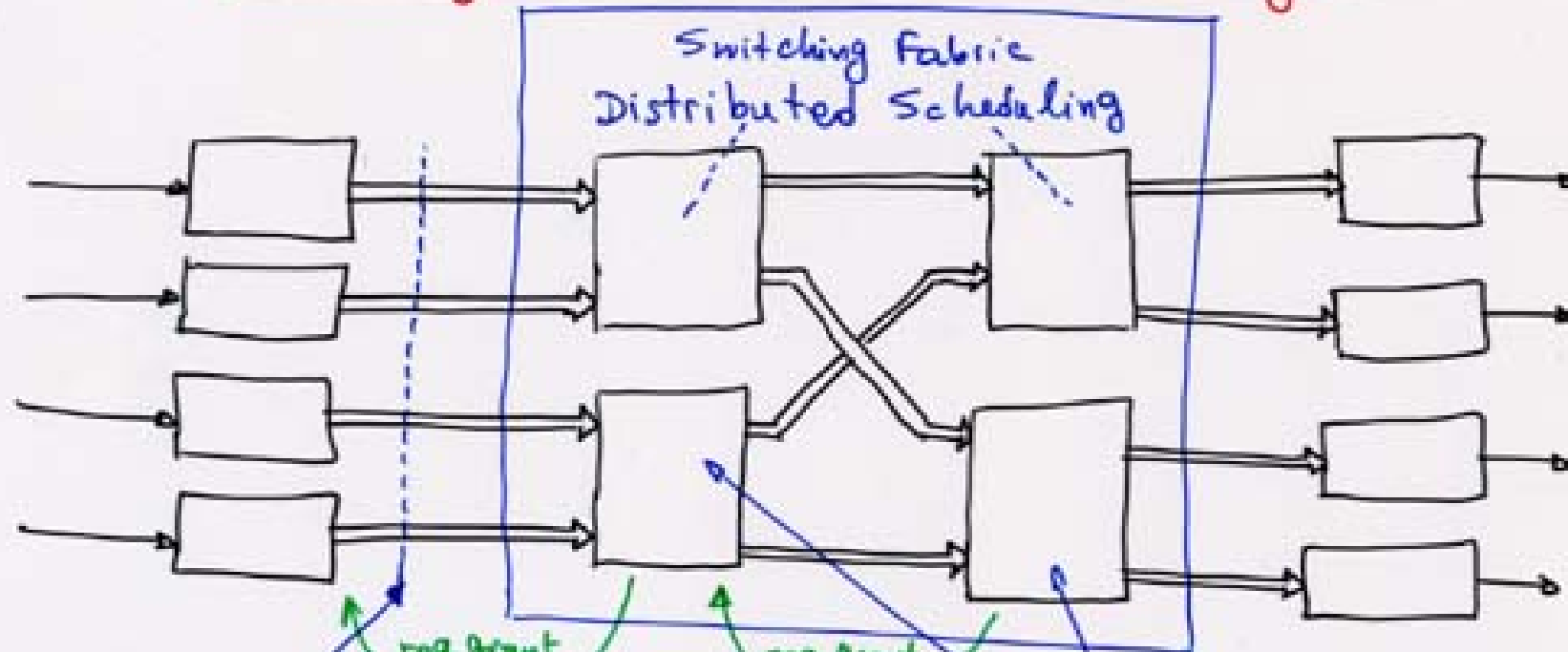


6.1 Credit-Based Flow Control (Backpressure)

- Summary to be added here

Central Scheduler is Impractical for large N

Solution 2: Switching Fabrics with Internal Buffering & Backpressure

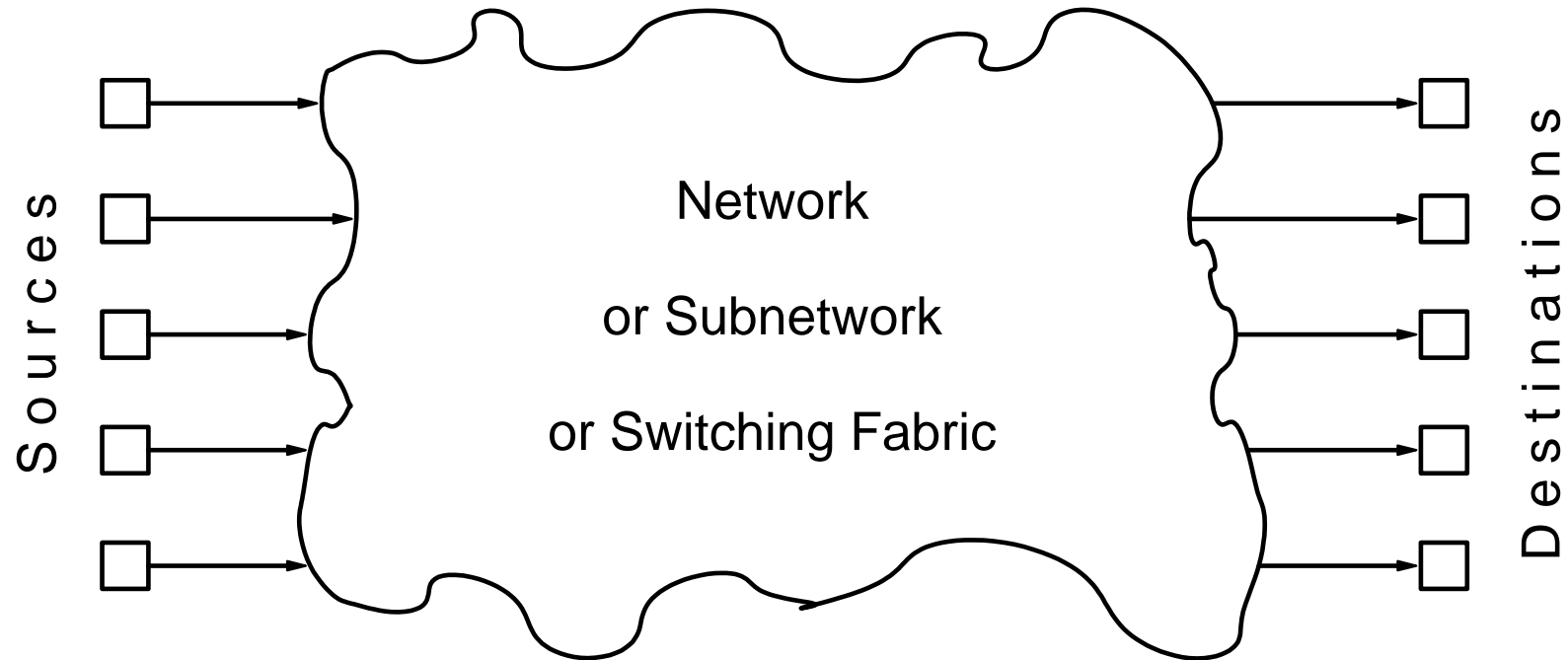


the traffic here may have packets that are short-term-conflicting in the switching fabric, but are long-term-non-conflicting in the fabric

small internal buffers
owing to backpressure and distributed scheduling
handled by these

FLOW CONTROL

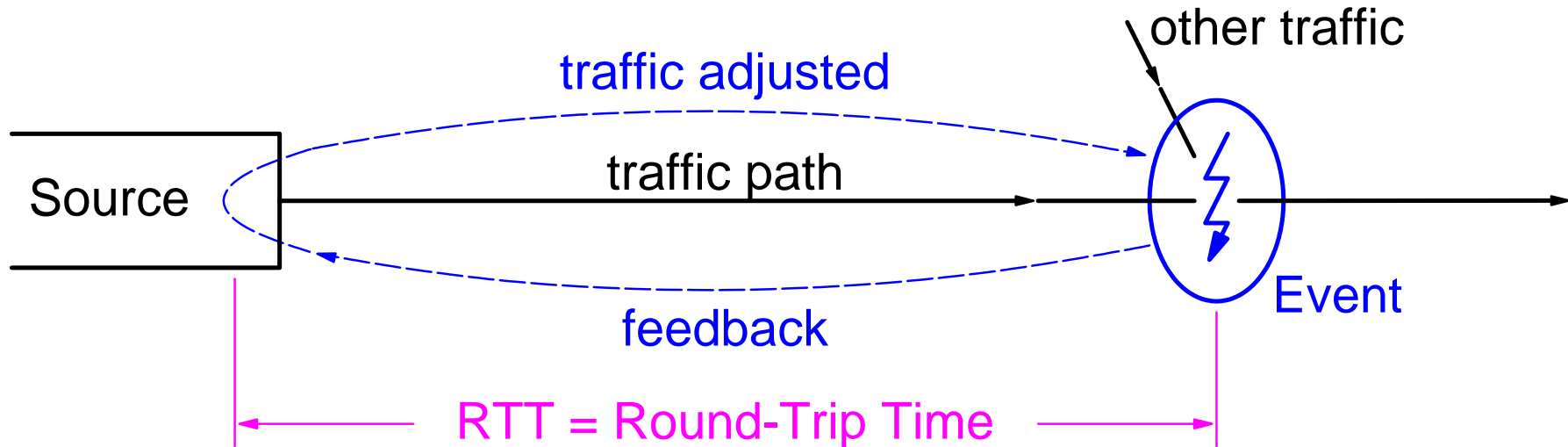
a feedback control problem



- How do the sources know at which rate to transmit?
- How do the sources know when their (collective!) demands exceed the network or the destination capacity?

⇒ Answer: FEEDBACK from the contention point(s), in the network or at the destinations, to the sources

RTT: the Fundamental Time-Constant of Feedback



The traffic is "blind" during a time interval of RTT:

- the source will only learn about the effects of a transmission RTT after this transmission has started (or RTT after a request for such transmission has been issued)
- the (corrective) effects of a contention event will only appear at the contention point RTT after the event occurrence

``Blind Mode" bits in faster Networks

- For faster networks: {
- start transmitting earlier
 - start transmitting at a higher rate



*For networks to get faster,
an increasing number of bits must be sent in ``blind mode"*

initial (``blind mode") rate of transmission	amount of data xmitted in ``blind mode"	
	distance = 8 km RTT ≈ 0.08 ms	distance = 8,000 km RTT ≈ 80 ms
1 Mbit/s	10 Bytes	10 KBytes
1 Gbit/s	10 KBytes	10 Mbytes

Lossy versus Lossless Flow Control

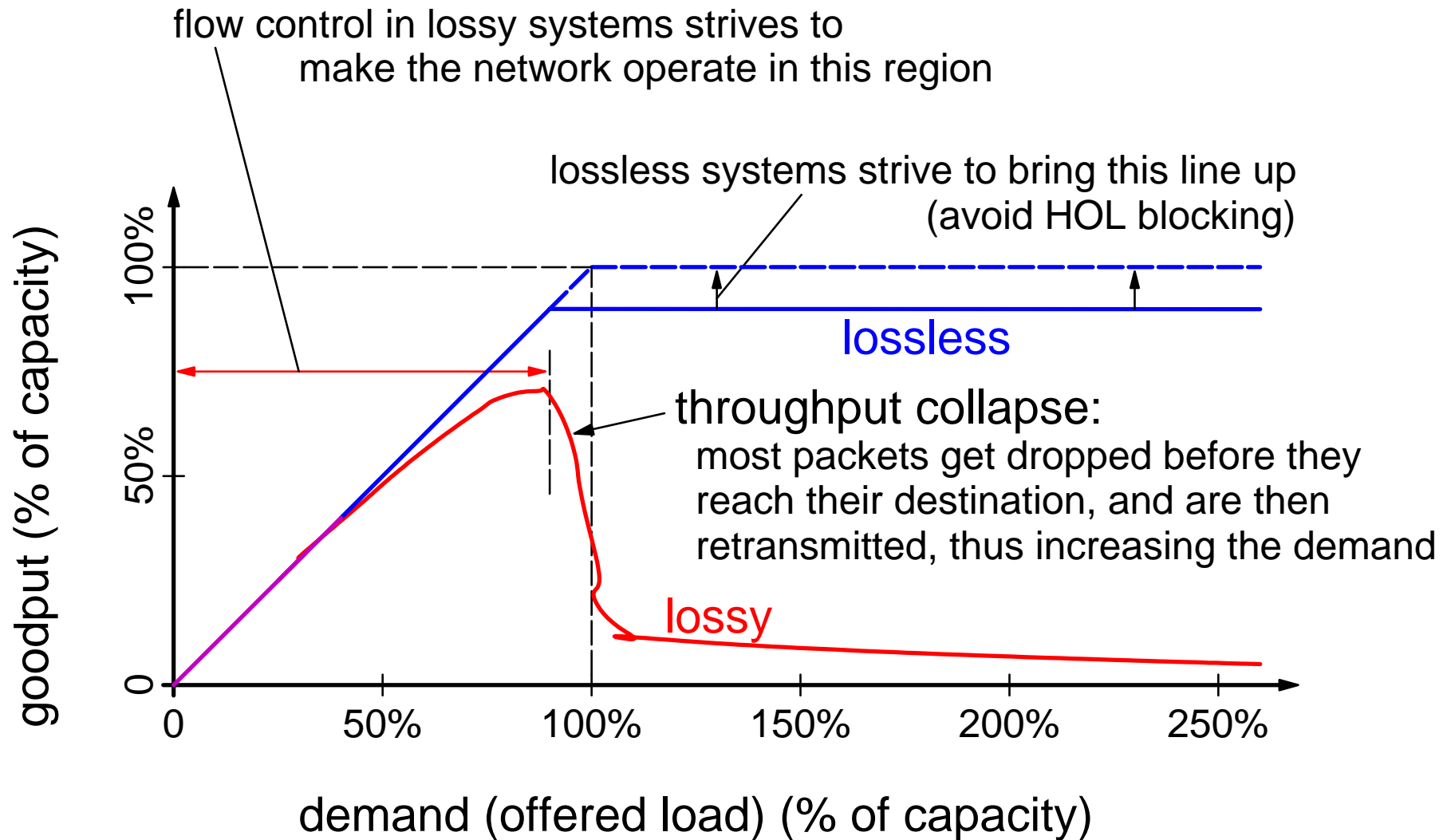
Lossy: *flow control may fail to prevent buffer overflows: packets can be dr*

- inherited from ``communications engineers": same as electrical noise
- simple switches
- for data: need retransmissions => long delays, complex if in H/W
- wastes communication capacity: ``goodput" versus throughput
- need carefully designed protocols to sustain satisfactory goodput

Lossless: *flow control guarantees that buffers will never overflow*

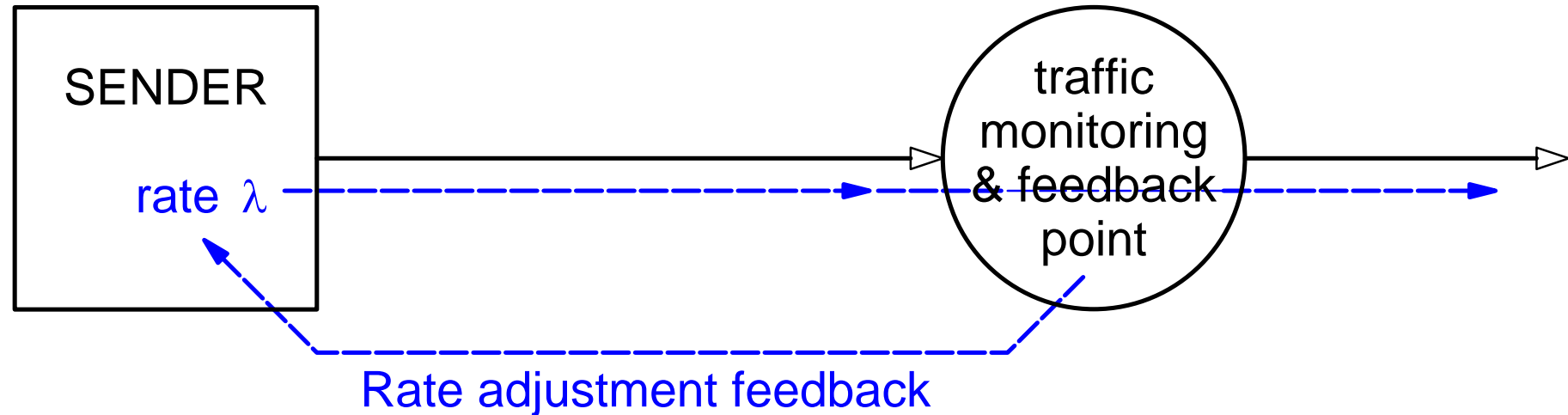
- inherited from ``hardware engineers": processors never drop data
- no wasted communication capacity, minimizes delay
- need multilane protocols to avoid HOL blocking & deadlocks
- switches are more complex, need H/W support for high speed

Goodput versus Demand in Lossy and Lossless Flow Control



(this slide intentionally left blank)

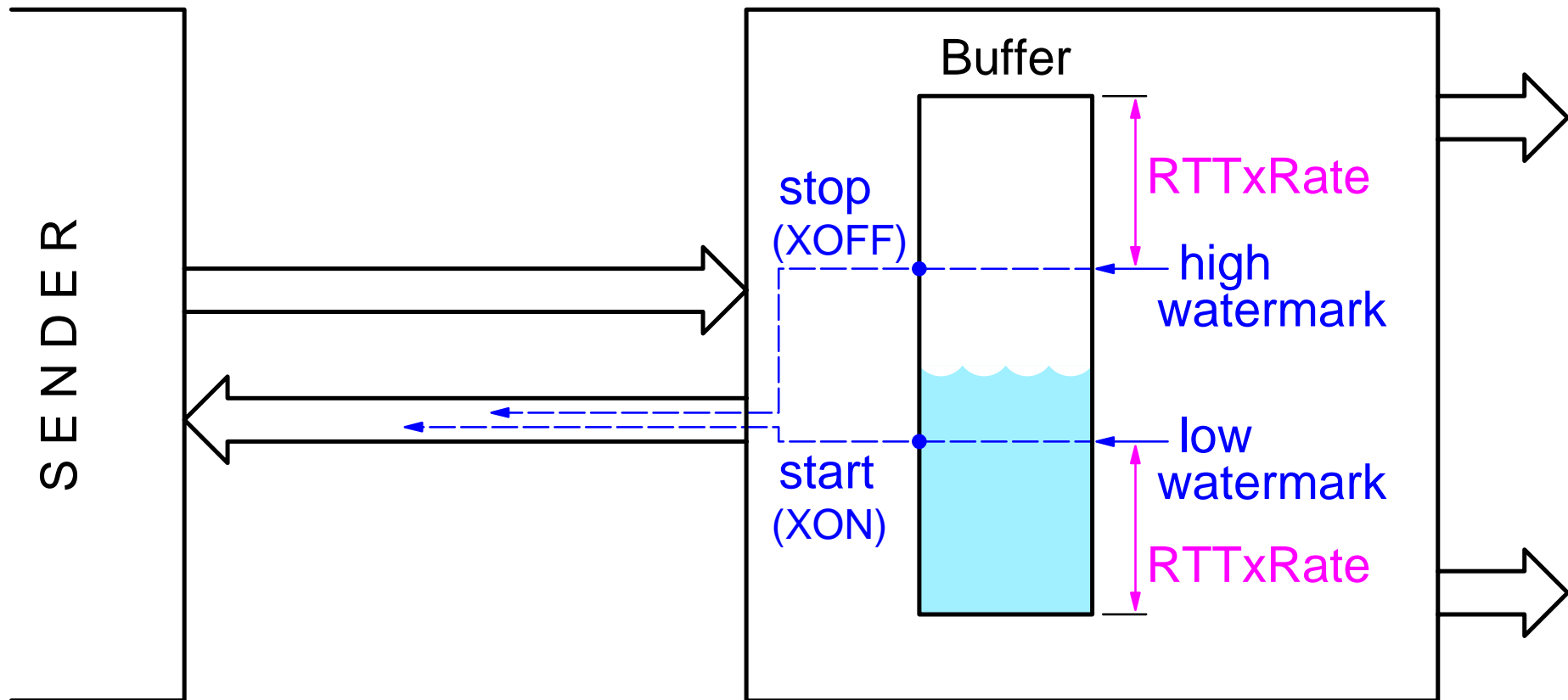
Rate-Based Flow Control



- differential (speed-up / slow-down), or
- absolute (new rate := value)

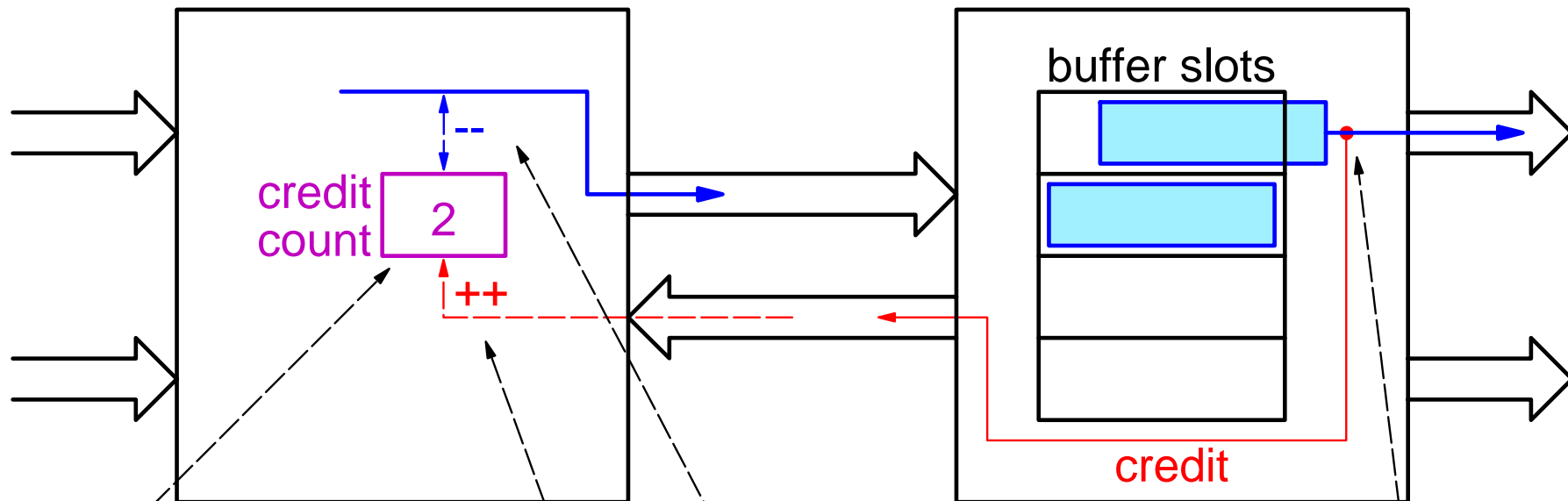
Note: oftentimes, the sender uses a variable-size window mechanism in order to control its rate

ON/OFF (start/stop) (XON/XOFF): simplistic Rate-based FC



- ``start'' \equiv (rate := peak); ``stop'' \equiv (rate := 0)
- rate-based flow control used for lossless transfers
- less than half the buffer efficiency of credit-based flow control

Credit-based (window) (backpressure) Flow Control

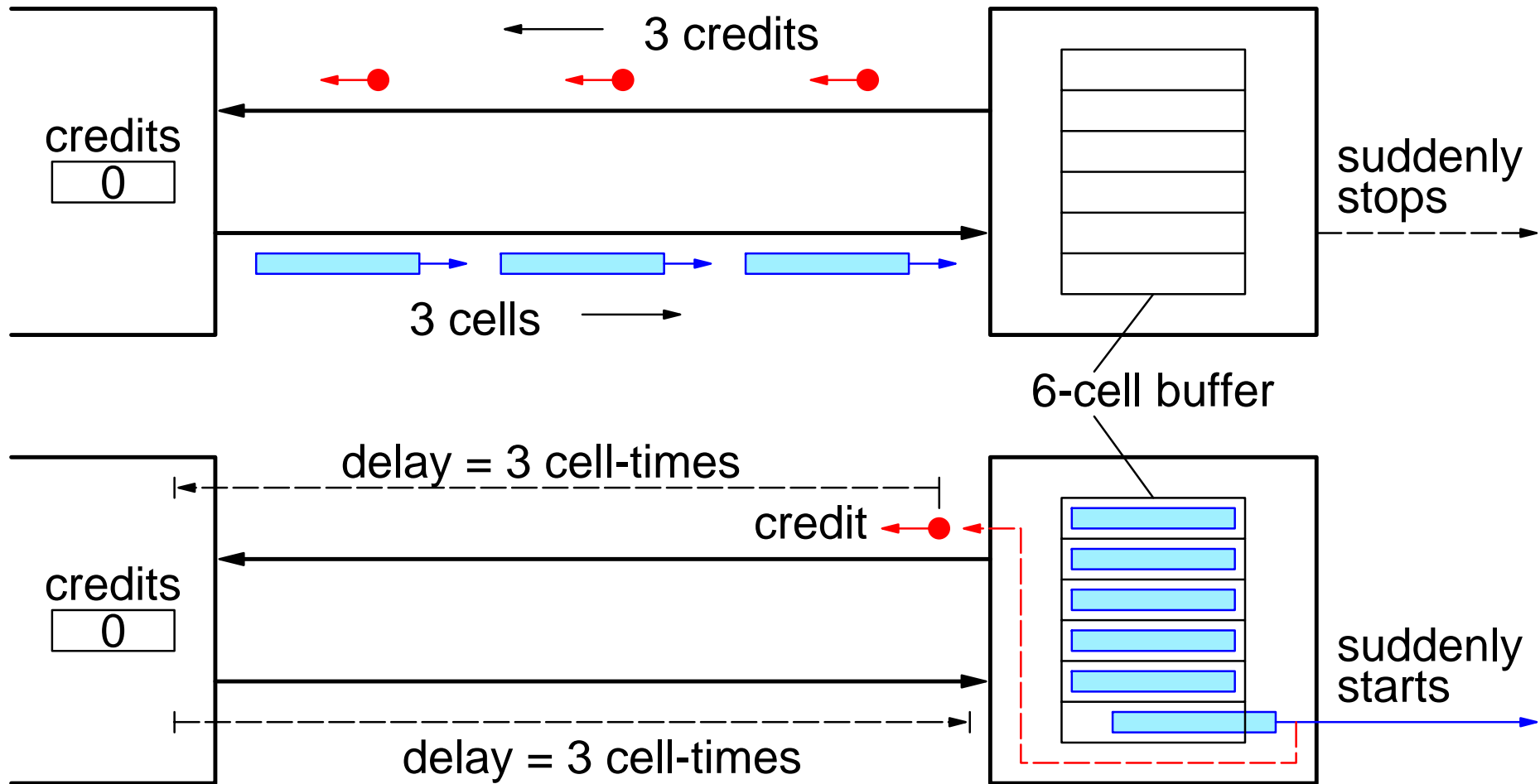


- count of buffer slots known to be available at the downstream site (not allowed to go negative)
- arriving credits increment the credit count
- traffic can only depart if and when it acquires (decrements) the credit(s) that correspond to the buffer slot(s) needed
- when new buffer slots are made available, corresponding credits are sent upstream

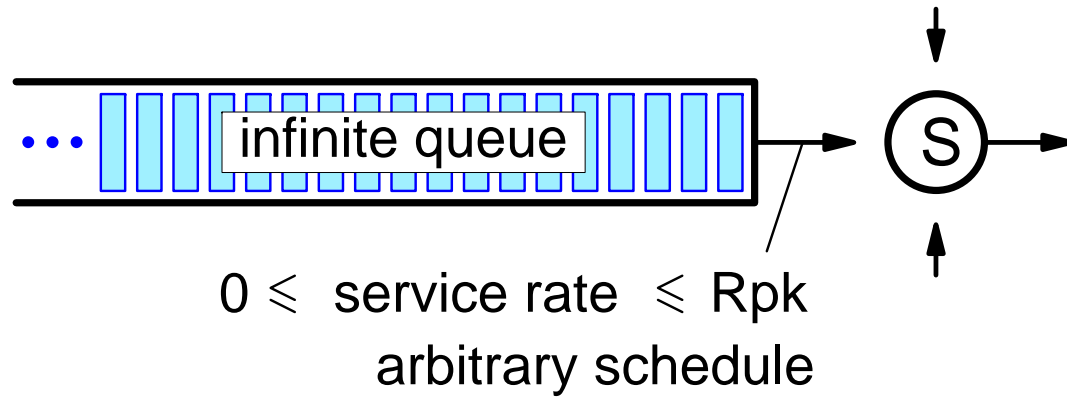
⇒ **Lossless Flow Control**

Buffer Space = Peak Throughput x Round-Trip Time

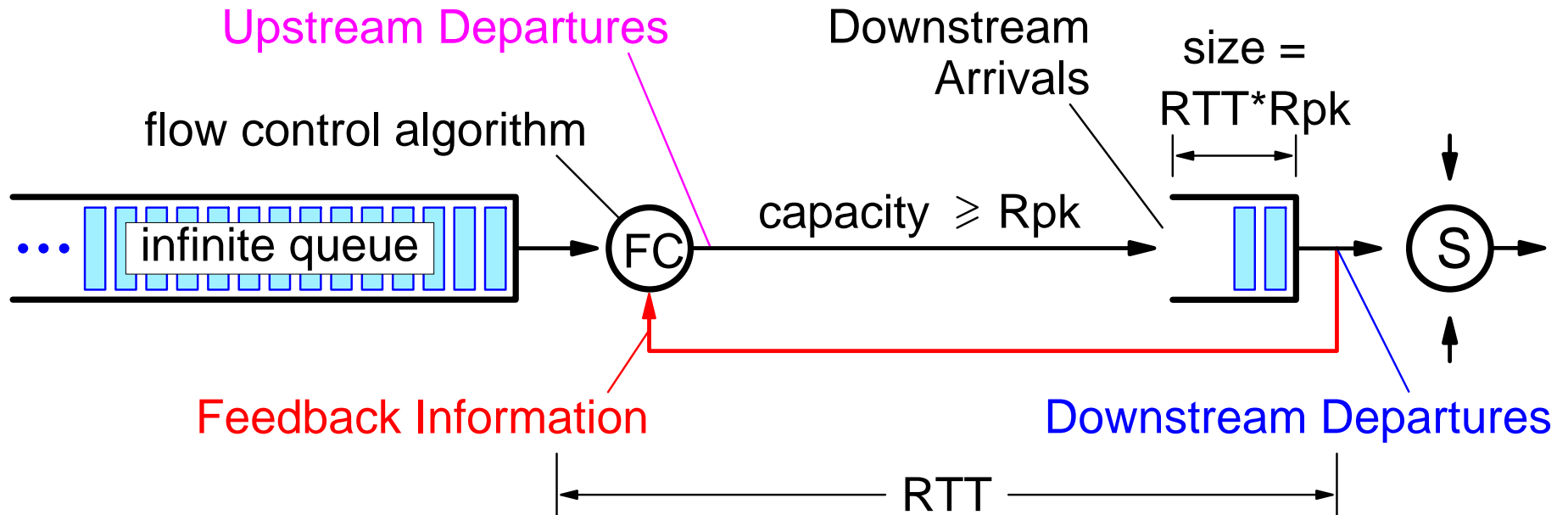
necessary & sufficient

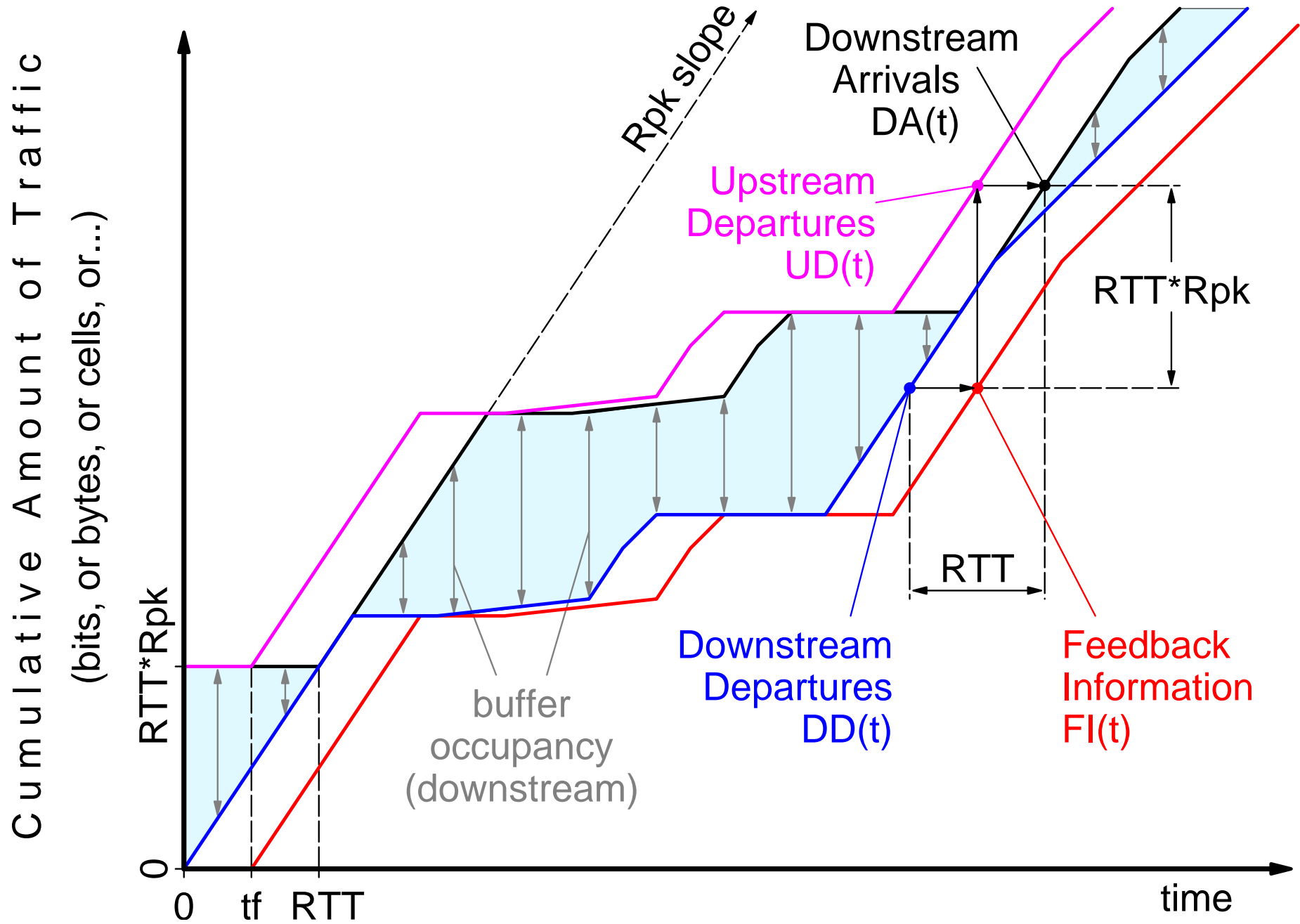


Theorem: Infinite Queue Push-Back



is equivalent to:





- Downstream Departures $DD(t)$ (cumulative):

arbitrary function of time, provided that its slope satisfies:

$$0 \leq \text{service rate of } S \leq R_{pk} \quad (1)$$

$$\Rightarrow 0 \leq DD(t+d) - DD(t) \leq d \cdot R_{pk} \quad (2)$$

- Upstream Departures $UD(t)$ (cumulative):

$UD(t) = DD(t-t_f) + RTT \cdot R_{pk}$; this is always feasible, since:

link capacity $\geq R_{pk} \geq$ service rate of S

- Downstream Arrivals: $DA(t) = DD(t-RTT) + RTT \cdot R_{pk} \quad (3)$

- Buffer Occupancy (downstream): $BO(t) = DA(t) - DD(t)$

$$(3) \Rightarrow BO(t) = RTT \cdot R_{pk} - [DD(t) - DD(t-RTT)]$$

with (2) \Rightarrow $0 \leq BO(t) \leq RTT \cdot R_{pk}$

↑
feasibility of arbitrary
departure schedule
(provided (1) holds)

↑
downstream buffer
never overflows

Feedback Format Options:

how to make the function $DD(t)$ known to the upstream neighbor

- ① QFC Credit-Based Flow Control
- ② Classical (incremental) Credit-Based Flow Control
- ③ Rate-Based Flow Control

① Quantum Flow Control (QFC) <http://www.qfc.org>

Every time $DD(t)$ changes by more than a given threshold N relative to the last time a feedback message was sent, transmit the current value of:

$$DD(t) \text{ modulo } 2^{28} \quad (\text{or modulo } 2^8 \text{ for short links})$$

- credit-based flow control: lossless
- robust: even if a feedback message is corrupted (lost), the next one will restore the error

② Classical (incremental) Credit-Based Flow Control

every N downstream departures ($DD(t_1) = DD(t_0) + N$),
transmit a credit back (N is an implicit parameter);
the upstream node maintains a credit count CC equal to:

$$CC = RTT * R_{pk} + DD(t-t_f) - UD(t);$$

this is incremented by N on every credit arrival, and
decremented by 1 on every departure of a unit of traffic

- shorter feedback (credit) messages than QFC
- non-robust: loss of a credit leads to buffer and transmission capacity underutilization;
accumulated losses of credits lead to deadlock!

Equivalence of Rate and Credit Based Flow Control

③ Rate-Based Flow Control

On every change of the slope of $DD(t)$ (rate of downstream departures), send back the new value of the rate (slope); upon reception of such feedback, the upstream node adjusts its rate of transmission (slope of $UD(t)$) to the value received; thus, $UD(t)$ is (almost!) a delayed and shifted-up copy of $DD(t)$.

- Rate-based flow control
- Could be made lossless, but this would not be robust: slight mismatches between real & measured rate accumulate to large differences between the values of $DD(t-t_f)$, $UD(t)$; similarly, variations in t_f (delay of feedback messages) lead to $UD(t)$ value errors.