- client deployment: SUSS on server, different kinds of clients in NZ or Sweden

Throughput

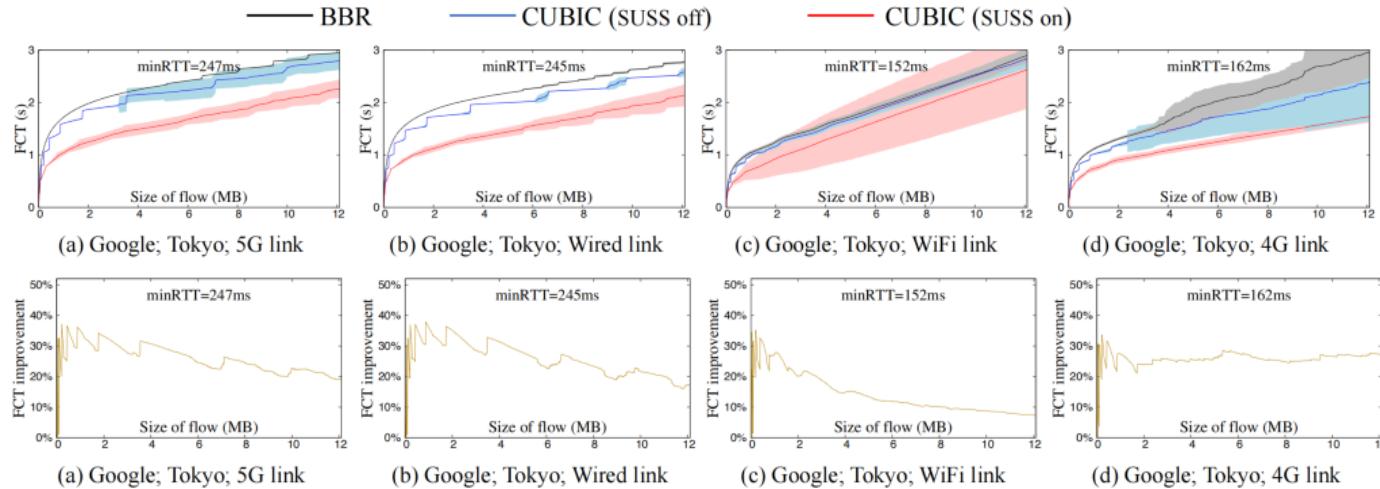


- CUBIC with SUSS exhibits a faster and smoother increase in cwnd and does not incur extra end-to-end delay
- SUSS has brought substantial improvements
 - SUSS can reach $cwnd^*$ more quickly than the traditional slow-start

Performance Evaluation

FCT (Flow Complete Time)

Small Flow

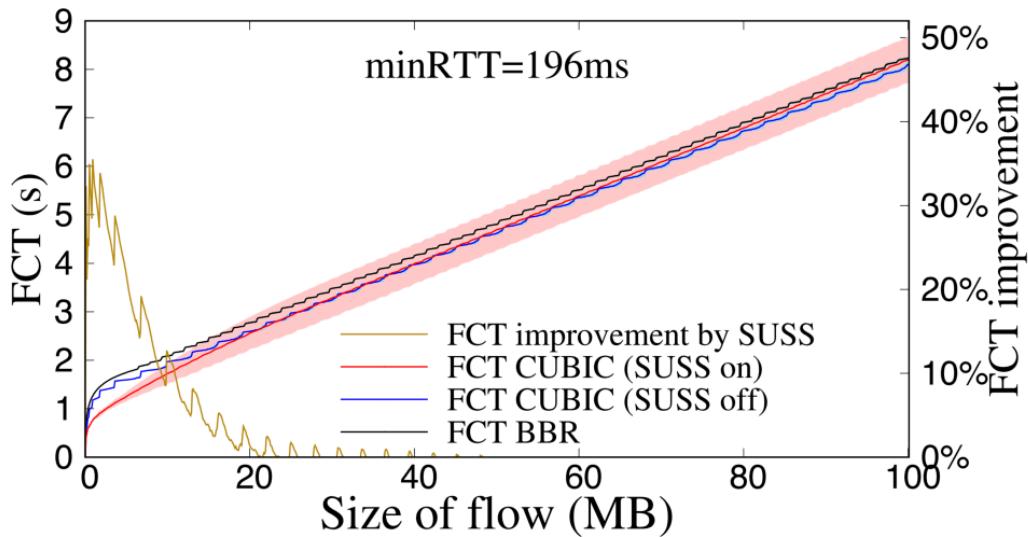


- accelerating slow-start can yield significant gain for small flows
 - as many small flows predominantly reside within the slow-start phase

Performance Evaluation

FCT (Flow Complete Time)

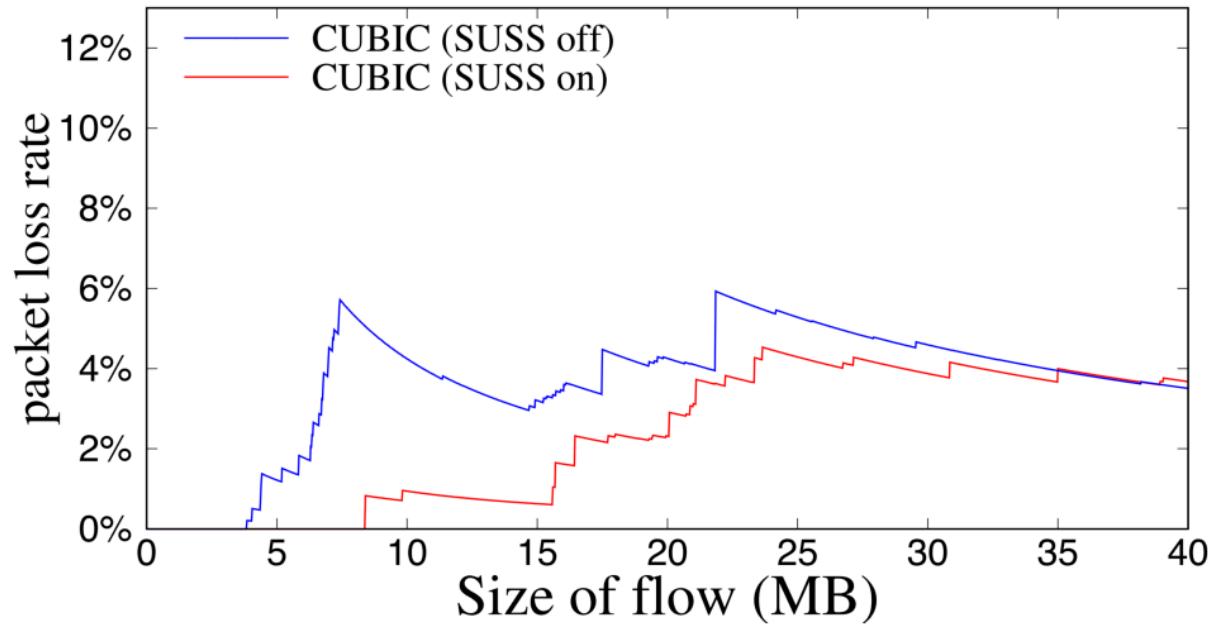
Large Flow



- SUSS enhances the performance of small-sized flows, it does not increase cwnd beyond its optimal value, nor does it impact the FCT of larger flows

Performance Evaluation

Packet Loss

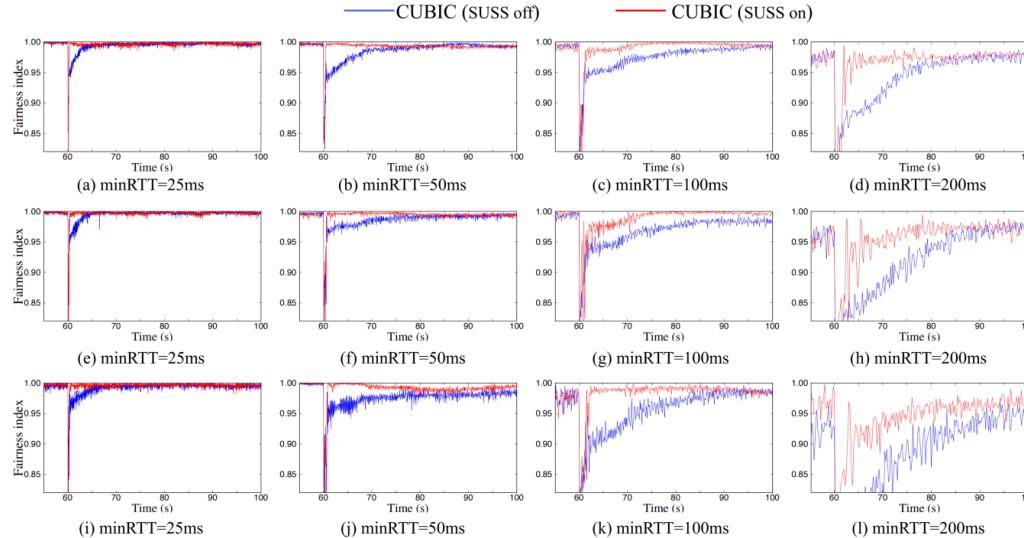


- Pacing can significantly reduce packet density and thereby reduce the packet loss rate
- It will converge when the size of flow increase

Performance Evaluation

Fairness

- $F = \frac{\sum(x_i)^2}{n * \sum(x_i^2)}$, n is the number of flows and x_i is the goodput of the i-th flow $\rightarrow F \rightarrow 1$, fairness \uparrow

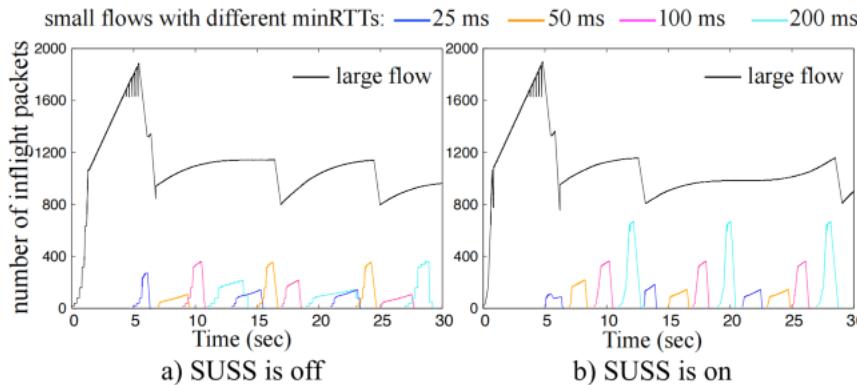


- when a level of fairness is achieved among the four flows, the fifth flow starts downloading
- the initiation of the fifth flow results in an immediate decrease in the value of F , exhibiting a prolonged delay in approaching 1 when SUSS is off

Performance Evaluation

Stability

- A large flow facing the initiation of multiple concurrent small flows



Buffer size	minRTT of large flow	SUSS off		SUSS on		FCT of large flow	FCT of small flows	SUSS off		SUSS on		FCT of large flow	FCT of small flows	SUSS off		SUSS on		Improvement for small flows
		FCT of large flow	FCT of small flows	FCT of large flow	FCT of small flows			FCT of large flow	FCT of small flows	FCT of large flow	FCT of small flows			FCT of large flow	FCT of small flows	FCT of large flow	FCT of small flows	
1 BDP	25 ms	30.4	1.33	31.2	0.94	29%		24.1	20.07	24.3	17.54	29%		28.5	1.44	29.2	1.04	29%
	50 ms	37.9	1.18	35.6	0.98	17%		24.6	13.64	25.2	8.28	17%		27.3	1.49	28.4	1.16	17%
	100 ms	29.0	1.42	28.5	1.22	14%		25.7	8.11	26.2	4.30	14%		27.0	1.69	27.2	1.34	14%
	200 ms	27.2	2.69	26.9	2.03	25%		27.6	4.43	28.2	2.52	25%		27.5	3.13	28.0	1.83	25%
Average		28.8	2.28	28.2	1.55	32%		25.7	7.40	26.0	5.32	28%		27.0	2.35	27.6	1.73	26%

(a) The large flow employs CUBIC

(b) The large flow employs BBRv1

(c) The large flow employs BBRv2

- SUSS improves FCT of small CUBIC flows without compromising the stability of the large flow

Thanks for Listening

[Paper](#) · [GitHub](#) · [Conference video](#)