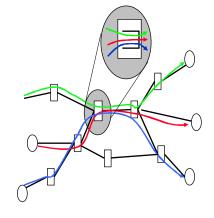
# CS 268: Lecture 8 Router Support for Congestion Control

Ion Stoica
Computer Science Division
Department of Electrical Engineering and Computer Sciences
University of California, Berkeley
Berkeley, CA 94720-1776

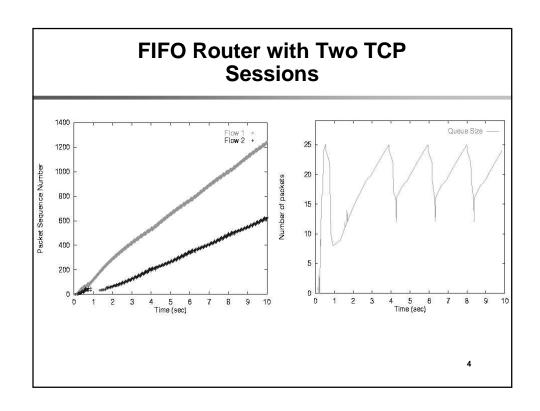
# Router Support For Congestion Management

- Traditional Internet
  - Congestion control mechanisms at end-systems, mainly implemented in TCP
  - Routers play little role
- Router mechanisms affecting congestion management
  - Scheduling
  - Buffer management
- Traditional routers
  - FIFO
  - Tail drop



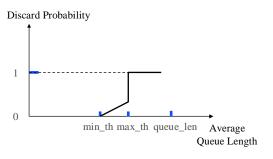
# **Drawbacks of FIFO with Tail-drop**

- Buffer lock out by misbehaving flows
- Synchronizing effect for multiple TCP flows
- Burst or multiple consecutive packet drops
  - Bad for TCP fast recovery



#### **RED**

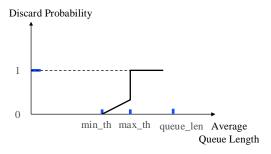
- FIFO scheduling
- Buffer management:
  - Probabilistically discard packets
  - Probability is computed as a function of average queue length (why average?)



5

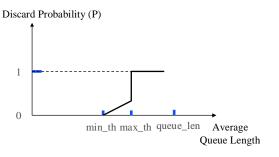
# RED (cont'd)

- min\_th minimum threshold
- max\_th maximum threshold
- avg\_len average queue length
  - avg\_len = (1-w)\*avg\_len + w\*sample\_len



# RED (cont'd)

- If (avg\_len < min\_th) → enqueue packet</p>
- If (avg\_len > max\_th) → drop packet
- If (avg\_len >= min\_th and avg\_len < max\_th) → enqueue packet with probability P



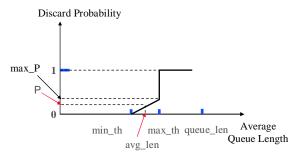
-

# RED (cont'd)

- P = max\_P\*(avg\_len min\_th)/(max\_th min\_th)
- Improvements to spread the drops

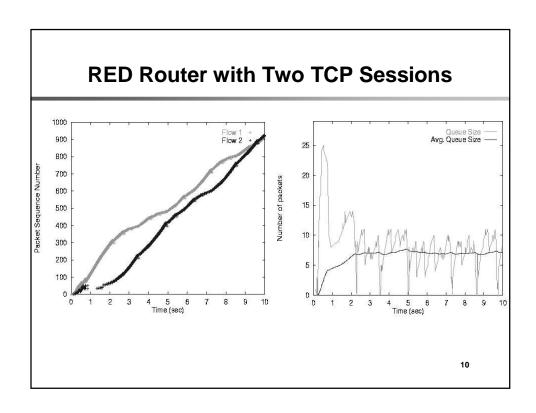
 $P' = P/(1 - count^*P)$ , where

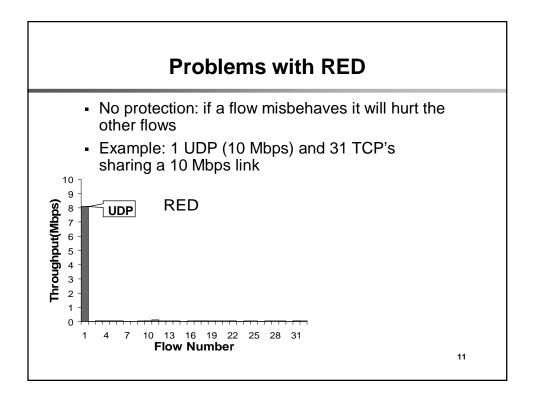
• count – how many packets were consecutively enqueued since last drop

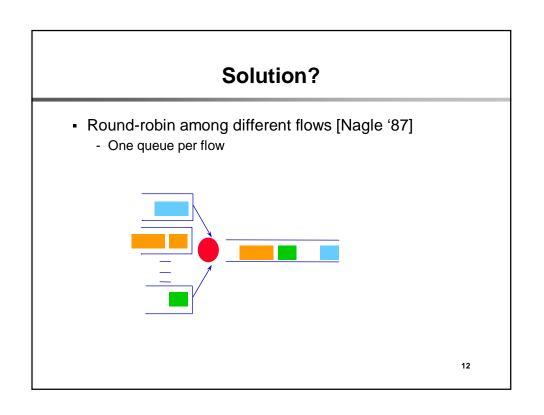


# **RED Advantages**

- Absorb burst better
- Avoids synchronization
- Signal end systems earlier







#### **Round-Robin Discussion**

- Advantages: protection among flows
  - Misbehaving flows will not affect the performance of well-behaving flows
  - FIFO does not have such a property
- Disadvantages:
  - More complex than FIFO: per flow queue/state
  - Biased toward large packets a flow receives service proportional to the number of packets (When is this bad?)

13

## Solution?

- Bit-by-bit round robin
- Can you do this in practice?
- No, packets cannot be preempted (why?)
- ...we can only approximate it

# Fair Queueing (FQ) [DKS'89]

- Define a fluid flow system: a system in which flows are served bit-by-bit
- Then serve packets in the increasing order of their deadlines
- Advantages
  - Each flow will receive exactly its fair rate
- Note:
  - FQ achieves max-min fairness

15

### **Max-Min Fairness**

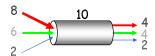
- Denote
  - C link capacity
  - N number of flows
  - $r_i$  arrival rate
- Max-min fair rate computation:
  - 1. compute C/N
  - 2. if there are flows *i* such that  $r_i \le C/N$ , update C and N

$$C = C - \sum_{i \, s.t \, r_i \le C} r_i$$

- 3. if no, f = C/N; terminate
- 4. go to 1
- A flow can receive at most the fair rate, i.e.,  $min(f, r_i)$

# **Example**

- C = 10;  $r_1 = 8$ ,  $r_2 = 6$ ,  $r_3 = 2$ ; N = 3
- $C/3 = 3.33 \rightarrow C = C r3 = 8$ ; N = 2
- C/2 = 4; f = 4



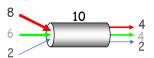
f = 4: min(8, 4) = 4 min(6, 4) = 4 min(2, 4) = 2

17

# **Alternate Way to Compute Fair Rate**

• If link congested, compute *f* such that

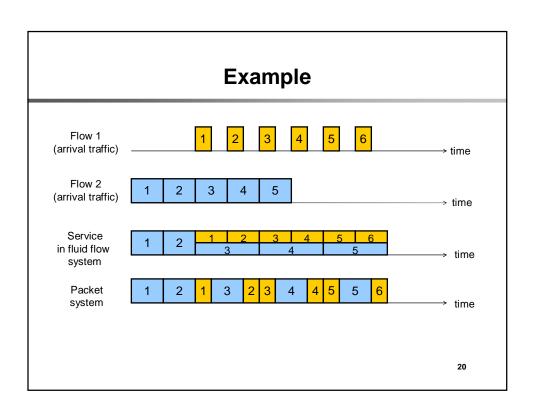
$$\sum_{i} \min(r_i, f) = C$$



f = 4: min(8, 4) = 4 min(6, 4) = 4 min(2, 4) = 2

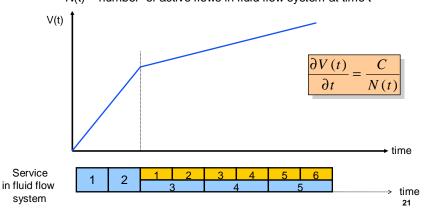
# **Implementing Fair Queueing**

 Idea: serve packets in the order in which they would have finished transmission in the fluid flow system



# System Virtual Time: V(t)

- Measure service, instead of time
- V(t) slope rate at which every active flow receives service
  - C link capacity
  - N(t) number of active flows in fluid flow system at time t



# **Fair Queueing Implementation**

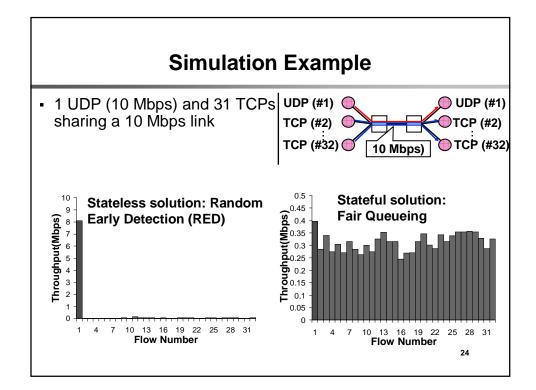
- Define
  - $F_i^k$  finishing time of packet k of flow i (in system virtual time reference system)
  - $a_i^k$  arrival time of packet k of flow i  $L_i^k$  length of packet k of flow i
- The finishing time of packet k+1 of flow i is

$$F_i^{k+1} = \max(V(a_i^k), F_i^k) + L_i^{k+1}$$

# "Weighted Fair Queueing" (WFQ)

- What if we don't want exact fairness?
  - E.g.,: file servers
- Assign weight w<sub>i</sub> to each flow i
- And change virtual finishing time

$$F_i^{k+1} = \max(V(a_i^k), F_i^k) + \frac{L_i^{k+1}}{w_i}$$



#### **Core-Stateless Fair Queueing (CSFQ)**

- Fair Queueing requires per flow state in routers
  - Maybe impractical for very high speed routers
- Core Stateless Fair Queueing eliminates the state at core routers ...
- ... but only approximates FQ's behavior

25

# Insight

 If each packet of a flow with arrival rate r is forwarded with probability

$$P = \min\left(1, \frac{f}{r}\right)$$

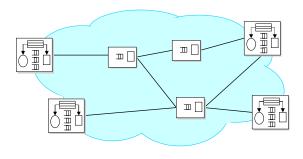
• the rate of flow's forwarded traffic r' is

$$r' = r \times P = r \times \min\left(1, \frac{f}{r}\right) = \min(r, f)$$

- No need to maintain per-flow state if r is carried in the packet
  - Need to update rate in packet to r

# **CSFQ**

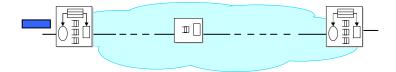
- A contiguous and trusted region of network in which
  - Edge nodes perform per flow operations
  - Core nodes do not perform any per flow operations



27

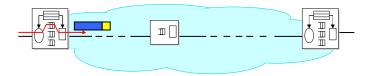
# **Algorithm Outline**

Ingress nodes: estimate rate r for each flow and insert it in the packets' headers



# **Algorithm Outline**

 Ingress nodes: estimate rate r for each flow and insert it in the packets' headers



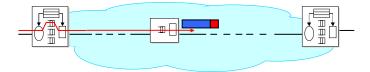
29

# **Algorithm Outline**

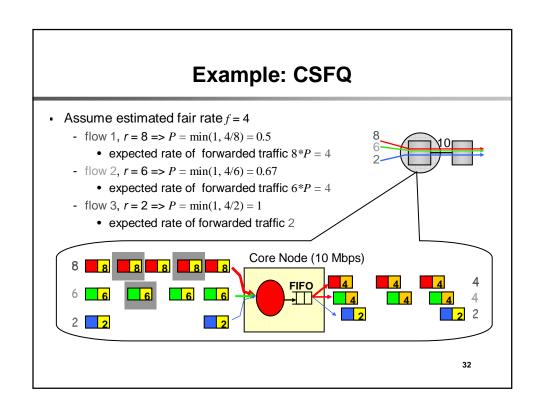
- Core node:
  - Compute fair  $\operatorname{rate} f$  on the output link
  - Enqueue packet with probability

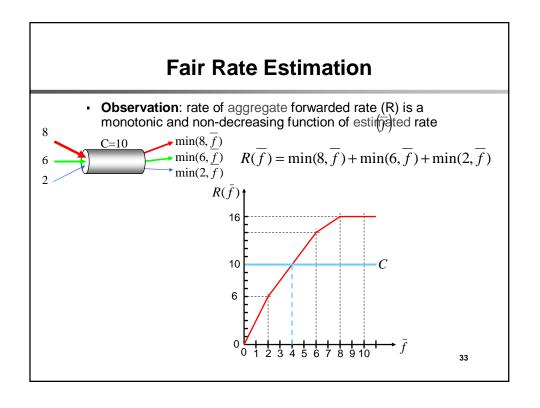
$$P = \min(1, f/r)$$

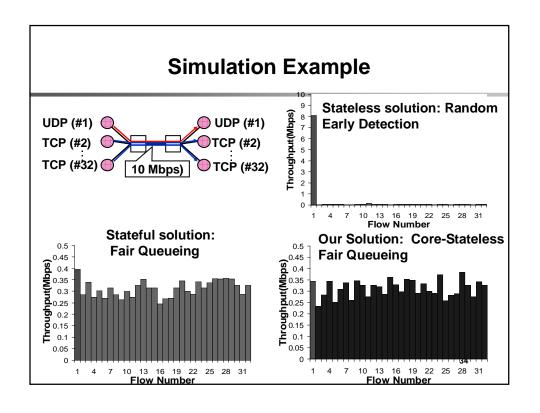
- Update packet label to  $r = \min(r, f)$ 



# Algorithm Outline • Egress node: remove state from packet's header







#### Summary

- FQ does not eliminate congestion → it just manages the congestion
- You need both end-host congestion control and router support for congestion control
  - End-host congestion control to adapt
  - Router congestion control to protect/isolate
- Don't forget buffer management: you still need to drop in case of congestion. Which packet's would you drop in FQ?
  - One possibility: packet from the longest queue

35

#### **Announcements**

- Project feedback
  - Tuesday, Feb 14, 12:30-2pm
  - Wednesday, Feb 15, 11:30-1pm