

ECE 356/COMPSI 356

Computer Network Architecture

Internet QoS

Wednesday November 13th, 2019

Recap

- Previous lecture: queuing and congestion avoidance
- Readings for this lecture: **PD 6.5.1, 6.5.3**

Lecture Outline

- Multimedia communications
- Internet QoS
- Coarse-grained QoS: differentiated services

Motivation

- Internet currently provides one single class of **“best-effort”** service
 - No assurance about delivery
- But different application classes have different needs
 - Should we aim to offer differentiation and guarantees?

Traditional Data Applications

- Web browsing, file transfer, e-mail
- “*Elastic*”
 - Can work without a guarantee of timely delivery of data
 - Benefit from shorter delays, but do not become unstable as delays increase
- Not loss-tolerant

Multimedia Networking

The Netflix logo, consisting of the word "NETFLIX" in red, uppercase letters on a black rectangular background.The YouTube logo, featuring a red play button icon inside a white rounded square, followed by the word "YouTube" in black.The Hulu logo, with the word "hulu" in a lowercase, green, sans-serif font.The Skype logo, with the word "skype" in a white, lowercase, sans-serif font inside a blue, cloud-like shape.The Cisco Webex logo, featuring a blue and green circular icon above the word "Cisco" in small black text and "webex" in a larger, lowercase, blue and green font.

Properties of Video (1/3)

- One of primary properties: **high bit rate**

	Bit rate	Bytes transferred in 67 min
Facebook browsing	160 kbps	80 MB
Spotify audio streaming	128 kbps	64 MB
Video streaming	2 Mbps	1 Gb

- Facebook browsing: a new photo every 10 s, photos are 200KB in size on average
- Requirements get higher and higher as video improves

Properties of Video (1/3)

- Even higher bit rates: virtual reality
 - 360 degree videos
 - Much higher frame rates
- RGB-D video
 - Depth data: another stream



Video Deployments: Current

- One camera installed for every 29 people on the planet
 - One for every 8 people in mature markets
- Wide range of applications
 - Traffic control
 - Surveillance in public and private spaces
 - ...



9

Properties of Video: Compression

- Video: a sequence of images displayed at a constant rate, e.g., 24 or 30 images per second
- Digital image: array of pixels
 - Each pixel represented by bits
- Coding: use redundancy within and between images to decrease # bits used to encode image
 - Spatial (within image)
 - Temporal (from one image to next)
- Can compress the video to almost any bit rate
 - The higher the bit rate, the better user viewing experience

10

Video Compression Examples

Spatial coding example: instead of sending N values of same color (all purple), send only two values: color value (*purple*) and number of repeated values (N)



frame i

Temporal coding example: instead of sending complete frame at $i+1$, send only differences from frame i



frame $i+1$

11

Properties of Video: Multiple Versions of the Same Video

- Use compression to create multiple versions of a video, with different quality levels
 - E.g., 300 kbps, 1 Mbps, 3 Mbps
- Users can decide which quality to choose
- Applications adapt quality to available bandwidth

12

Properties of Audio

- Also can be compressed to multiple levels
 - Human speed is intelligible when compressed to under 10 kpbs
 - Common encoding rate: 128 kpbs
- Users are more sensitive to audio glitches than to video glitches
 - E.g., a video conference can be OK if video feed is lost once in a while, but would likely be terminated if audio is not getting through

13

Types of Multimedia Network Applications: Streaming *Stored* Audio and Video

- **Streaming**: can begin playout before downloading the entire file
- *Stored* (at a server): can transmit faster than audio/video will be rendered (implies storing/buffering at client)
- Interactivity: user may pause or reposition content
 - Need to react to the user with sufficiently low latency
- Continuous playout: data must be received from the server in time for its playout at the client

14

Types of Multimedia Network Applications: Conversational Voice and Video-over-IP

- **Highly delay-sensitive**
 - Interactive nature of human-to-human conversation limits delay tolerance
 - A few 100 ms at most
 - E.g., for voice, 150 ms is not perceived, 150 – 400 ms is acceptable, 400 ms + is frustrating and potentially unintelligible
- **Loss-tolerant**
 - In contrast with elastic data applications



15

Types of Multimedia Network Applications: Streaming *Live* Audio and Video

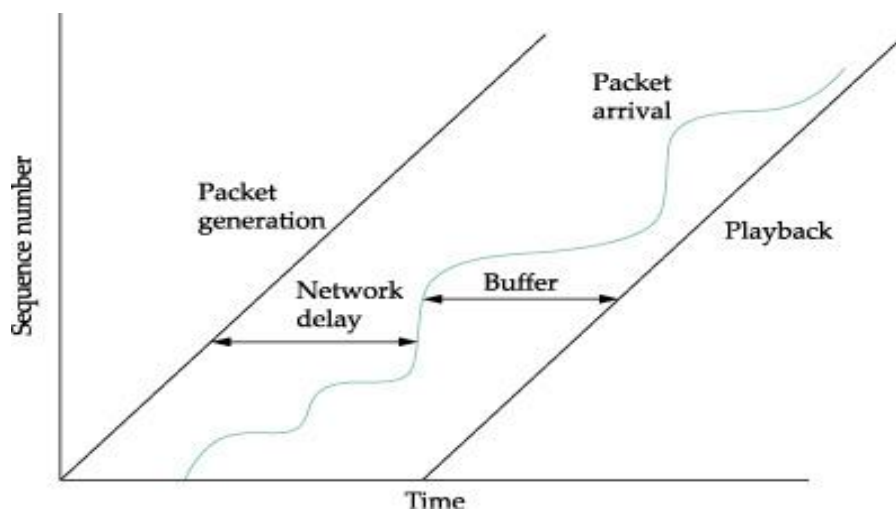
- Sports event, news event
- Usually transmitted to many users simultaneously
- Less stringent requirements than conversational multimedia
- Delays can be an issue
 - Delays of up to ~ 10s can be tolerated

16

Playback Applications



- Sample signal \rightarrow packetize \rightarrow transmit \rightarrow buffer \rightarrow playback
 - Fits most multimedia applications
- Performance concern:
 - Jitter: variation in end-to-end delay
 - Delay = fixed + variable = (propagation + packetization) + queuing
- Solution:
 - Playback point – delay introduced by buffer to hide network jitter



Characteristics of Playback Applications

- In general lower delay is preferable
- Doesn't matter when packet arrives as long as it is before playback point
- Network guarantees (e.g., bound on jitter) would make it easier to set playback point
- Applications can tolerate some loss

Lecture Outline

- Multimedia communications
- **Internet QoS**
- Coarse-grained QoS: differentiated services

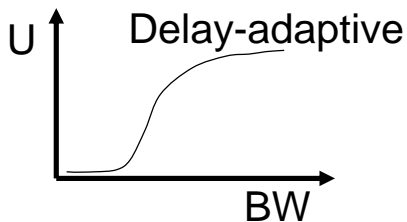
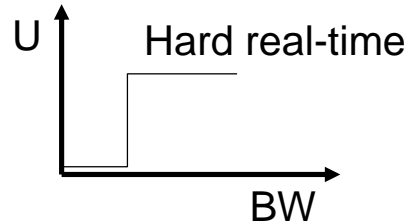
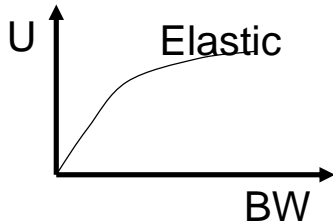
Inelastic Applications

- Continuous media applications
 - **Lower and upper limit** on acceptable performance
 - Below which video and audio are not intelligible
 - Internet telephones, teleconferencing with high delay (200 - 300ms) impair human interactions
- Hard real-time applications
 - Require **hard limits on performance**
 - E.g., industrial control applications
 - Internet surgery

Design question #1: Why a New Service Model?

- What is the **basic objective** of network design?
 - Maximize total bandwidth? Minimize latency? Maximize ISP's revenues?
 - **The designer's choice: maximize social welfare:** the total **utility** given to users (why not profit?)
- What does utility vs. bandwidth look like?
 - Must be non-decreasing function
 - Shape depends on application

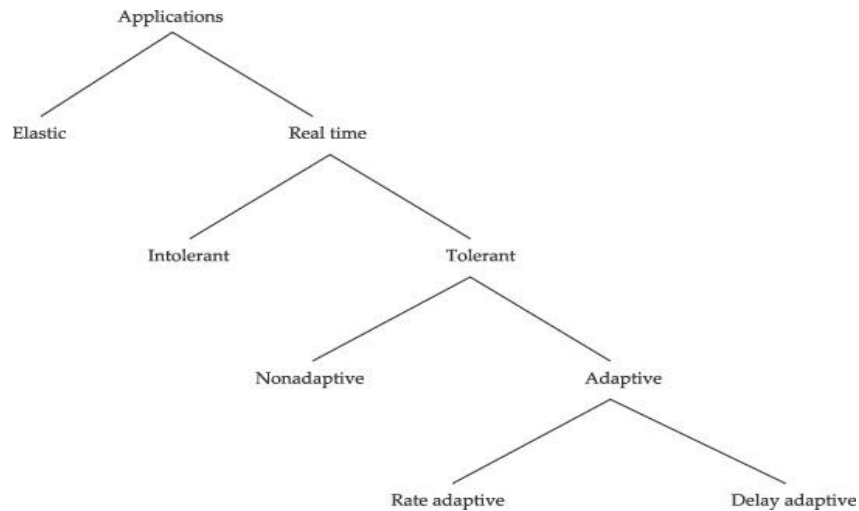
Utility Curve Shapes



- Stay to the right and you are fine for all curves

Applications Variations

- Rigid and adaptive applications
 - Delay adaptive
 - Rigid: set fixed playback point
 - Adaptive: adapt playback point
 - E.g. Shortening silence for voice applications
 - Rate adaptive
- Loss tolerant and intolerant applications
- Four combinations



Applications Variations

Really only two classes of applications

- 1) Intolerant and rigid
- 2) Tolerant and adaptive

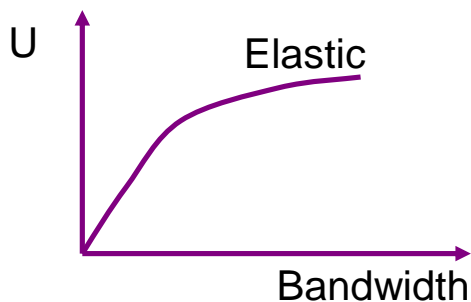
Other combinations make little sense

- 3) Intolerant and adaptive
 - Cannot adapt without interruption
- 4) Tolerant and rigid
 - Missed opportunity to improve delay

Design Question 2: How to maximize $V = \sum U(s_i)$

- Choice #1: add more pipes
- Choice #2: fix the bandwidth but offer different services
 - Q: can differentiated services improve V ?

If all Users' Utility Functions are Elastic

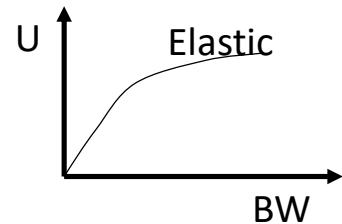


Does equal allocation of bandwidth maximize total utility?

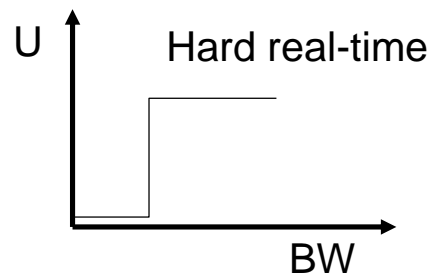
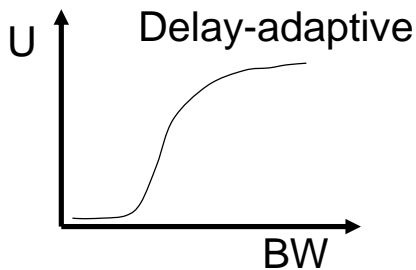
- $\sum s_i = B$
- $\text{Max } \sum U(s_i)$

Design Question: Is Admission Control Needed?

- If $U(\text{bandwidth})$ is concave
 - elastic applications
 - Incremental utility is decreasing with increasing bandwidth
 - $U(x) = \log(x^p)$
 - $V = n \log(B/n)^p = \log B^p n^{1-p}$
 - Is always advantageous to have more flows with lower bandwidth
 - No need of admission control;
- This is why the Internet works! And fairness makes sense



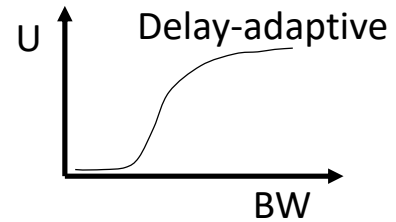
Utility Curves – Inelastic traffic



Does equal allocation of bandwidth maximize total utility?

Is Admission Control needed?

- If U is convex \rightarrow inelastic applications
 - U (number of flows) is no longer monotonically increasing
 - Need admission control to maximize total utility
- **Admission control** \rightarrow deciding when the addition of new people would result in reduction of utility
 - Basically avoids overload



Incentives

- Who should be given what service?
 - Users have incentives to cheat
 - Pricing seems to be a reasonable choice
 - But usage-based charging may not be well received by users

Over-provisioning

- Pros: simple
- Cons
 - Not cost effective
 - Bursty traffic leads to a high peak/average ratio
 - E.g., normal users versus leading edge users
 - It might be easier to block heavy users

Comments

- End-to-end QoS has not happened
- Why?
- Can you think of any mechanism to make it happen?

Approaches to QoS

- Fine-grained:
 - Integrated services
 - RSVP
- Coarse-grained:
 - Differentiated services

Lecture Outline

- Multimedia communications and Internet QoS
- **Coarse-grained QoS: differentiated services**

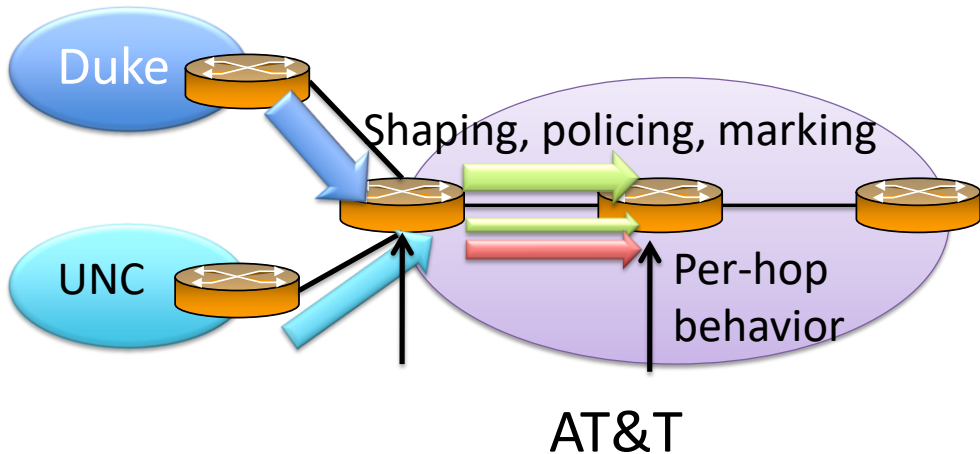
Motivation of DiffServ

- Analogy:
 - Airline service, first class, coach, various restrictions on coach as a function of payment
- Economics and assurances
 - Pay more, and get better service
 - Best-effort expected to make up bulk of traffic,
 - Revenue from first class important to economic base
 - Not motivated by real-time or maximizing social welfare

Basic Architecture

- Agreements/service provided within a domain
 - Service Level Agreement (SLA) with ISP
- Edge routers do traffic conditioning
 - Shaping, Policing, and Marking
- Core routers
 - Process packets based on packet marking and defined per hop behavior (PHB)
- More scalable than IntServ
 - No per flow state or signaling

DiffServ Architecture Example



Per-hop Behaviors (PHBs)

- Define behavior of individual routers rather than end-to-end services; there may be many more services than behaviors
 - No end-to-end guarantee
- Multiple behaviors – need more than one bit in the header
- Six bits from IP TOS field are taken for Diffserv code points (DSCP)

Per-hop Behaviors (PHBs)

- Two PHBs defined so far
- Expedited forwarding aka premium service (type P)
 - Possible service: providing a virtual wire
- Assured forwarding (type A)
 - Possible service: strong assurance for traffic within profile and allow source to exceed profile

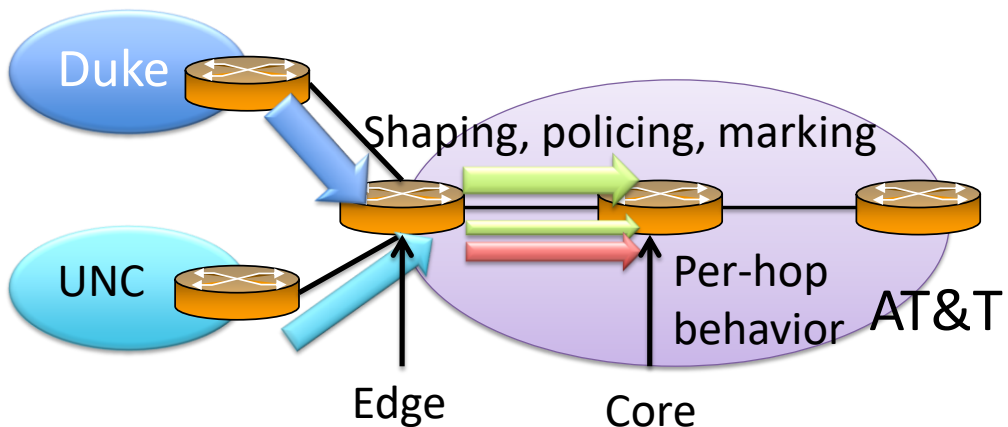
Expedited Forwarding PHB

- Goal: EF packets are forwarded with minimal delay and loss
- Mechanisms:
 - User sends within profile and network commits to delivery with requested profile
 - Rate limiting of EF packets at edges only, using token bucket to shape transmission
 - Priority or Weighted Fair Queuing

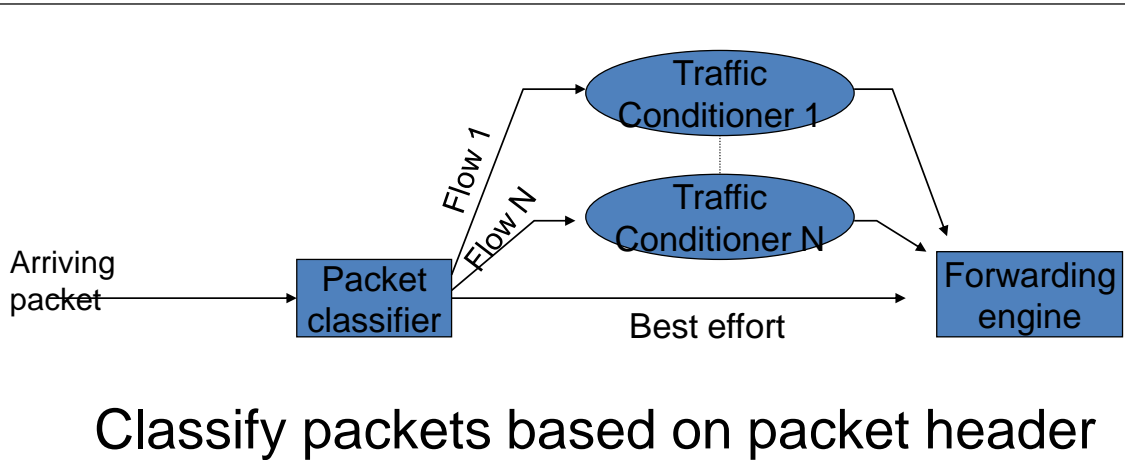
Assured Forwarding PHB

- Goal: good services for in-profile traffic
- Mechanisms:
 - User and network agree to some traffic profile
 - How to define profiles is an open/policy issue
 - Edges mark packets up to allowed rate as “in-profile” or low drop precedence
 - Other packets are marked with one of two higher drop precedence values
 - Random Early Detection in/out queues

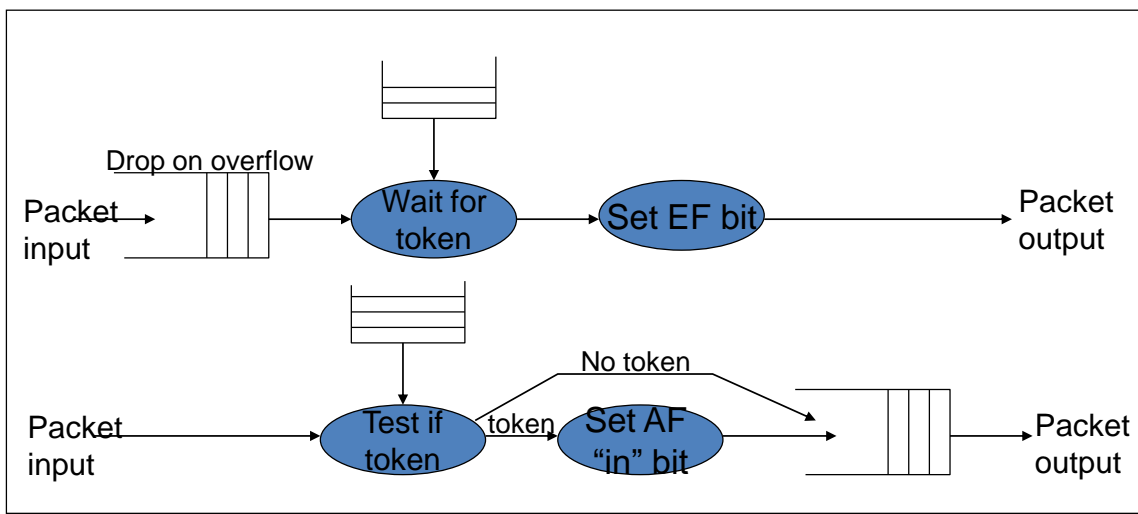
DiffServ Architecture Example



Edge Router Input Functionality

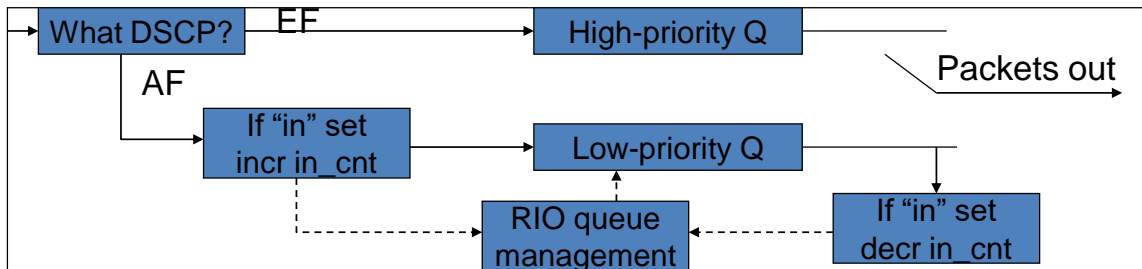


Traffic Conditioning



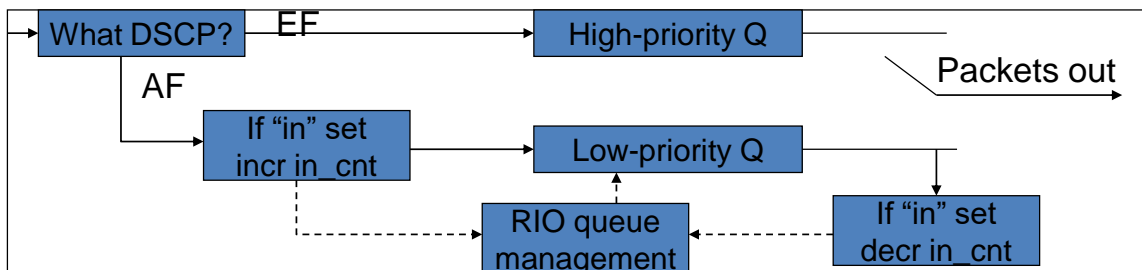
Router Output Processing

- Two queues: EF packets on higher priority queue
- Lower priority queue implements RED “In or Out” scheme (RIO)



Router Output Processing

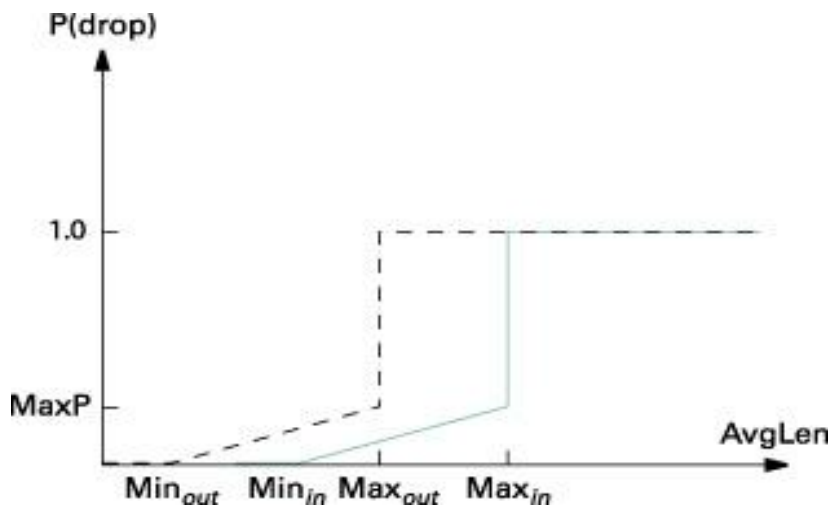
- Two queues: EF packets on higher priority queue
- Lower priority queue implements RED “In or Out” scheme (RIO)



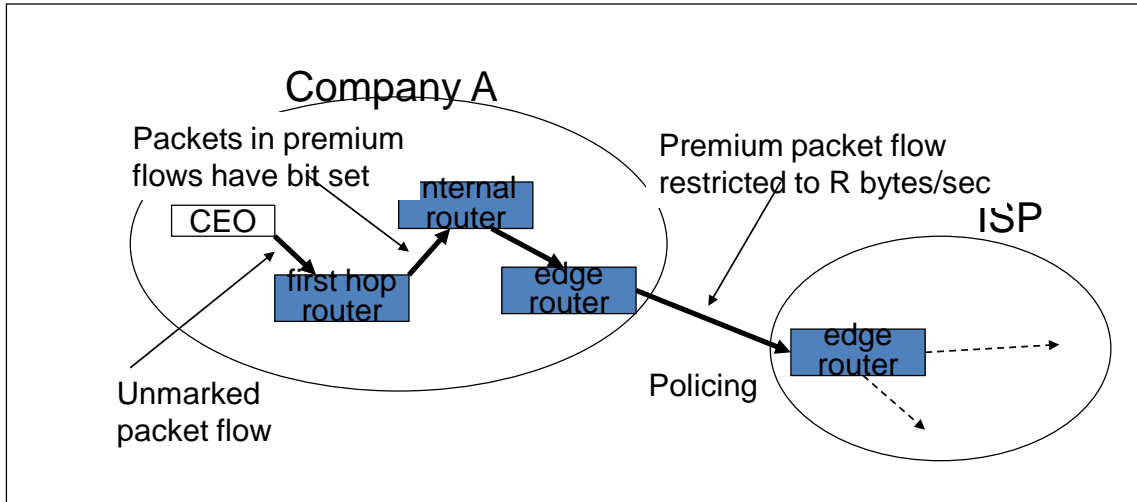
Red with In or Out (RIO)

- Similar to RED, but with two separate probability curves
- Has two classes, “In” and “Out” (of profile)
- “Out” class has lower $\text{Min}_{\text{thresh}}$, so packets are dropped from this class first
 - Based on queue length of all packets
- As avg queue length increases, “in” packets are also dropped
 - Based on queue length of only “in” packets

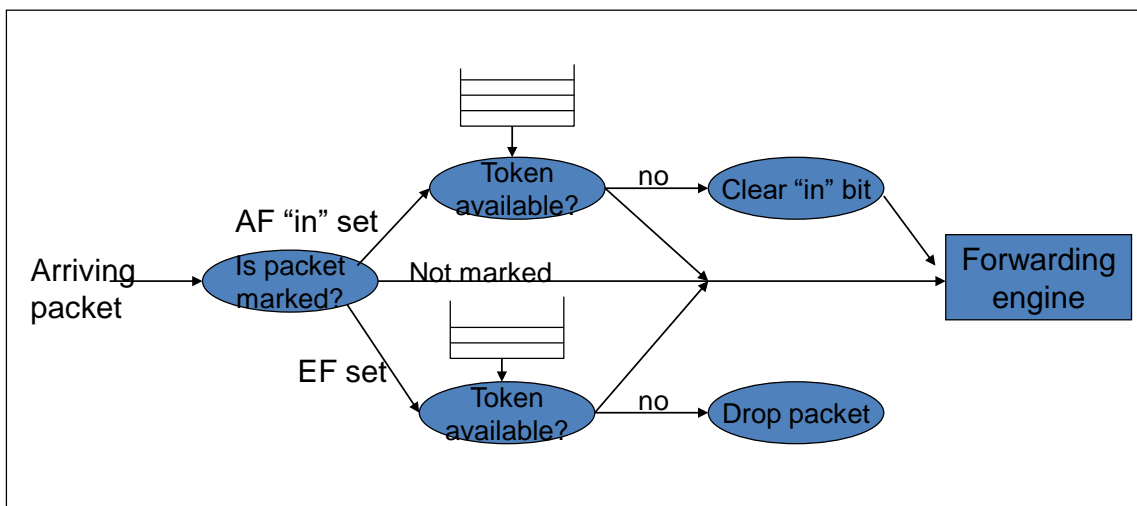
RIO Drop Probabilities



Pre-marking and Traffic Conditioning



Edge Router Policing



Diffserv Service Model: Observations

- End-to-end service must be fashioned from multiple ISPs
 - ISPs need to cooperate
- With Diffserv in place, if networks run at a moderate load, *most of the time there would be no perceived difference between a best-effort service and a Diffserv service*
 - End-to-end delays are usually dominated by access rates and router hops rather than router queuing delays
 - Not a great business model if you want to charge extra for priority service

DiffServ: Key Points to Remember

QoS Deployment

- “Dead” at the Internet scale
- Areas of success
 - Enterprise networks
 - Residential uplinks
 - Datacenter networks
- Ideas keep surfacing for

Lecture Summary

- QoS
 - Why do we need it?
 - Differentiated Services
 - Motivated by business models

Next Lecture

- DNS and content distribution

57