# A Network Calculus Model for Congestion Control in Data Center Network

Natchanon Luangsomboon
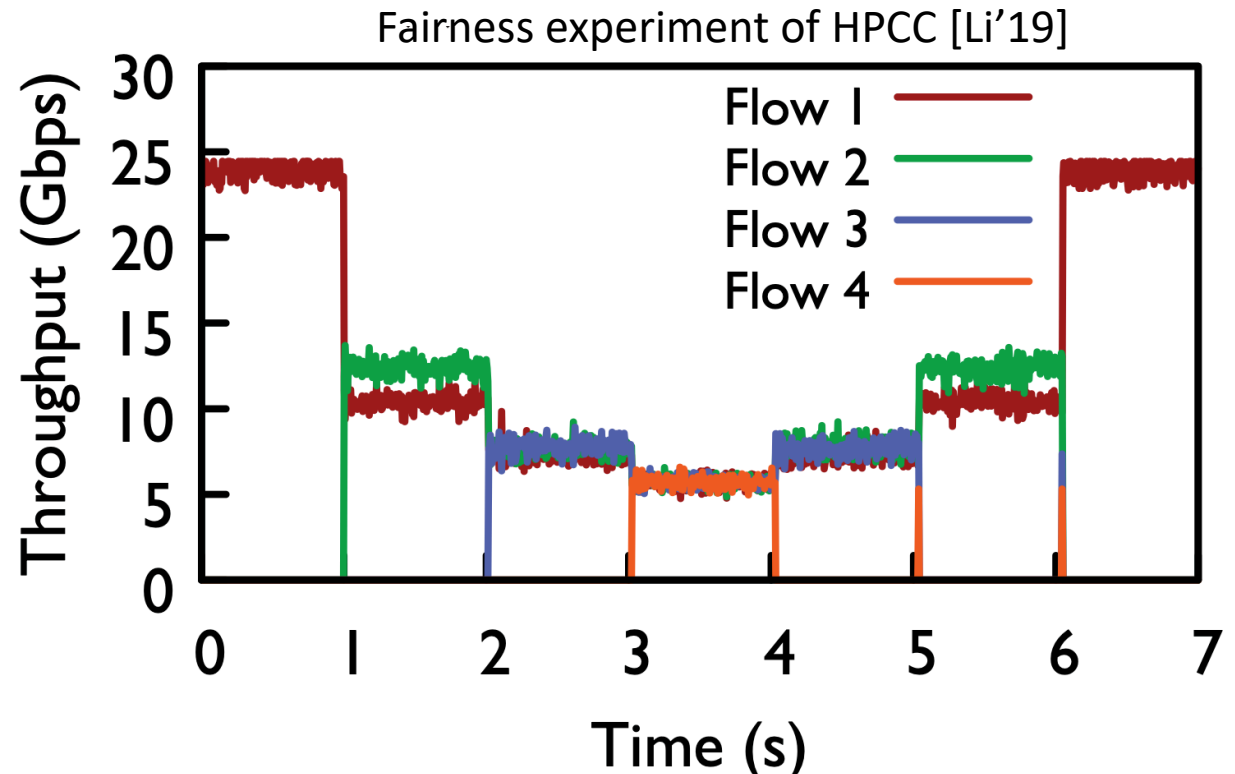
The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
UNIVERSITY OF TORONTO

# This talk

- Use network calculus to model data center network as multi-flow window flow control

- Analyze its bounds on transmission rate and fairness given the time-variable congestion windows

- Goal: a congestion control algorithm (CCA) that is fair and fully utilizes the link bandwidth based on the multi-flow model
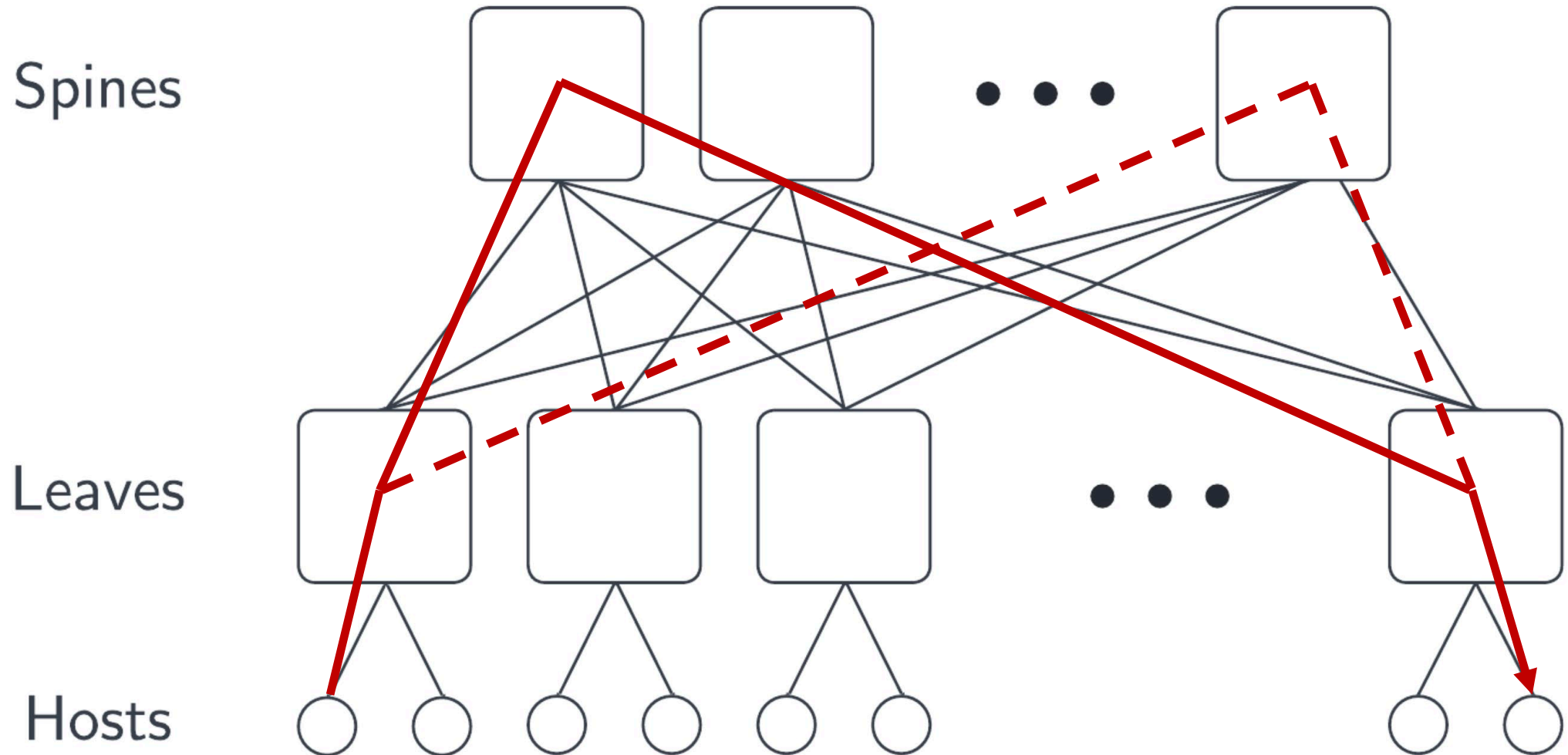
# Fairness is an important aspect of CCAs

- The fainess experiment is a common test for CCAs

- BBRv1 is criticized for being unfair when used with other CCAs, e.g., CUBIC, Reno
  - Its inter-protocol fairness has been improved in subsequent versions
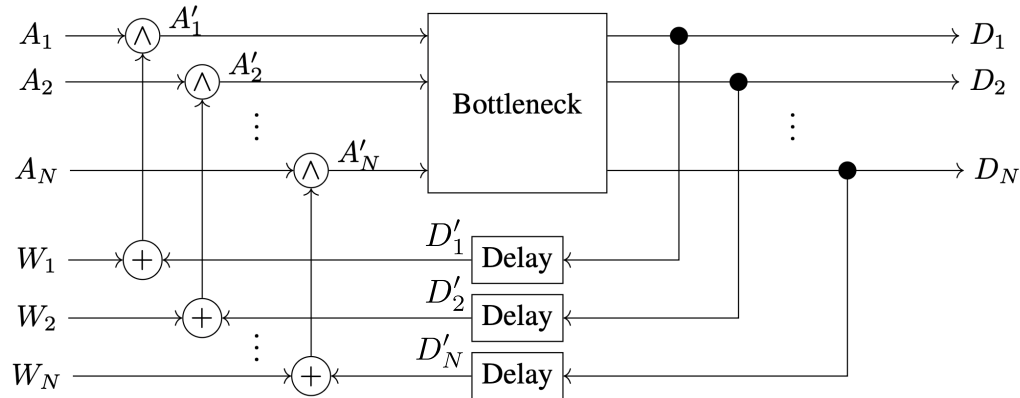
Fairness experiment of HPCC [Li'19]

# Related works

- Fairness and Stability of End-to-End Congestion Control [Kelly'03]
- Congestion Control for Large-Scale RDMA Deployments (DCQCN) [Zhu'15]
- Model-based Insights on the Performance, Fairness, and Stability of BBR [Scherrer'22]
- Toward Stability Analysis of Data Transport Mechanisms: A Fluid Model and Application (Reno & CUBIC) [Vardoyan'18]
- …

# Data center network has homogeneous delay

# Network model and notations



$$N_{\mathrm{net}}(t) := \{i \in N \mid B_i(t) = W_i(t)\}$$

For each flow ,

- $A_i(t)$: Available data from the application
- $A'_i(t)$: Data transmitted from the sender
- $D_i(t)$: Data arrived at the receiver
- $D'_i(t)$: Acknowledgement
- $B_i(t)$: Inflight bytes
- $W_i(t)$: Congestion window
- : feedback delay
- : List of all flows
- $N_{\mathrm{net}}(t)$: List of network limited flows

# Additional bottleneck switch assumptions

- The bottleneck switch transmits data in a FIFO order, i.e., for all ,

$$A'(s) \le D(t) \iff \forall i, A'_i(s) \le D_i(t)$$

- The switch transmits one packet at a time

$$\forall s \exists t, \ A'(s) = D(t) \qquad \text{and} \qquad \forall s \exists t, A'_i(s) = D_i(t)$$

<span style="color:red">Swith transmits at s</span>

<span style="color:red">Remains idle for ε</span>

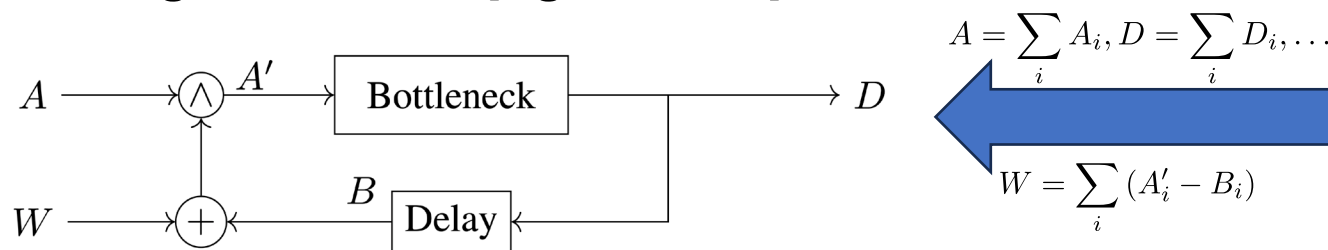- Each transmission is at least apart

$$\forall s, D(s) < \lim_{v \to s^+} D(v) \implies \lim_{v \to s^+} D(v) = D(s + \epsilon)$$
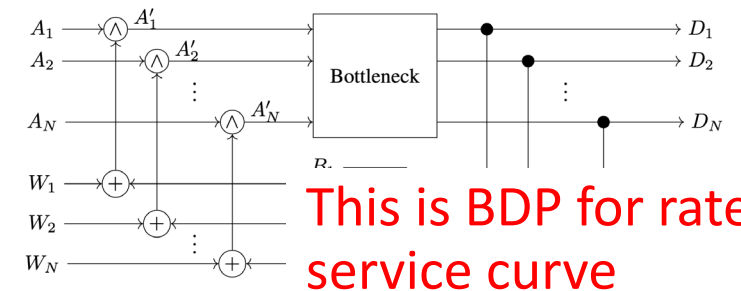
# Minimum $W_i$ to "fill the pipe" (BDP)

Given upper/lower service curves at the switch, what is the minimum $W_i$ s.t. the entire network has the same upper/lower service curves?

Single-flow case [Agrawal'99]

Multi-flow case



$$A = \sum_i A_i, D = \sum_i D_i, \dots$$

$$W = \sum_i (A_i' - B_i)$$

This is BDP for rate service curve

$$\forall t, W(t) \geq \sum_{t \in \mathbb{R}} \{\underline{S}(t) - \underline{S} \otimes \underline{S}(t-d)\} \implies \underline{S} = \underline{S}_{\mathrm{net}}$$

$$\forall t, W(t) \geq \sum_{t \in \mathbb{R}} \{\overline{S}(t) - \overline{S} \otimes \overline{S}(t-d)\} \implies \overline{S} = \overline{S}_{\mathrm{net}}$$
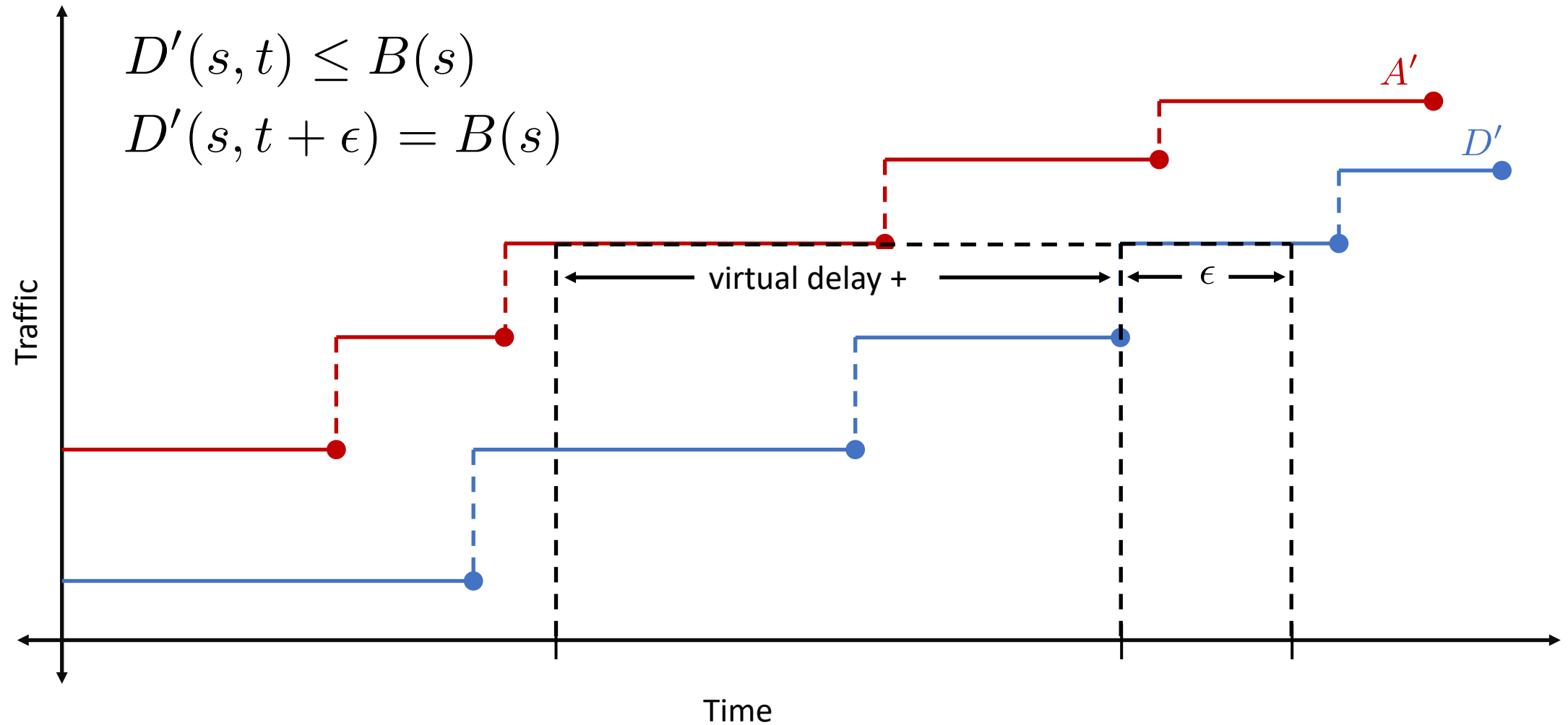
# Max-min fairness

- Each flow   is associated with a weight $\phi_i$
- If two flows  ,   are network-limited throughout  an interval

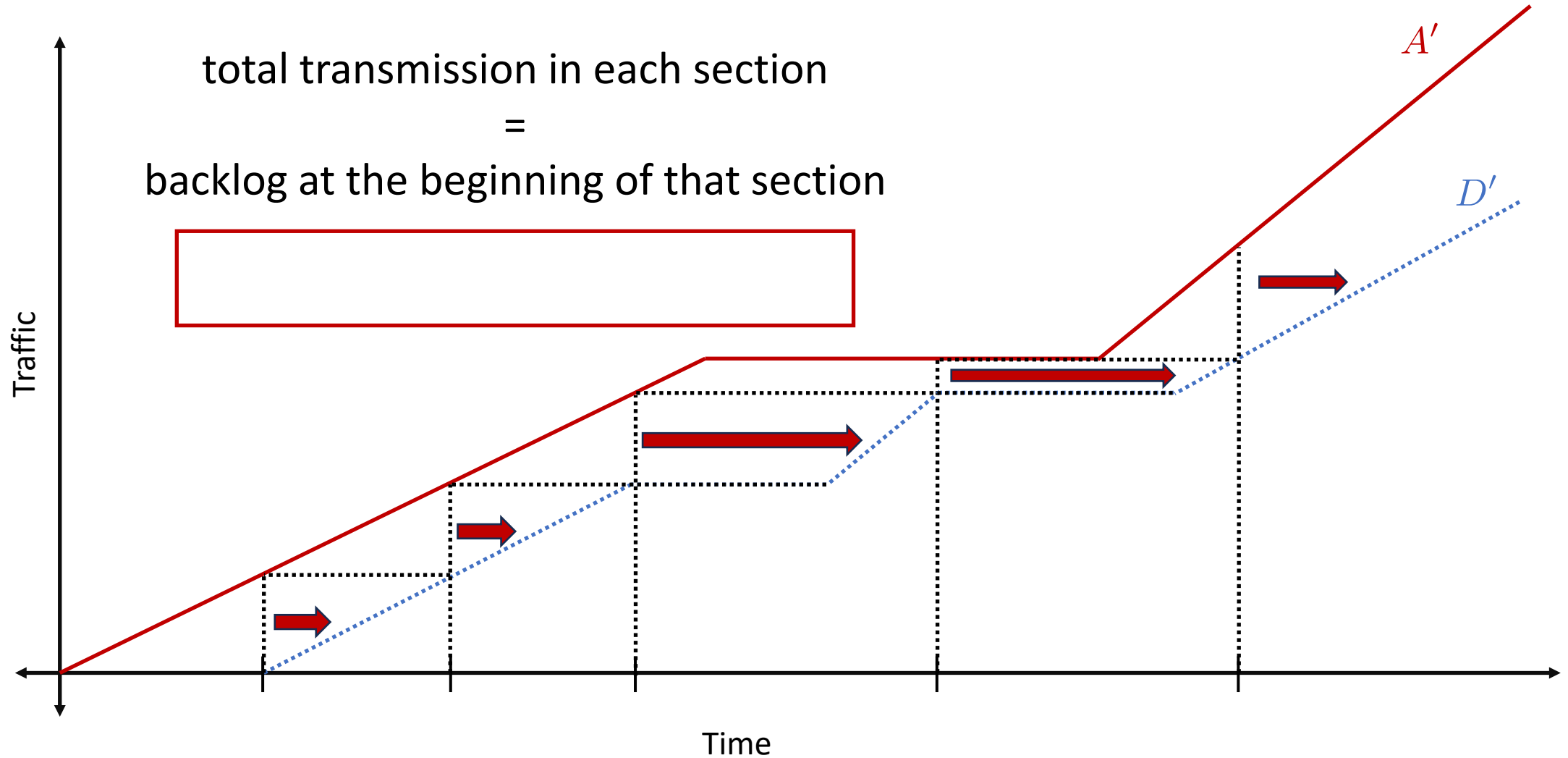$$\left| \frac{D_i(s,t)}{\phi_i} - \frac{D_j(s,t)}{\phi_j} \right| \leq X$$

$$\left| \frac{D_i'(s-d,t-d)}{\phi_i} - \frac{D_j'(s-d,t-d)}{\phi_j} \right| \leq X$$

# Relationship between backlog and transmission over a virtual delay



$$D'(s,t) \leq B(s)$$
$$D'(s,t+\epsilon) = B(s)$$

$A'$

$D'$

virtual delay +

$\epsilon$

Traffic

Time

# Splitting the timeline into sections



total transmission in each section
=
backlog at the beginning of that section

Traffic

Time

$A'$

$D'$

# Pivot definition

- Each flow defines pivots $p_0^i, p_1^i, p_2^i, \ldots$ and $P_i = \{p_k^i\}_{k \in \mathbb{N}_0}$ s.t.

$$A_i(p_0^i) = 0 \qquad p_{k+1}^i = \inf\{\tau \geq p_k^i \mid A'(p_k^i) \leq D(\tau)\} + \epsilon + d$$

- From the definition,

$$\boxed{p_k^i \text{ + virtual delay}}$$

$$\Longrightarrow \quad D_i'(p_{k+1}^i) = A_i'(p_k^i)$$

$$D_i'(p_k^i, p_{k+1}^i) = A_i'(p_k^i) - D_i'(p_k^i) = B_i(p_k^i)$$

$$\forall m \leq n, \ D_i'(p_m^i, p_n^i) = \sum_{k=m}^{n-1} D_i'(p_k^i, p_{k+1}^i) = \sum_{k=m}^{n-1} B_i(p_k^i)$$

# Pivots from different flows interleave each other

- For any flows , and indices , where $p_k^i \leq p_l^j \leq p_{k+1}^i$,

$$D'(p_{k+1}^i) = A'(p_k^i) \leq A'(p_l^j) = D'(p_{l+1}^j) \quad \Longrightarrow \quad p_{k+1}^i \leq p_{l+1}^j$$
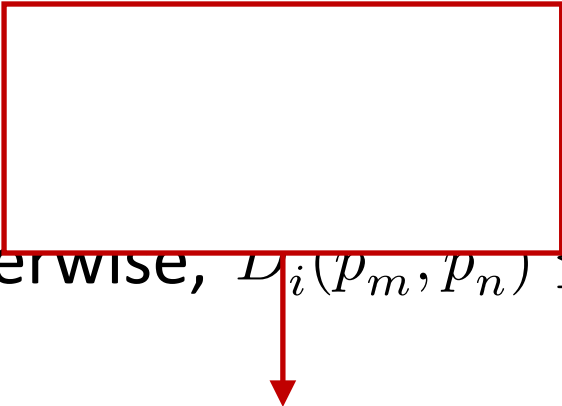
$$\Longrightarrow \quad p_k^i \leq p_l^j \leq p_{k+1}^i \leq p_{l+1}^j \leq p_{k+2}^i \leq p_{l+2}^j \leq p_{k+3}^i \leq \cdots$$

$$P_i(s,t) := \{u \in [s,t] \mid u \in P_i\}$$
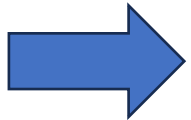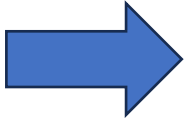
# upper- and lower-bounds over an interval

- Given interval      , let                    and

- If      contains no pivot, it is between          and

- Otherwise,  $D_i(p_m, p_n) \leq D'_i(s,t) \leq D'_i(p^i_{m-1}, p^i_{n+1})$  and so
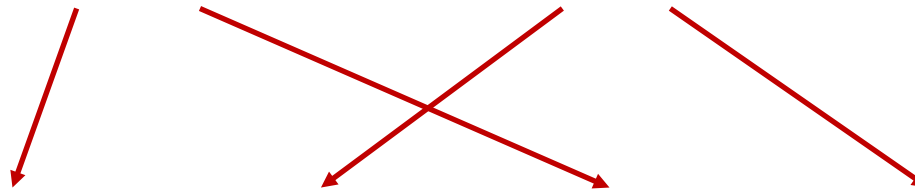
# Bound on              difference over an interval
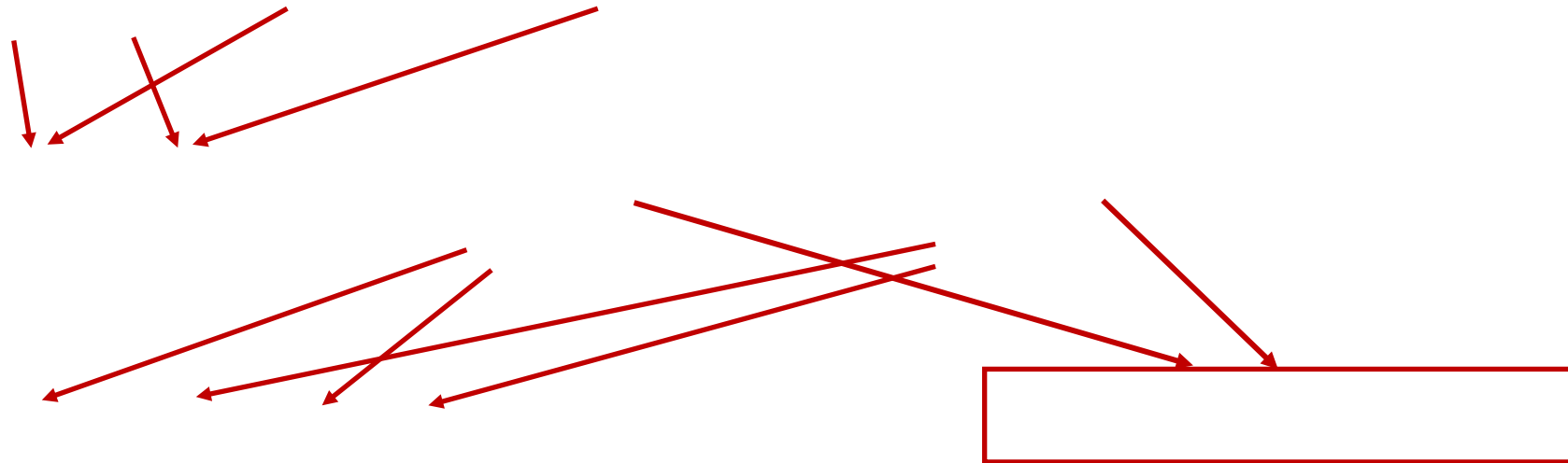
- Let                                and

# Bound on max-min fairness

- If for some constants

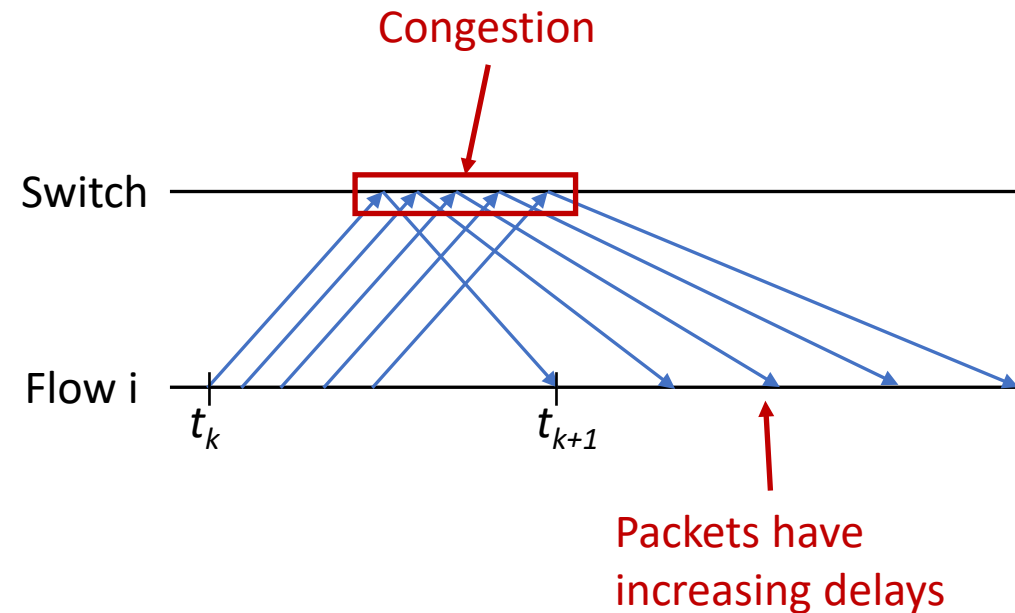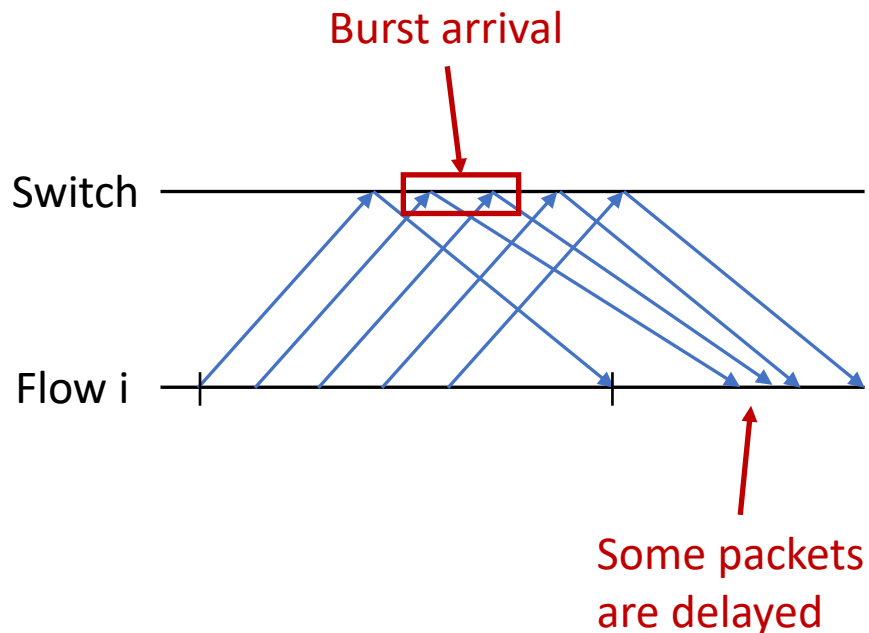   then the network is max-min fair

# Bound on max-min fairness (Con't)

# Using the ACK rate to differentiate short-term and long-term bursts

- We can utilize that fairness bound only depends on     at pivots
  - Flow i transmits data at a fixed rate for
  - Acknowledgement rate can identify the network status



Burst arrival

Switch

Flow i

Some packets are delayed

Congestion

Switch

Flow i

$t_k$        $t_{k+1}$

Packets have increasing delays

# Approximating incast traffic using ACK rate

- For a work-conserving bottleneck switch with rate , If the arrival $A'_i$ is a constant bit rate (CBR), so is the departure $D'_i$ [Yashar'09]

$$A'_i(t) = C_i t \implies D'_i(t) = \begin{cases} C_i t & , \text{if } C \geq \sum_i C_i \\ \frac{C_i}{\sum_i C_i} Ct & , \text{if } C < \sum_i C_i \end{cases}$$

- If $A'_i(t) = C_i t$ and $D'_i(t) = R_i[t - d]^+$ for some R$_i$

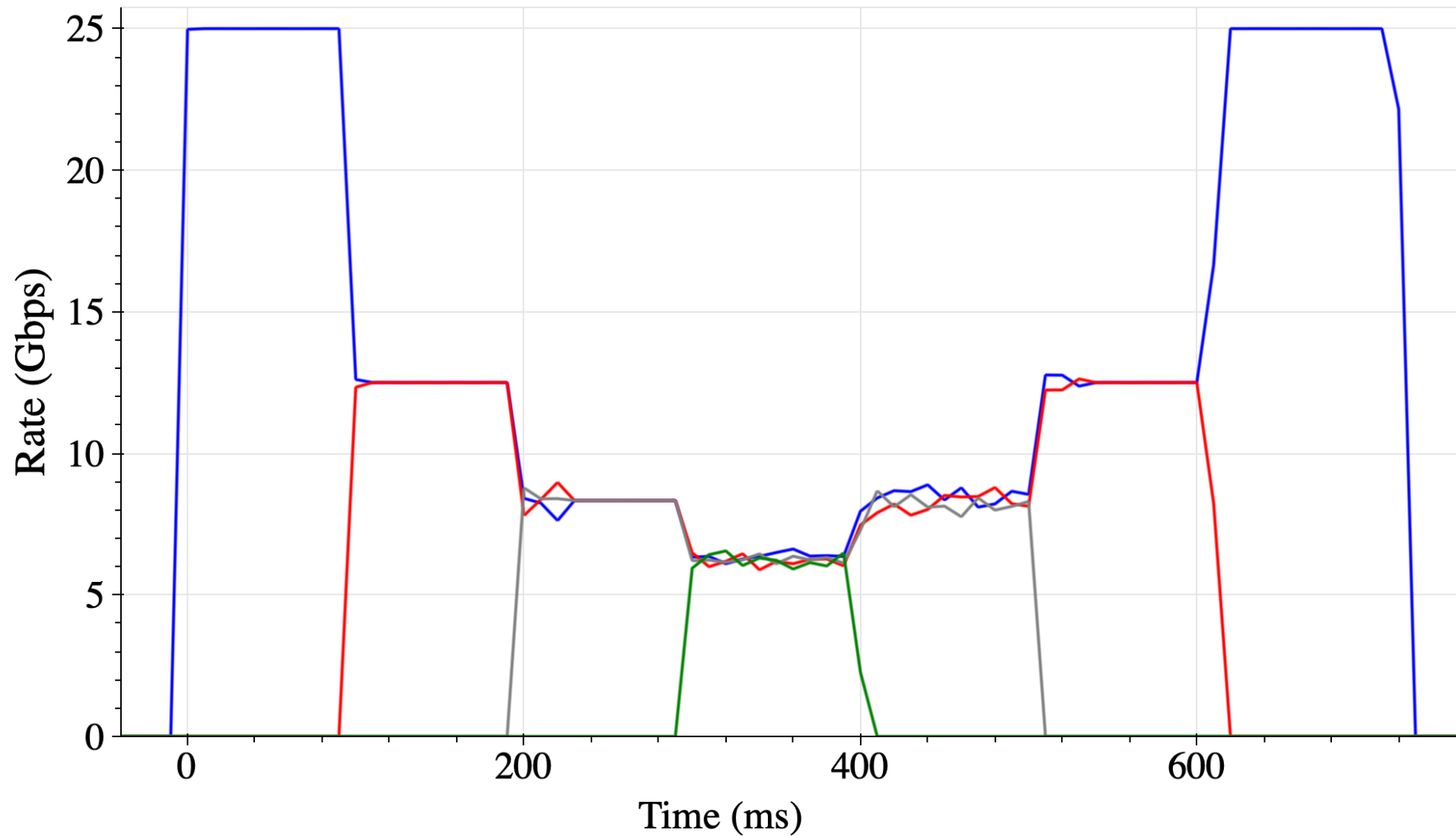$$R_i = C_i \implies C \geq \sum_i C_i$$

$$R_i < C_i \implies R_i = \frac{C_i}{\sum_i C_i} C$$

$$\sum_i C_i = \frac{C_i C}{R_i}$$

19

# Rate-check congestion control

- The sender operates in rounds
- The first packet transmission of each round is a pivot
- The sender transmits at rate $C_r$ in each round r for a duration of d
- Inflight bytes at the beginning of round r is $C_{r-1}d$
- At the beginning of round r, approximate $R_{r-2}$

# Fairness experiment

# Conclusion

- We modeled data center network as multi-flow window flow control

- CCAs "fill the pipe" if the total congestion windows over network-limited flows exceed BDP

- CCAs is max-min fair if the congestion windows <u>at pivots</u> converge toward values proportional to their weights
  - The convergence does not depend on congestion windows outsize of the pivots

# Future work

- Allow the feedback delay for each flow to differ from one another
- Relax the absolutely converge requirements

# Thank you!