



# ADVANCED COMPUTER NETWORKS

[DKS89] Demers, Keshav, and Shenker, "Analysis and Simulation of a Fair Queueing Algorithm," *Proc. of ACM SIGCOMM '89*, 19(4):1-12, Sep. 1989

## Queueing and Scheduling

When would you see a queue?

- when resource is limited
- and there is a contention for the resource
- in networking context: when packet incoming rate is faster than outgoing (service) rate

Want to be fair in service and to protect against misbehaving connections

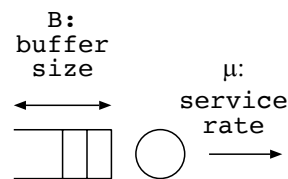
Scheduling discipline: which packet to serve next?

## Queue Management

Queue management design issues:

- Fairness
  - at the minimum want protection against malicious sources
- Delay bound
- Drop policy
- Cost of operation

FIFO/FCFS Scheduler:



A resource centric view: max-min fair share, assuming users have **equal rights** to resource

## Max-Min Fair Share

Let:

$\mu_{total}$ : total resource (e.g., bandwidth) available

$\mu_i$ : total resource given to (flow)  $i$

$\mu_{fair}$ : fair share of resource

$\rho_i$ : request for resource by (flow)  $i$

Max-Min fair share is  $\mu_i = \text{MIN}(\rho_i, \mu_{fair})$

$$\mu_{total} = \sum_{i=1 \text{ to } n} \mu_i$$

# Max-Min Fair Share

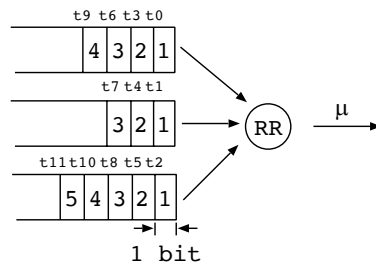
In words: max-min fair share maximizes minimum shares of flows whose demands have not been fully satisfied

1. no flow gets more than its request
2. no other allocation satisfying condition 1 has a higher minimum allocation
3. condition 2 remains recursively true as we remove the minimal request and reduce total resource accordingly

# Bit-by-Bit Round Robin

1 round,  $R()$ , is defined as all non-empty queues have been served 1 quantum

- $R(t_5) = 2$
- time at Round 3? Round 4?



BbB-RR achieves max-min fair share

Max-min fair-share isolates flows

BbB-RR protects against misbehaving flows

A.k.a. Generalized Processor Sharing (GPS)

# Max-Min Fair Share Example

Let:

$i$	$\rho_i$	$\mu_i$
A	12	11
B	11	11
C	8	8

$\mu_{total} = 30$

Initially  $\mu_{fair} = 10$

$\rho_C = 8$ , so unused resource ( $10 - 8 = 2$ ) is divided evenly between the remaining  $i$ 's whose demands are not yet fully met

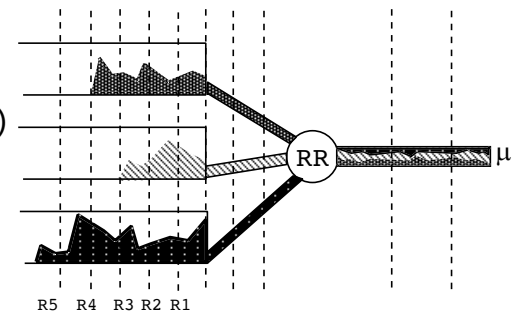
Thus,  $\mu_{fair} = 10 + 2/2 = 11$

# Fluid-Flow Approximation

A continuous service model

- instead of thinking of each quantum as serving discrete bits in a given order
- think of each connection as a stream of fluid, described by the **speed** and **volume** of flow

At each quantum the same amount of fluid from each (non-empty) stream flows out concurrently

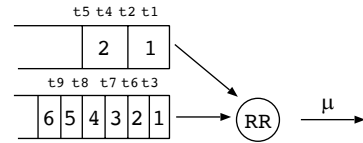


# Packetized Scheduling

Packet-by-packet round robin:

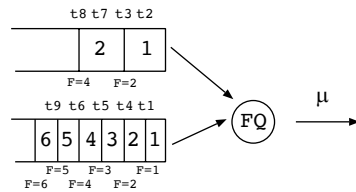
Problem:

Gives bigger share to flows with big packets



Packet-by-packet fair-queueing:

- $F$ : finish round, the round a packet finishes service
- simulates RR in the computation of  $F$ 's
- serve packets with the smallest  $F$  first



# Round Computation

Recall: 1 round is defined as all active flows have been served 1 quantum

Round's rate of change:  $\partial R / \partial t = \mu / N_{ac}(t)$ , where

- $\mu$ : link bandwidth
- $N_{ac}(t)$ : number of active flows at time  $t$

Flow  $\alpha$  is active at time  $t$  if  $R(t) \leq F_i^\alpha$ , where  $i$  is the last bit of flow  $\alpha$  in queue

The speed of 1 round through all active flows depends on the number of active flows: faster if there are less active flows

# Start and Finish Rounds

At what round does packet  $i$  of flow  $\alpha$  start seeing service?

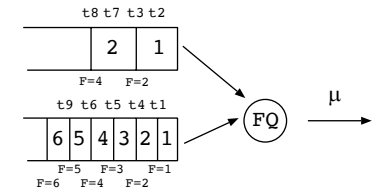
$$S_i^\alpha = \text{MAX}(F_{i-1}^\alpha, A_i^\alpha)$$

- $A_i^\alpha = R(t_i^\alpha)$ : round at the time packet  $i$  arrives
- $S_i^\alpha = F_{i-1}^\alpha$  if there is a queue,  $A_i^\alpha$  otherwise

When does packet  $i$  finish service?

$$F_i^\alpha = S_i^\alpha + P_{ii}^\alpha$$

where  $P_{ii}^\alpha$  is the size (service time) of packet  $i$



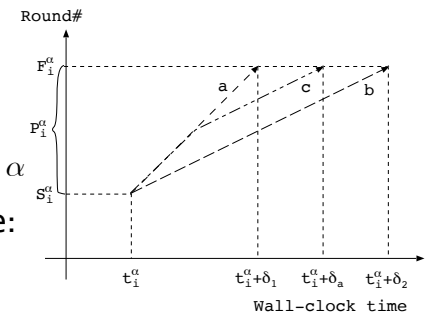
# Round# vs. Wall-Clock Time

Let:

- **time**: wall-clock time
- **round**: virtual-clock time
- $\mu = 1$  unit
- $t_i^\alpha$ : arrival time of packet  $i$  of flow  $\alpha$

Computing the rate of change:

- a**:  $N_{ac} = 1, \partial R / \partial t = 1$ ,
- b**:  $N_{ac} = 2, \partial R / \partial t = 1/2, \delta_2 = 2 * \delta_1$
- c**: at the beginning,  $N_{ac} = 1, \partial R / \partial t = 1$ , halfway serving packet  $i$ , a packet belonging to another flow arrives,  $N_{ac} = 2, \partial R / \partial t = 1/2$



As  $N_{ac}(t)$  changes, finish round stays the same, actual time stretches

# Round Computation Example

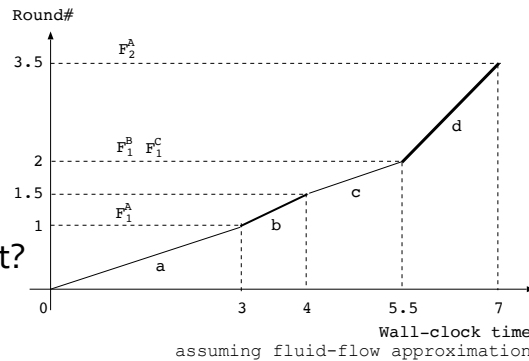
Scenario:

- flows A has 1 packet of size 1 arriving at time 0
- flows B and C each has 1 packet of size 2 arriving at time 0
- flow A has another packet of size 2 arriving at time 4

Slope ( $\partial R / \partial t$ ):

$a = 1/3, b = 1/2,$   
 $c = 1/3, d = 1$

What is the arrival round of A's 2<sup>nd</sup> packet?  
 $R(t^A_2) = 1.5$



# Iterated Deletion

Variables:

- $t_i$ : the time when the round was last computed
- $R(t_i)$ : the round at that time
- global array of finish rounds of inactive flows, sorted

Want to compute the current Round by taking into account each of the finish rounds in the array

Finish round rate at time  $t$ :  $\partial F^\alpha_i / \partial t = \mu / N_{ac}(t)$

We know  $N_{ac}(t)$  hasn't changed since  $t_i$ :

- Finish round, given finish time:  $F^\alpha_i = R(t_i) + \mu (t(F^\alpha_i) - t_i)$
- Finish time, given finish round:  $t(F^\alpha_i) = t_i + [F^\alpha_i - R(t_i)]N_{ac} / \mu$

Now walk the global array and compute "round catchup" to the present time

# Arrival Round Computation

When packet  $i$  of an active flow arrives, its finish time is computed as  $F^\alpha_i = F^\alpha_{i-1} + P^\alpha_i$ , where  $F^\alpha_{i-1}$  is the finish time of the last packet in  $\alpha$ 's queue

If flow  $\alpha$  is inactive, there's no packet in its queue,  $F^\alpha_i = A^\alpha_i + P^\alpha_i$ , how do we compute  $A^\alpha_i$ ?

If flow  $\alpha$  has been inactive for  $\Delta t$  time and there has been  $N_{ac}$  flows during the whole time, we can compute  $A^\alpha_i = F^\alpha_{i-1} + \Delta t(1/N_{ac}) \Rightarrow$  inactive flow must store  $F^\alpha_{i-1}$

But what if  $N_{ac}$  has changed, several times, over  $\Delta t$ ?

# Iterated Deletion: the Algorithm

Let  $t^C_i$  be the arrival time of packet  $i$  of flow C

For each flow  $\alpha$  with packets enqueued or has last finish round in the global array {

```

compute  $t(F^\alpha_i)$  // finish time of flow  $\alpha$ 
if ( $t(F^\alpha_i) \leq t^C_i$ ) { // if flow is inactive by time  $t^C_i$ 
     $R(t(F^\alpha_i)) = F^\alpha_i$  // finish round when flow went inactive
     $t_l = t(F^\alpha_i)$  // the last round computation time
     $N_{ac}--$ 
} else { break }

```

Time complexity:  $O(N_{ac})$

}  
 $R(t^C_i) = R(t_l) + \mu / N_{ac}(t^C_i - t_l), // A^\alpha_i$

assuming  $N_{ac} \neq 0$ , else may reset  $R$

$t_l = t^C_i$

## Variation: Credit Accumulation

Allow a flow to have a bigger share if it has been idle

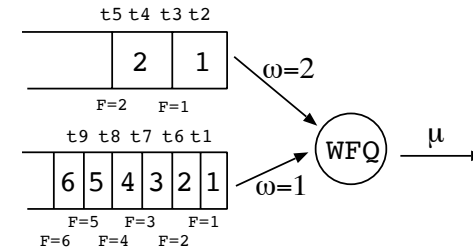
Compute bid per packet:

$$B_i^\alpha = \text{MAX}(B_{i-1}^\alpha, A_i^\alpha - \delta) + P_i^\alpha$$

- if  $\delta = 0$ , no credit accumulated,  $B_i^\alpha = F_i^\alpha$
- if  $\delta = \infty$ ,  $B_i^\alpha = F_{i-1}^\alpha + P_{ii}^\alpha$  regardless of packet  $i$ 's arrival time

Credit accumulation is discouraged because it can be abused: accumulate credits for a long time, then send a big burst of data

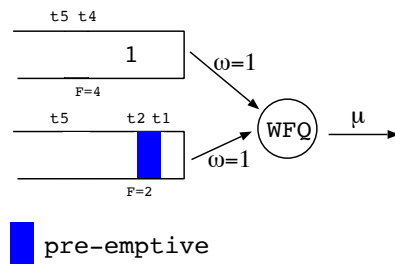
## Weighted Fair Queueing



$$F_{i-1}^\alpha = S_i^\alpha + P_i^\alpha / \omega^\alpha$$

where  $\omega^\alpha$  is the weight (reserved rate) of flow  $\alpha$

## Pre-emptive WFQ



Non-preemptive worst-case: packet arrives as another packet of larger  $F$  just started service

Packet has to wait  $P_{max}$  round before seeing service

## Limitations of WFQ

Round computation expensive: must re-compute  $R$  every time  $N_{ac}(t)$  changes

Packetization causes two kinds of unfairness:

- shift from fluid-flow round-robin: fairness "quantized" by minimum packet size (Absolute Fairness Bound)
- unfairness to some flows: if flows have different packet sizes (Relative Fairness Bound)