



Lic.(Tech.) Marko Luoma (5/38)

### Queues

- Congestion situations demand queue management to decide
  - When packets should be discarded
  - Which are the packets that should be discarded
- Prevalent solutions
  - Tail Drop
  - Random Early Detection (RED)
  - Random Early Detection In/Out (RIO)
  - Weighted Random Early Detection (WRED)



Lic.(Tech.) Marko Luoma (7/38)

# Random Early Detection

- RED is an active queue management algorithm (AQM), which aims to
  - Prevent global syncronisation
  - Offer better fairness among competing connections
  - Allow transient burst without packet loss
- Algorithm operates on the knowledge of current Qsize and average Qsize (avg)
  - Avg is updated on every arrival and departure from the actual queue



HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory

Lic.(Tech.) Marko Luoma (6/38)

# Tail Drop

- Simple algorithm:
  - If arriving packets sees a full queue it is discarded
  - Otherwise it is accepted to the end of queue
- Problem:
  - Poor fairness in distribution of buffer space
  - Unable to accommodate short transients when queue is almost full
    - Bursty discarding process leading to global syncronisation
- Global syncronisation is a process where large number of TCP connections syncronise their window control due to concurrent packet losses.
  - Packet losses tend to be bursty, therefore window decreases to one and halts the communication



HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory

Lic.(Tech.) Marko Luoma (8/38)

# RED

• Qsize is used to calculate average length of the queue:

Initial condition: avg(0) = 0 Count = -1When Qsize = 0:  $T_{idle} = T_{now}$ After every packet arrival: if Qsize(n) > 0:  $avg(n+1) = (1-\epsilon) \cdot avg(n)$ 

 $avg(n+1)=(1-\epsilon)\cdot avg(n)+\epsilon \cdot Qsize(n)$ else:

 $avg(n+1)=avg(n)\cdot(1-\epsilon)^{f(T_{now}-T_{idle})}$ 

If queue is empty, averaging is done based on the assumption that N packets have passed the algorithm before actual packet arrival. -> Decay of average during idle times







- He sends packets with rate Y bps
- If Y is greater than X, some packets are marked as out of profile.
  - Out of profile packets usually experience harsh treatment on contending situations
- Calculation of the average queue length is modified to take into accout number of packets with different markings:
  - In (green): Only green packets
  - In/Out (yellow): Green and yellow packets
  - Out (red): All packets in the queue



• Generally they can be divided into

- Work-conserving vs non-work-

- Time-based vs frame-based

- Continuous vs packetized - Priority vs no priority

Work- conserving Non- work- conserving

SCHEDULERS

Frame- based

Packet- per- packet

Sorted-priority

Fluid- flow

categories of

conserving

• Policy defines the amount of resources which are allocated to the connections / classes / aggregates for which single packets belong to.

with predefined policy.

- allocated to the connections.
- Other end is that no allocation is done and resources are shared on the basis of the need



Lic.(Tech.) Marko Luoma (17/38)

## Scheduling

- **Conservation of work** means that scheduler is executing its task as long as it has some work to do.
- Technically this means that there are packets in the queue which has to be sent into the link before scheduler can take a break i.e. change to the idle state.
- Non-work conserving scheduler can idle even though it has packets in the queue.

- Why we would want to have nonwork conserving scheduler ?
- Conservation of work means that packets are sent to the link even though receiver would prefer them to come a little bit later.
- This can happen with real-time applications which send packets with constant time intervals. However, network can multiplex them so that they form bursts. Non-work conserving scheduler may delay packets so that intervals structure is maintained throughout the network.



Lic.(Tech.) Marko Luoma (19/38)

# Scheduling

#### Time based scheduling

- Uses either arrival time. finishing time or both as a criteria for ordering
- Time may be virtual or real-time depending on scheduler time
- Virtual time is usually finishing time in ideal scheduler i.e. scheduler which is not packetized

#### Frame based scheduling

- Uses fixed frame which is partitioned for the scheduled packets based on their weights.
- During a rotation,
  - If there are enough tokens (partition + left overs), then packet is served.
  - · Otherwise tokens are added for the next round.
- A number of packets may be served from a single class if frame is big.

HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory

Lic.(Tech.) Marko Luoma (18/38)

## Scheduling

- **Continuous time** 
  - Scheduling decissions and calculations are done based on continuous time units
  - Fluid-Flow modeling packets are infinitesimally small
  - Assumes that number of packets could be served on same time (not possible)

- Packetized
  - Scheduling decissions and calculations are based on packet per packet analysis
  - Distorts fluid flow model







HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory

Lic.(Tech.) Marko Luoma (20/38)

# Scheduling

- Scheduling can happen:
  - Within one queue, sorting packets inside queue to appropriate transmission order
  - Between several queues, dispatching head of line packets from different queues
  - Hierarchically over several schedulers, combination of previous ones
- Many of scheduling algorithms can be used to produce QoS in each of these cases



# Scheduling

- First Come First Served (FCFS) is prevalent scheduling method in routers.
- FCFS uses arrival time information as sorting criteria for packet dispatching.
- FCFS is not able to offer any QoS as arrival time is the only parameter that has influence to the order of packets.



# Scheduling

- Prioritized ordering may lead to starvation of resources in low priority classes if traffic in high priority classes is not limited.
- This can be accomplished by using
  - Connection admission control
  - Over provisioning
  - Rate control
  - Modifying priority scheduler to take class rates into account (token based operation)
    - Hierarchical Token Bucket (HTB)

High priority







Lic.(Tech.) Marko Luoma (25/38)

Lic.(Tech.) Marko Luoma (27/38)

## Scheduling

- **Delay based** scheduling schemes (e.q. PDD, WPT, HPD) are based on the calculation of queueing delay
  - Long term
  - Short term
  - Combination of both
- Packets are transmitted on the order of
  - Absolute queueing delay
  - Relative queueing delay
    - · Queueing delays are normalized with differentiation factor

HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory Lic.(Tech.) Marko Luoma (26/38) Scheduling Generalized Processor Sharing is ideal fair queueing algorithm which is based on fluid flow model. • GPS provides service to the individual connections based on their weights. • GPS is work conserving scheduler and thus distributes excess capacity to connections which are able to utilize it. C2 Class Arrival TimeService Time C3 C1C2 C1 C3 C2 Departure tim 3 L . . -C1 C3 Scheduling time Arrival time Weights are all equal C2 C2 C3 HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory Lic.(Tech.) Marko Luoma (28/38) Scheduling

- Advantages of GPS are:
  - Fairness which it provides for the sharing connections

$$\frac{[Service(t, t+\Delta t)]_i}{[Service(t, t+\Delta t)]_j} \ge \frac{Weight_i}{Weight_j}$$

 Strict delay bound caused by scheduling when traffic is constrained by a token bucket of token rate r and bucket depth b

Service rate for connection i: 
$$r_i \ge \frac{Weight_i}{\sum_j Weight_j}$$
. Link Rate  
Delay for connection i:  $D_i \le \frac{b_i}{r_i}$ 

Remember these results were derived from the assumption that packets flow like fluid through the system i.e. there would be a dedicated link with capacity *r* between endpoints.

## Scheduling

• Disadvantages of GPS are:

HELSINKI UNIVERSITY OF TECHNOLOGY

Networking laboratory

- Departures from GPS are colliding which makes the use of GPS based scheduler impossible
  - However it may be used as background scheduler if collisions are resolved in some manner
- Heavy calculation of departure times
  - Departure time of every packet in scheduler changes whenever a packet arrives or departs the scheduler





Lic.(Tech.) Marko Luoma (33/38)

# Scheduling

- **Deficit Round Robin** is extention of WRR which takes account the packet size
- DRR uses a rotation where a frame of *N* bits is divided to indivivual connections in relation to their weights (quantums).
- Quantums which individual connections receive serve packets
  - If the quantum is small, many rotations are required to serve backlogged connection -> approximated WFQ
  - If the quantum is big, many packets can be served on one rotation ->
    resource usage differs from the policy
- DRR uses special counter for each backlogged connection which stores the information of received bits.
  - If connection gets to non backlogged state counter is cleared



Lic.(Tech.) Marko Luoma (34/38)

Lic.(Tech.) Marko Luoma (36/38)

rmediat

Class #2

Leaf Class #2.1.

Class #2.

## Scheduling

- Class Based Queueing is one form hierarchical scheduling
- In CBQ scheduling is divided into two cases:
  - Unregulated: When a class is scheduled by general scheduler
  - Regulated: When a class is scheduled by link share scheduler
- Class is regulated in situations when network is persistently contended and class has run over its limits
- · Actual implementation of scheduling is uniform

HELSINKI UNIVERSITY OF TECHNOLOGY

· Link sharing guidelines are based on

- Intermediate Classes form

 Leaf classes are actual queues with distinct traffic

logical groupings

· Protocols

Organisations

Link resources are on Root Class

Networking laboratory

tree like structure

- Both schedulers manipulate HOL packets <u>time to send</u> information which is then examined by actual dispatcher.
- CBQ uses different variants of round robin schedulers as a general scheduler

Scheduling

Top Level

Class #1.1

Class #1.1

Class #

Class #1.2.

• Link share scheduler is based on general rules supplied by user



Lic.(Tech.) Marko Luoma (35/38)

# Scheduling

- Advantage of CBQ is that scheduling during contention is easily manipulated to produce outcome which is not only based on time and priority information
- Disadvantage is that CBQ requires a lot of processing time when there are a lot
  of independent connections / classes
  - HTB is an option for CBQ
    - · Almost the same functionalities with less overhead



Lic.(Tech.) Marko Luoma (37/38)

# Scheduling

- CBQ has concept of **borrowing**:
  - If class has run over its limit but it has parent class which is not over its limit, it may borrow capacity from the parent
  - Borrowing may be limited to some level in link sharing tree (Top Level)
- Formal definition between regulated and un regulated follows from borrowing:
  - Class is unregulated if:
    - It is under its limit
      - or
    - It has parent below Top Level which is under its limit

HELSINKI UNIVERSITY OF TECHNOLOGY Networking laboratory

Lic.(Tech.) Marko Luoma (38/38)

### Summary

- There is a lot of room to make more intelligent and effective scheduling and queue management algorithms
  - Resource adaptation
    - Network status changes -> resource allocation policy changes
    - Delay control for real-time communication
    - P2P
  - Fairness issues
    - How to bring differentiation into the Internet traffic without too much complexity