

# How to Improve Traffic Management



- Improve TCP
  - Stay with end-point only architecture
- Enhance routers to help TCP
  - · Random Early Discard
- Enhance routers to control traffic
  - Rate limiting
  - Fair Queueing
- Provide QoS by limiting congestion

Pouter Mechanisms

Buffer management: when and which packet to drop?

Scheduling: which packet to transmit next?

Classifier management

Buffer management

Buffer management

#### Overview



- Queue management & RED
- Fair-queuing
- Why QOS?
- Integrated services

# **Queuing Disciplines**



- Each router must implement some queuing discipline
- Queuing allocates both bandwidth and buffer space:
  - Bandwidth: which packet to serve (transmit) next
  - Buffer space: which packet to drop next (when required)
- Queuing also affects latency

### **Typical Internet Queuing**



- FIFO + drop-tail
  - · Simplest choice used widely in the Internet
- FIFO (first-in-first-out)
  - · Implies single class of traffic
  - · All packets treated equally
- Drop-tail
  - Arriving packets get dropped when queue is full regardless of flow or importance
- Important distinction:
  - FIFO: scheduling discipline
  - Drop-tail: drop policy

# FIFO + Drop-tail Problems



- Leaves responsibility of congestion control completely to the edges (e.g., TCP)
- Does not separate between different flows
- No policing: send more packets → get more service
- Synchronization: multiple end hosts may react to same events at about the same time
  - E.g. due to busts of packet loss

10

# FIFO + Drop-tail Problems



- Full queues
  - Routers are forced to have have large queues to maintain high utilizations
  - TCP detects congestion from loss
    - Forces network to have long standing queues in steady-state
    - May no longer be able to absorb bursts
- Lock-out problem
  - Drop-tail routers treat bursty traffic poorly
  - Traffic gets synchronized easily → allows a few flows to monopolize the queue space

1

# Active Queue Management



- Design active router queue management to aid congestion control
- Why?
- Router has unified view of queuing behavior
  - Routers see actual queue occupancy (distinguish queue delay and propagation delay)
  - Routers can decide on transient congestion, based on workload

## **Design Objectives**



- · Keep throughput high and delay low
  - High power (throughput/delay)
- Accommodate bursts
- Queue size should reflect ability to accept bursts rather than steady-state queuing
- Improve TCP performance with minimal hardware changes

13

#### Lock-out Problem



- Random drop
  - Packet arriving when queue is full causes some random packet to be dropped
- Drop front
  - On full queue, drop packet at head of queue
- Random drop and drop front solve the lock-out problem but not the full-queues problem

14

# Full Queues Problem



- Drop packets before queue becomes full (early drop)
- Intuition: notify senders of incipient congestion
- Example: early random drop (ERD):
  - If qlen > drop level, drop new packets with fixed probability p
  - Does not control misbehaving users

6

# Random Early Detection (RED)



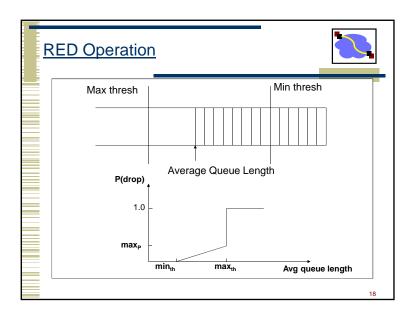
- Detect incipient congestion
- Assume hosts respond to lost packets
- Avoid window synchronization
  - Randomly mark packets
- · Avoid bias against bursty traffic

# **RED Algorithm**



- Maintain running average of queue length
- If avg < min<sub>th</sub> do nothing
  - Low queuing, send packets through
- If avg > max<sub>th</sub>, drop packet
  - · Protection from misbehaving sources
- Else drop packet in a manner proportional to queue length
  - Notify sources of incipient congestion
  - Can instead mark the packet if ECN is used (later)

17



# Explicit Congestion Notification (ECN) [Floyd and Ramakrishnan 98]



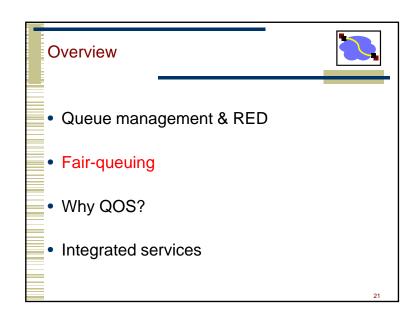
- Traditional mechanism: packet drop as implicit congestion signal to end systems
  - TCP will slow down
- · Works well for bulk data transfer
- Does not work well for delay sensitive applications
  - audio, WEB, telnet
- Explicit Congestion Notification (ECN)
  - · borrow ideas from DECBit
  - · use two bits in IP header
    - · ECN-Capable Transport (ECT) bit set by sender
    - · Congestion Experienced (CE) bit set by router

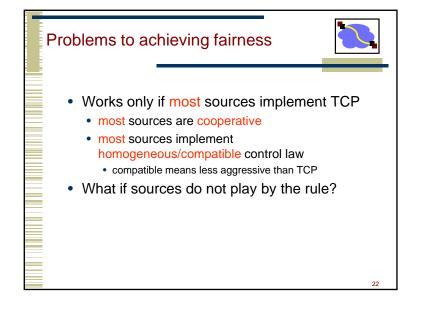
0

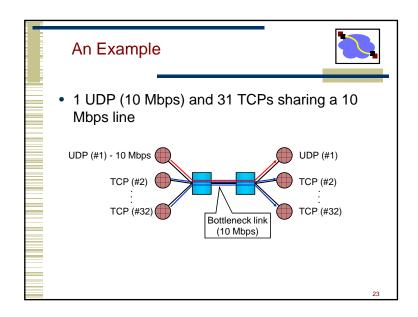
# **Congestion Control Summary**

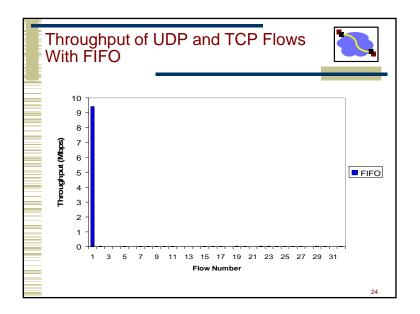


- Basis: end system detects congestion and slow down
  - slow start/congestion avoidance
    - packet drop detected by retransmission timeout RTO as congestion signal
    - AIMD: search for right bandwidth, considering fairness
- TCP Improvement:
  - fast retransmission/fast recovery
    - packet drop detected by three duplicate acks
  - · SACK: better feedback to source
- Router support
  - RED: early signaling
  - ECN: explicit signaling









#### Fairness Goals



- · Allocate resources fairly
- Isolate ill-behaved users
  - Router does not send explicit feedback to source
  - Still needs e2e congestion control
- · Still achieve statistical multiplexing
  - One flow can fill entire pipe if no contenders
  - Work conserving → scheduler never idles link if it has a packet

25

#### What is Fairness?



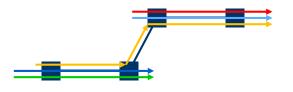
- At what granularity?
  - Flows, connections, domains?
- What if users have different RTTs/links/etc.
  - Should it share a link fairly or be TCP fair?
- Maximize fairness index?
  - Fairness =  $(\Sigma x_i)^2/n(\Sigma x_i^2)$  0<fairness<1
- · Basically a tough question to answer!
- Good to separate the design of the mechanisms from definition of a policy
  - User = arbitrary granularity

26

#### Max-min Fairness



- Allocate user with "small" demand what it wants, evenly divide unused resources to "big" users
- Formally:
  - · Resources allocated in terms of increasing demand
  - No source gets resource share larger than its demand
  - · Sources with unsatisfied demands get equal share of resource



27

# Implementing Max-min Fairness



- Generalized processor sharing
  - Fluid fairness
  - · Bitwise round robin among all queues
- Why not simple round robin?
  - Variable packet length → can get more service by sending bigger packets
  - Unfair instantaneous service rate
    - What if arrive just before/after packet departs?

# Bit-by-bit RR

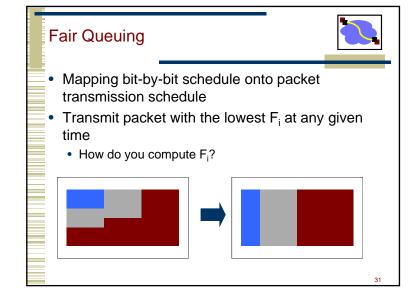


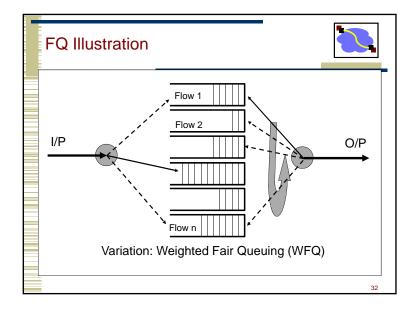
- Single flow: clock ticks when a bit is transmitted. For packet i:
  - P<sub>i</sub> = length, A<sub>i</sub> = arrival time, S<sub>i</sub> = begin transmit time, F<sub>i</sub> = finish transmit time
  - $F_i = S_i + P_i = \max(F_{i-1}, A_i) + P_i$
- Multiple flows: clock ticks when a bit from all active flows is transmitted → round number
  - Can calculate F<sub>i</sub> for each packet if number of flows is know at all times
    - Why do we need to know flow count? → need to know A → This can be complicated

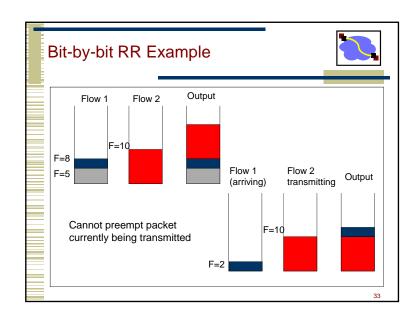
Pit-by-bit RR Illustration

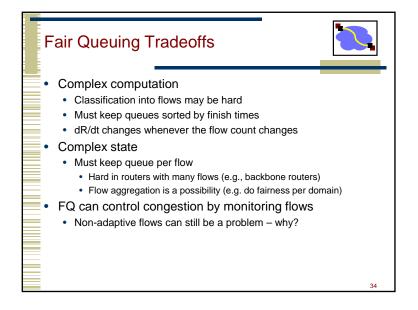
Not feasible to interleave bits on real networks

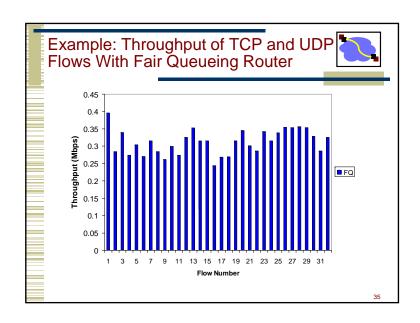
FQ simulates bit-by-bit RR

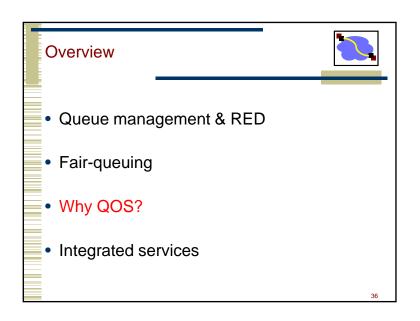












#### Motivation



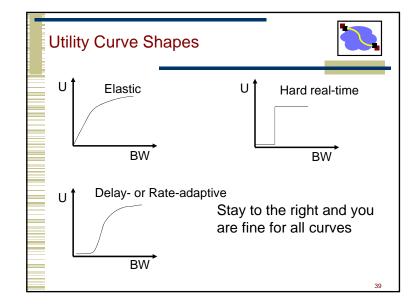
- Internet currently provides one single class of "best-effort" service
  - No assurances about delivery
  - All users (packets?) are equal
- At internet design most applications are *elastic* 
  - Tolerate delays and losses
  - · Can adapt to congestion
- Today, many "real-time" applications are inelastic

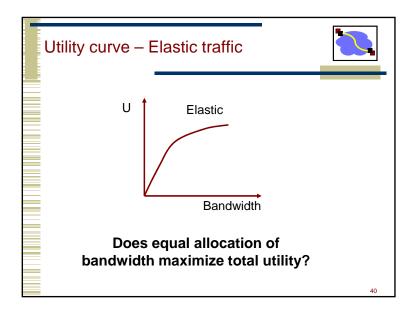
37

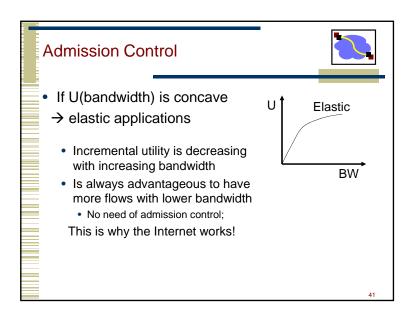
# Why a New Service Model?

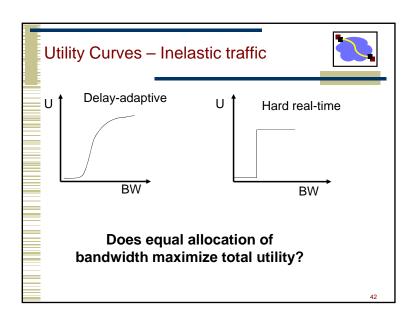


- What is the basic objective of network design?
  - Maximize total bandwidth? Minimize latency?
  - Maximize user satisfaction the total utility given to users
- What does utility vs. bandwidth look like?
  - Shape depends on application
  - Must be non-decreasing function

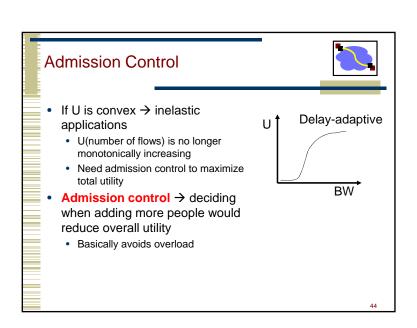








# Continuous media applications Lower and upper limit on acceptable performance. BW below which video and audio are not intelligible Internet telephones, teleconferencing with high delay (200 - 300ms) impair human interaction Sometimes called "tolerant real-time" since they can adapt to the performance of the network Hard real-time applications Require hard limits on performance E.g. control applications



#### Overview



- Queue management & RED
- Fair-queuing
- Why QOS?
- Integrated services

45

# Components of Integrated Services



- 1. Type of commitment
  - What does the network promise?
- 2. Packet scheduling

How does the network meet promises?

3. Service interface

How does the application describe what it wants?

4. Establishing the guarantee

How is the promise communicated to/from the network How is admission of new applications controlled?

46

# Type of Commitments



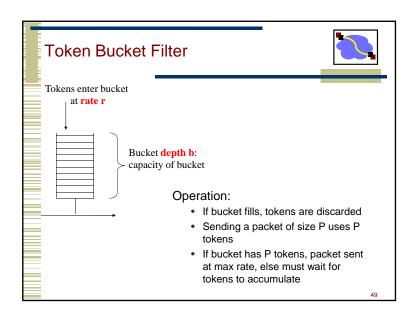
- Guaranteed service
  - For hard real-time applications
  - Fixed guarantee, network meets commitment if clients send at agreed-upon rate
- Predicted service
  - For delay-adaptive applications
  - Two components
    - · If conditions do not change, commit to current service
    - If conditions change, take steps to deliver consistent performance (help apps minimize playback delay)
    - Implicit assumption network does not change much over time
- Datagram/best effort service

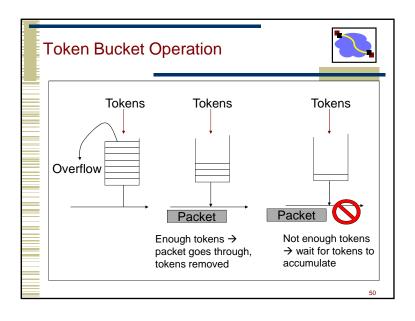
47

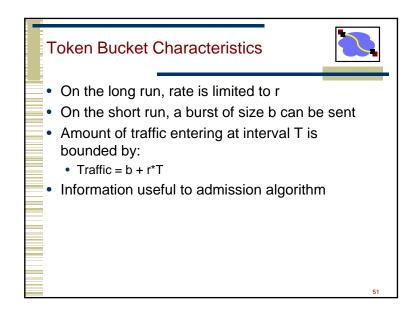
# Scheduling for Guaranteed Traffic

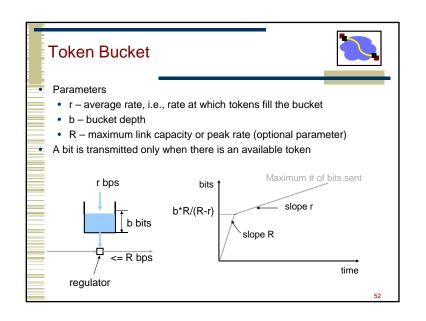


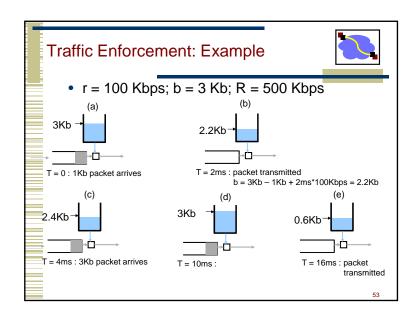
- Use token bucket filter to characterize traffic
  - Described by rate r and bucket depth b
- Use Weighted Fair-Queueing at the routers
- Parekh's bound for worst case queuing delay = b/r

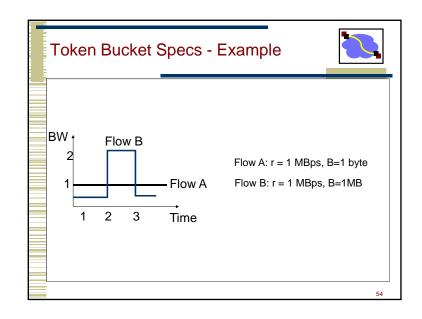












# Guarantee Proven by Parekh



- · Given:
  - Flow i shaped with token bucket and leaky bucket rate control (depth b and rate r)
  - · Network nodes do WFQ
- Cumulative queuing delay D<sub>i</sub> suffered by flow i has upper bound
  - **D**<sub>i</sub> < **b/r**, (where r may be much larger than average rate)
  - Assumes that  $\Sigma r < \text{link speed at any router}$
  - All sources limiting themselves to r will result in no network queuing

Sharing versus Isolation



- Impact of queueing mechanisms:
  - Isolation: Isolates well-behaved from misbehaving sources
  - · Sharing: Mixing of different sources in a way beneficial to all
- FIFO: sharing
  - · each traffic source impacts other connections directly
    - · e.g. malicious user can grab extra bandwidth
  - the simplest and most common queueing discipline
  - · averages out the delay across all flows
- Priority gueues: one-way sharing
  - high-priority traffic sources have impact on lower priority traffic only
  - has to be combined with admission control and traffic enforcement to avoid starvation of low-priority traffic
- WFQ: two-way isolation
  - provides a guaranteed minimum throughput (and maximum delay)

# Putting It All Together



- Assume 3 types of traffic: guaranteed, predictive, best-effort
- Scheduling: use WFQ in routers
- Each guaranteed flow gets its own queue
- All predicted service flows and best effort aggregates in single separate queue
  - · Predictive traffic classes
    - Worst case delay for classes separated by order of magnitude
    - When high priority needs extra bandwidth steals it from lower class
  - · Best effort traffic acts as lowest priority class

67

#### Lessons



- TCP can use help from routers
  - RED → eliminate lock-out and full-queues problems
  - FQ  $\rightarrow$  heavy-weight but explicitly fair to all
- QoS
  - What type of applications are there? → Elastic, adaptive real-time , and hard real-time.
  - Why do we need admission control → to maximize utility
  - How do token buckets + WFQ provide QoS guarantees?

#### Service Interfaces



- Guaranteed Traffic
  - · Host specifies rate to network
  - Why not bucket size b?
    - If delay not good, ask for higher rate
- Predicted Traffic
  - Specifies (r, b) token bucket parameters
  - · Specifies delay D and loss rate L
  - · Network assigns priority class
  - · Policing at edges to drop or tag packets
    - Needed to provide isolation why is this not done for guaranteed traffic?
      - · WFQ provides this for guaranteed traffic