

# Chapter 7

## Multimedia Networking

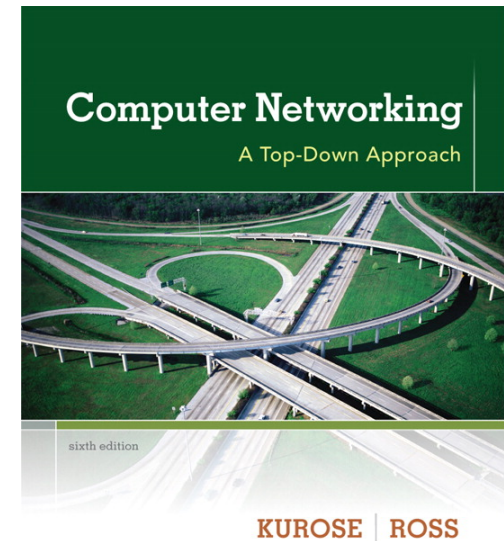
### A note on the use of these ppt slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a lot of work on our part. In return for use, we only ask the following:

- ❖ If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- ❖ If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

© All material copyright 1996-2012  
J.F Kurose and K.W. Ross, All Rights Reserved



*Computer  
Networking: A  
Top Down  
Approach  
6<sup>th</sup> edition  
Jim Kurose, Keith Ross  
Addison-Wesley  
March 2012*

# Multimedia networking: outline

- 7.1 multimedia networking applications
- 7.2 streaming *stored* video
- 7.3 voice-over-IP
- 7.4 protocols for *real-time* conversational applications
- 7.5 network support for multimedia

# Multimedia networking: outline

7.1 multimedia networking applications

7.2 streaming *stored* video

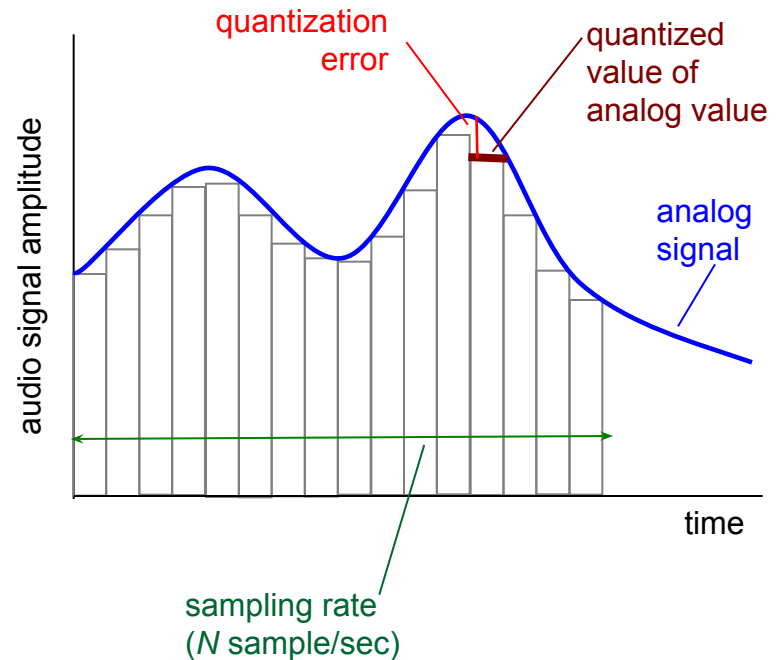
7.3 voice-over-IP

7.4 protocols for *real-time* conversational applications

7.5 network support for multimedia

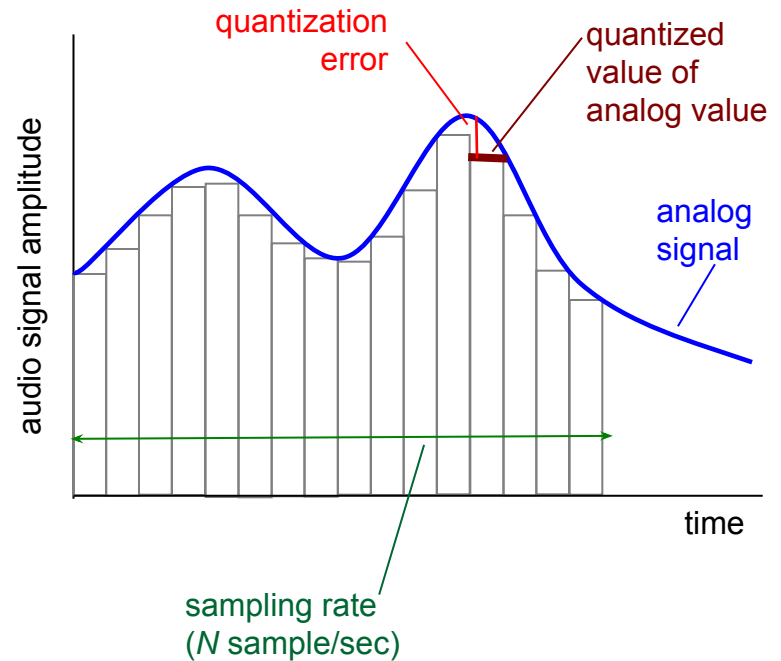
# Multimedia: audio

- ❖ analog audio signal sampled at constant rate
  - telephone: 8,000 samples/sec
  - CD music: 44,100 samples/sec
- ❖ each sample quantized, i.e., rounded
  - e.g.,  $2^8=256$  possible quantized values
  - each quantized value represented by bits, e.g., 8 bits for 256 values



# Multimedia: audio

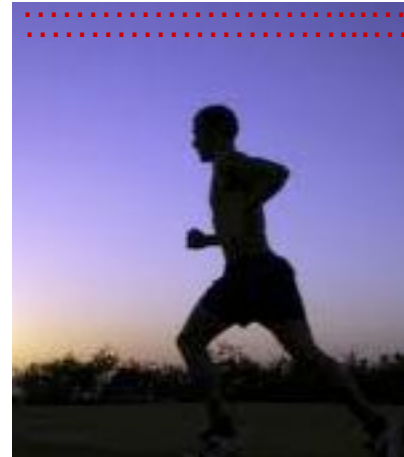
- ❖ example: 8,000 samples/sec, 256 quantized values: 64,000 bps
- ❖ receiver converts bits back to analog signal:
  - some quality reduction



# Multimedia: video

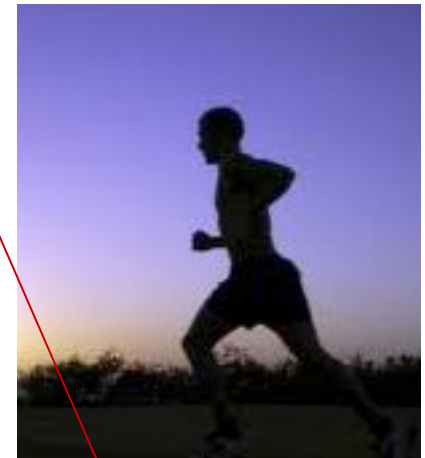
- ❖ video: sequence of images displayed at constant rate
  - e.g. 24 images/sec
- ❖ digital image: array of pixels
  - each pixel represented by bits
- ❖ coding: use redundancy *within* and *between* images to decrease # bits used to encode image
  - spatial (within image)
  - temporal (from one image to next)

*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  $i+1$

# Multimedia: video

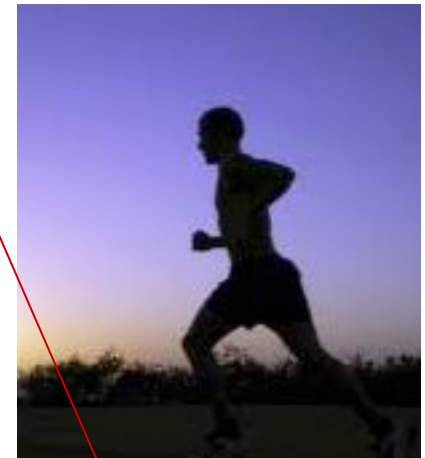
- ❖ **CBR: (constant bit rate):**  
video encoding rate fixed
- ❖ **VBR: (variable bit rate):**  
video encoding rate changes  
as amount of spatial,  
temporal coding changes
- ❖ **examples:**
  - MPEG I (CD-ROM) 1.5 Mbps
  - MPEG2 (DVD) 3-6 Mbps
  - MPEG4 (often used in Internet, < 1 Mbps)

*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  $i+1$

# Multimedia networking: 3 application types

---

- ❖ *streaming, stored* audio, video
  - *streaming*: can begin playout before downloading entire file
  - *stored (at server)*: can transmit faster than audio/video will be rendered (implies storing/buffering at client)
  - e.g., YouTube,
- ❖ *conversational* voice/video over IP
  - interactive nature of human-to-human conversation limits delay tolerance
  - e.g., Skype
- ❖ *streaming live* audio, video
  - e.g., live sporting event (futbol)



# Multimedia networking: outline

7.1 multimedia networking applications

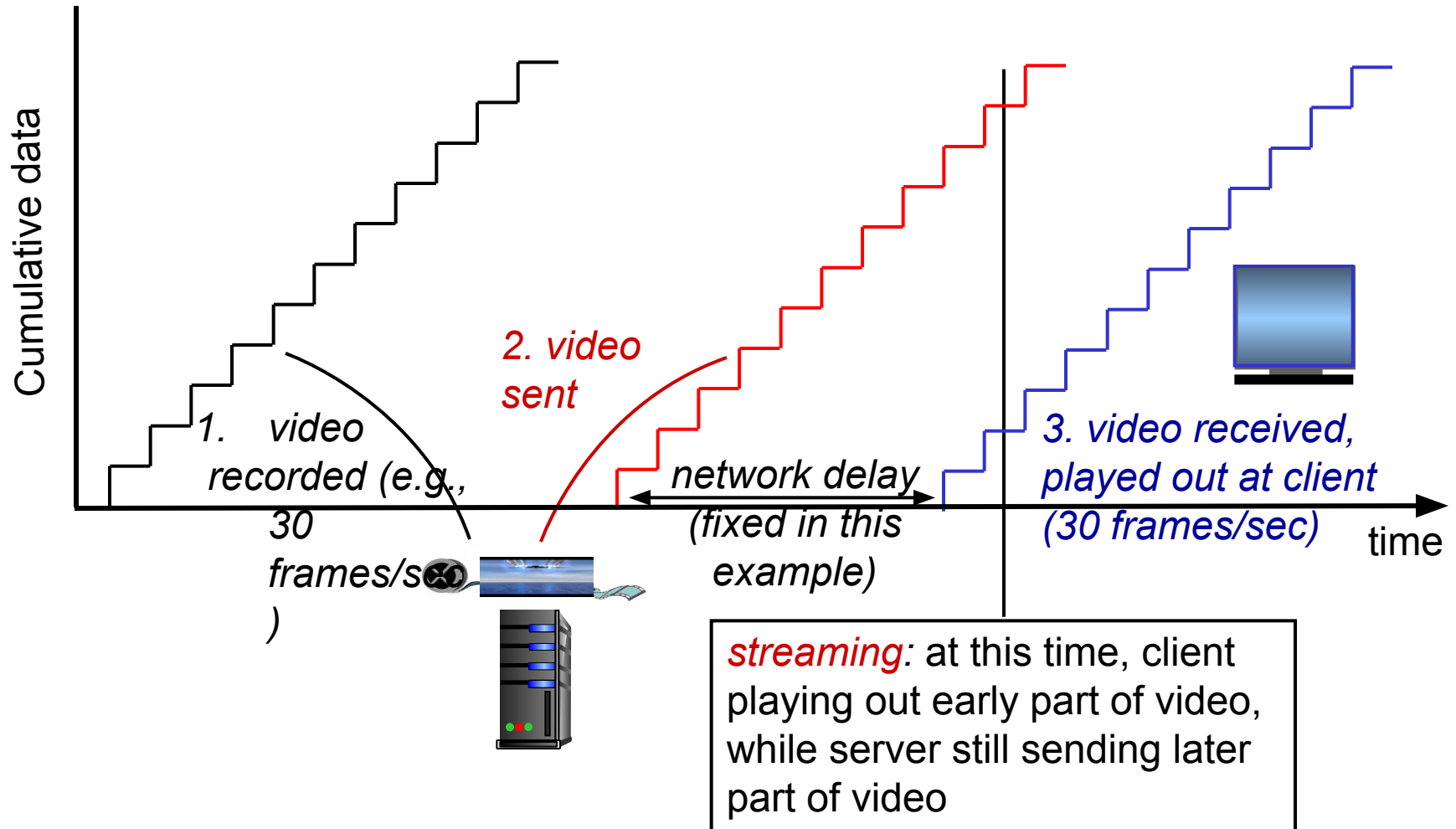
7.2 streaming *stored* video

7.3 voice-over-IP

7.4 protocols for *real-time* conversational applications

7.5 network support for multimedia

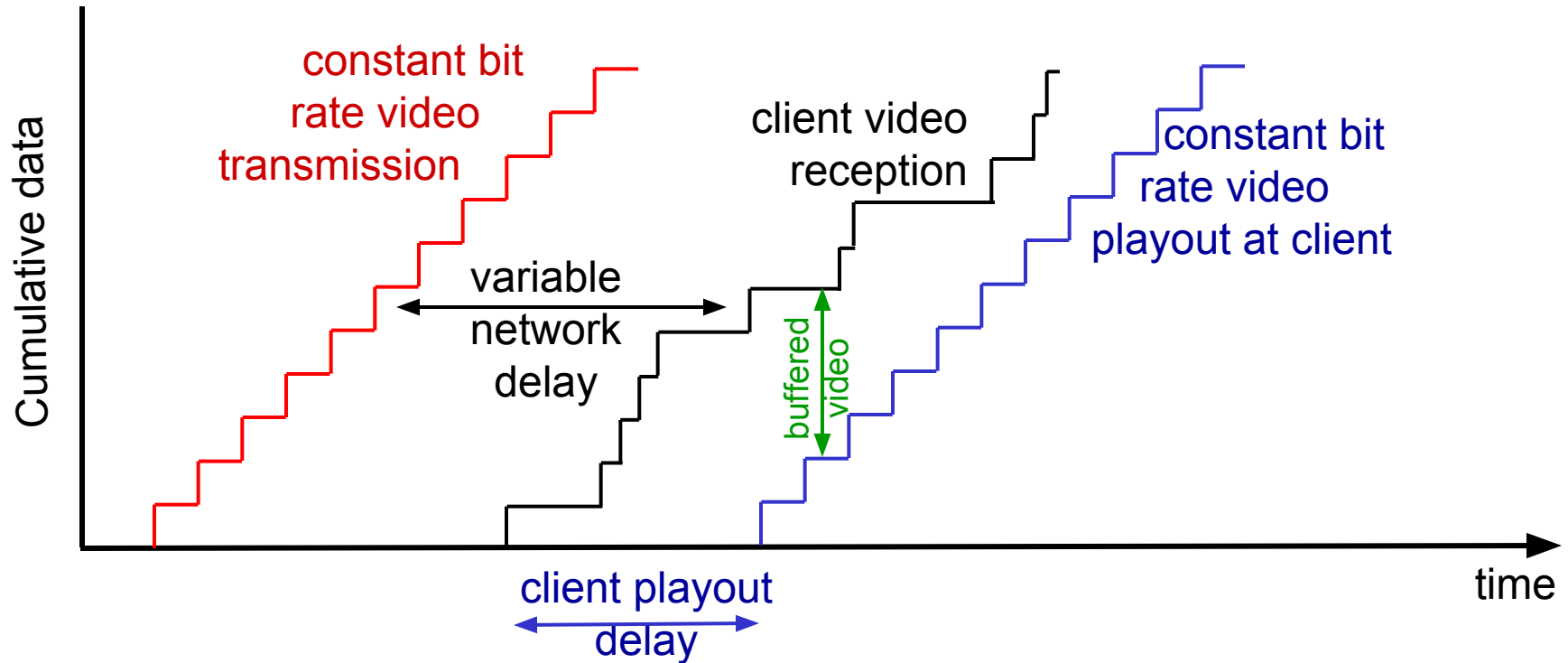
# Streaming stored video:



# Streaming stored video: challenges

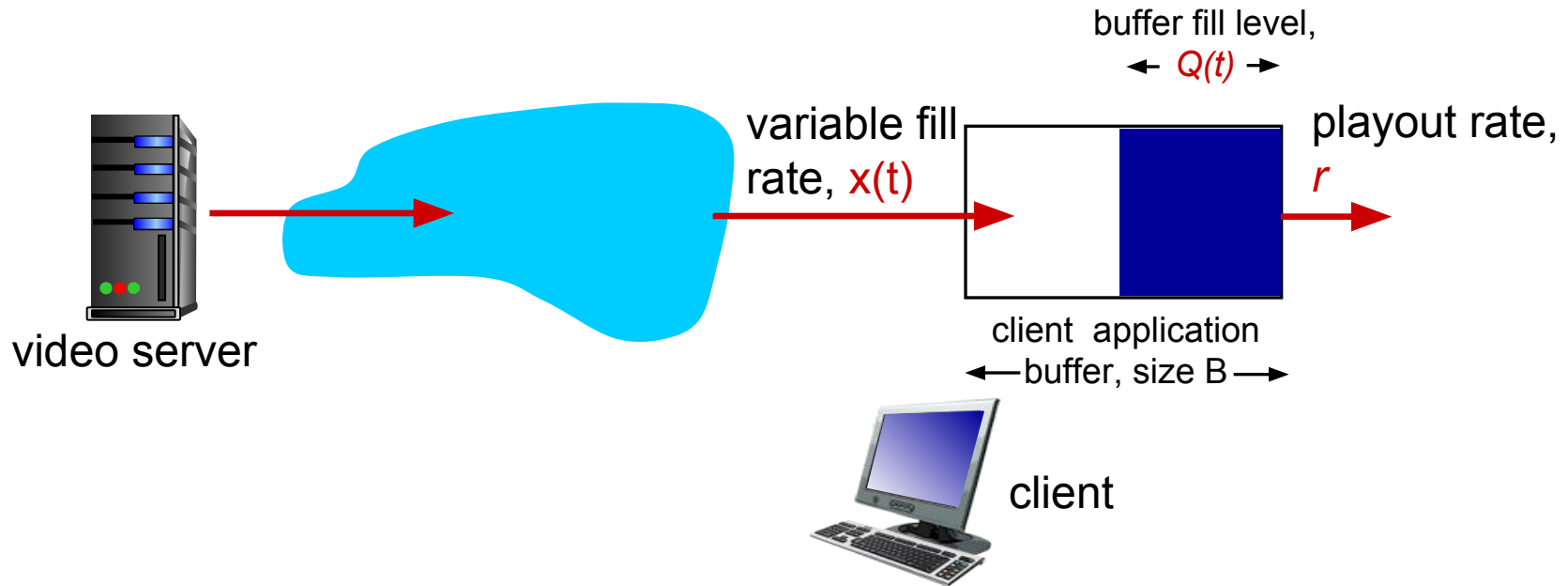
- ❖ *continuous playout constraint*: once client playout begins, playback must match original timing
  - ... but *network delays are variable* (jitter), so will need *client-side buffer* to match playout requirements
- ❖ other challenges:
  - client interactivity: pause, fast-forward, rewind, jump through video
  - video packets may be lost, retransmitted

# Streaming stored video: revisited

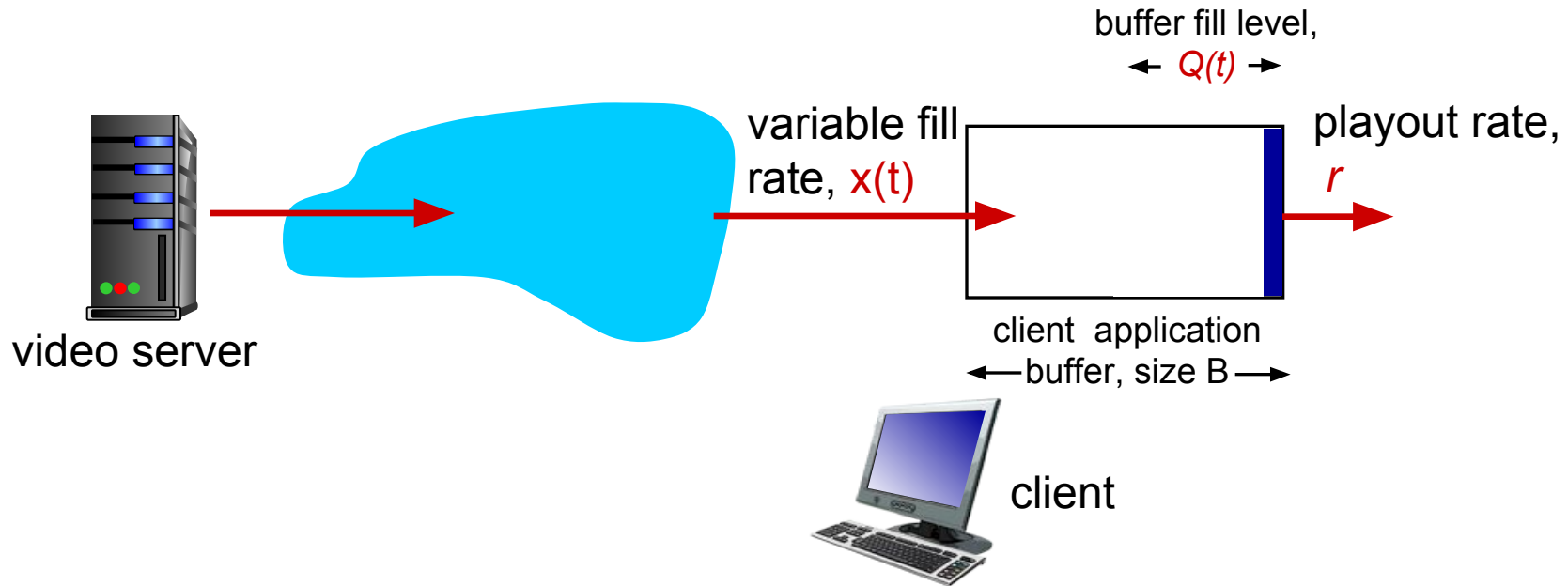


- ❖ *client-side buffering and playout delay:*  
compensate for network-added delay, delay jitter

# Client-side buffering, playout

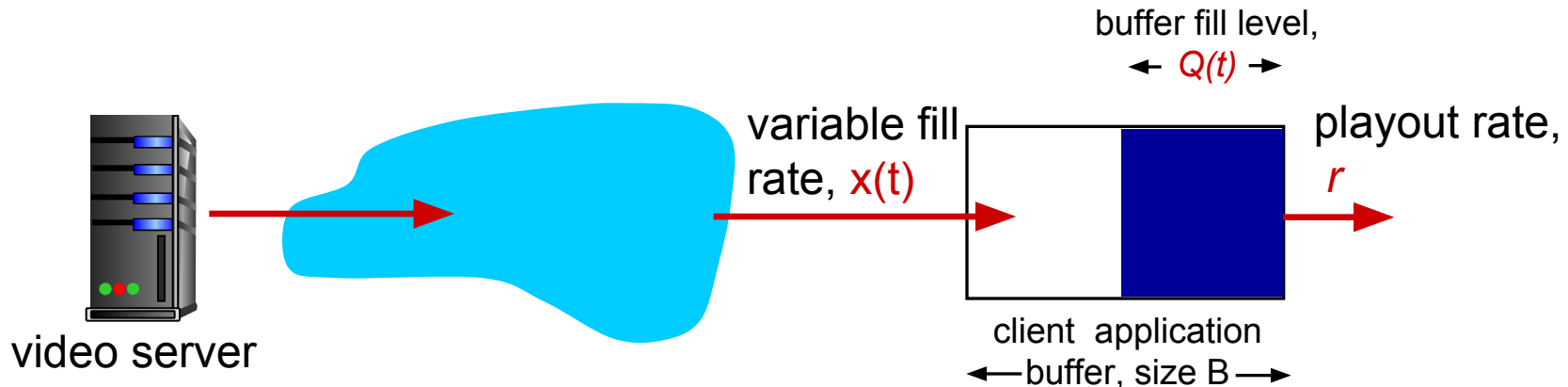


# Client-side buffering, playout



1. Initial fill of buffer until playout begins at  $t_p$
2. playout begins at  $t_p$ ,
3. buffer fill level varies over time as fill rate  $x(t)$  varies and playout rate  $r$  is constant

# Client-side buffering, playout



*playout buffering: average fill rate  $\bar{x}$ , playout rate  $r$ .*

- ❖  $\bar{x} < r$ : buffer eventually empties (causing freezing of video playout until buffer again fills)
- ❖  $\bar{x} > r$ : buffer will not empty, provided initial playout delay is large enough to absorb variability in  $x(t)$ 
  - *initial playout delay tradeoff*: buffer starvation less likely with larger delay, but larger delay until user begins watching

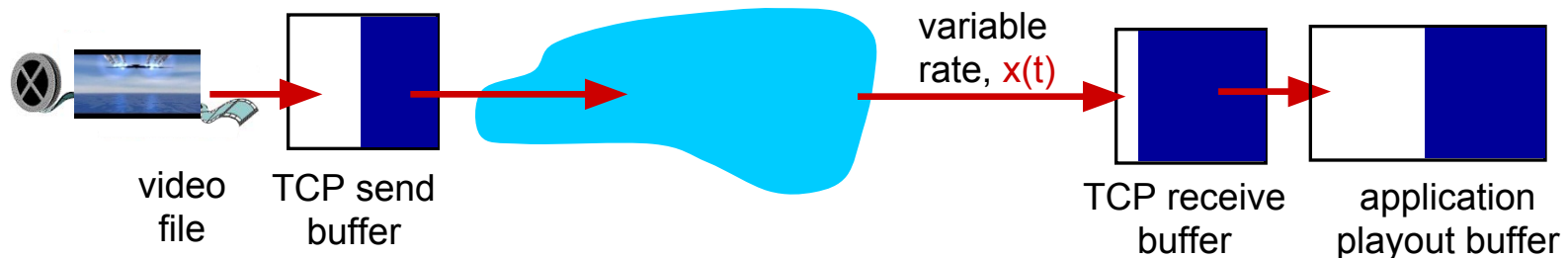
# Streaming multimedia: UDP

- ❖ server sends at rate appropriate for client
  - often: send rate = encoding rate = constant rate
  - transmission rate can be oblivious to congestion levels
- ❖ short playout delay (2-5 seconds) to remove network jitter
- ❖ error recovery: application-level, time permitting
- ❖ UDP may *not* go through firewalls



# Streaming multimedia: HTTP

- ❖ send at maximum possible rate under TCP



- ❖ fill rate *server* fluctuates due to TCP congestion control, retransmissions (in-order delivery)
- ❖ larger playout delay: smooth TCP delivery rate
- ❖ HTTP/TCP passes more easily through firewalls

# Content distribution networks

---

- ❖ *challenge*: how to stream content (selected from millions of videos) to hundreds of thousands of simultaneous users?
  
  - ❖ *option 1*: single, large “mega-server”
    - single point of failure
    - point of network congestion
    - long path to distant clients
    - multiple copies of video sent over outgoing link
- ....quite simply: this solution *doesn't scale*

# Content distribution networks

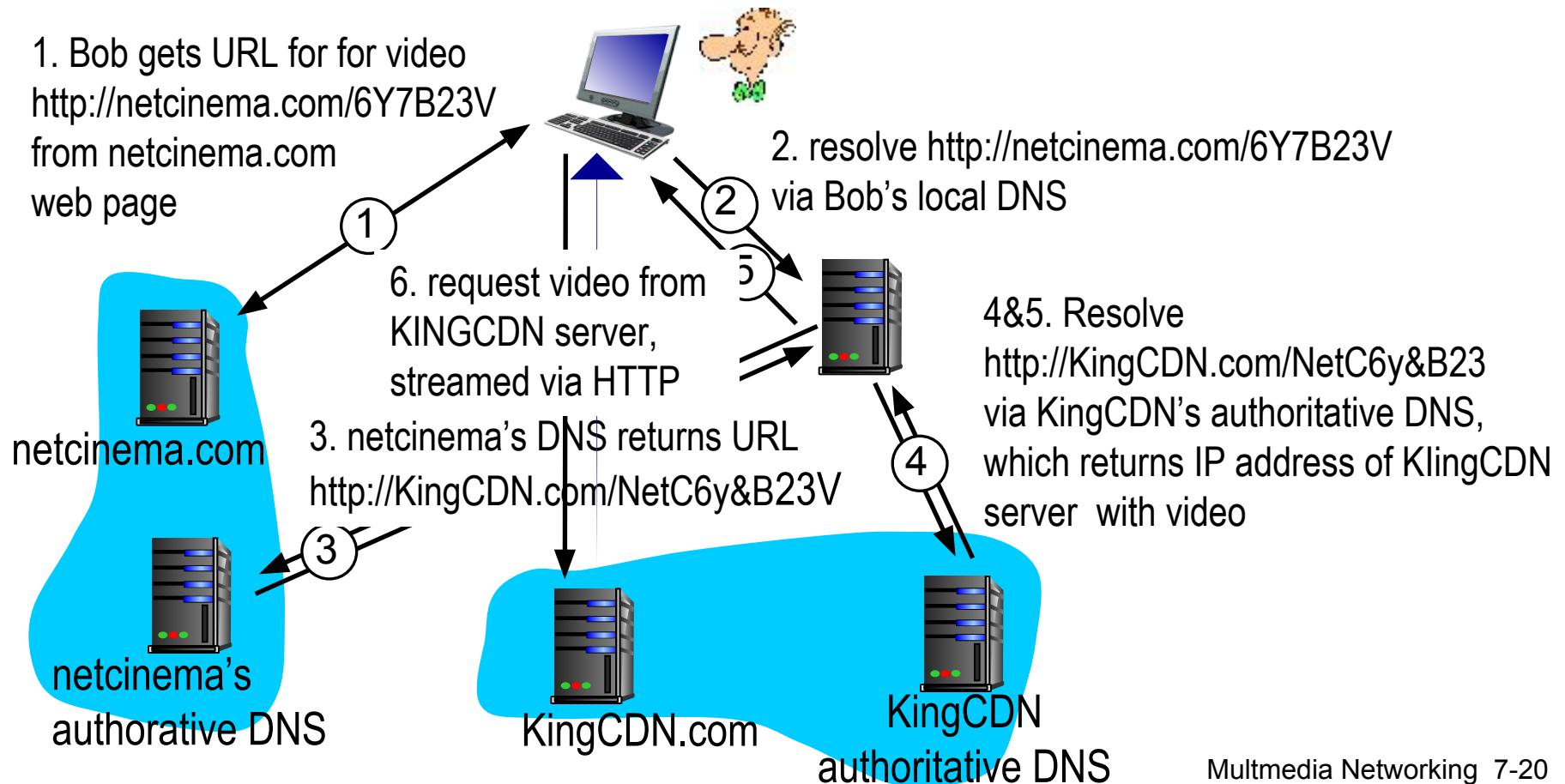
---

- ❖ *challenge:* how to stream content (selected from millions of videos) to hundreds of thousands of simultaneous users?
- ❖ *option 2:* store/serve multiple copies of videos at multiple geographically distributed sites (*CDN*)

# CDN: “simple” content access scenario

Bob (client) requests video <http://netcinema.com/6Y7B23V>

- video stored in CDN at <http://KingCDN.com/NetC6y&B23V>



# CDN cluster selection strategy

- ❖ *challenge:* how does CDN DNS select “good” CDN node to stream to client
  - pick CDN node geographically closest to client
  - pick CDN node with shortest delay (or min # hops) to client (CDN nodes periodically ping access ISPs, reporting results to CDN DNS)
  
- ❖ *alternative:* let *client* decide - give client a list of several CDN servers
  - client pings servers, picks “best”

# Case study: Netflix

- ❖ Informal homework: figure out how netflix outsources content delivery.

# Multimedia networking: outline

7.1 multimedia networking applications

7.2 streaming *stored* video

7.3 voice-over-IP

7.4 protocols for *real-time* conversational applications

7.5 network support for multimedia

# Voice-over-IP (VoIP)

- ❖ *VoIP end-end-delay requirement:* needed to maintain “conversational” aspect
  - higher delays noticeable, impair interactivity
  - < 150 msec: good
  - > 400 msec bad
  - includes application-level (packetization, playout), network delays
- ❖ *session initialization:* how does callee advertise IP address, port number, encoding algorithms?
- ❖ *value-added services:* call forwarding, screening, recording
- ❖ *emergency services:* 911



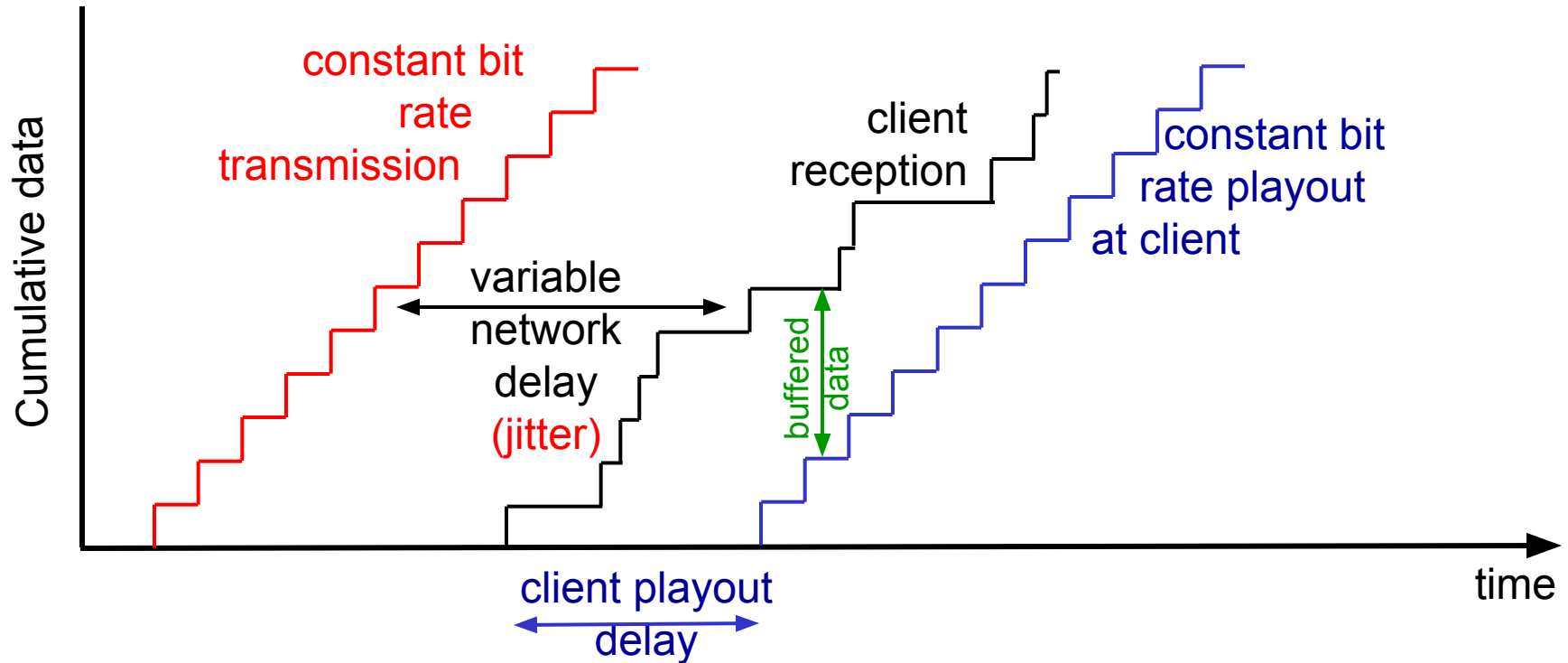
# VoIP characteristics

- ❖ speaker's audio: alternating talk spurts, silent periods.
  - 64 kbps during talk spurt
  - pkts generated only during talk spurts in chunks (app layer)
- ❖ chunk+header encapsulated into UDP or TCP segment
- ❖ application sends segment into socket every 20 msec during talkspurt

# VoIP: packet loss, delay

- ❖ *network loss*: IP datagram lost due to network congestion (router buffer overflow)
- ❖ *delay loss*: IP datagram arrives too late for playout at receiver
  - delays: processing, queueing in network; end-system (sender, receiver) delays
  - typical maximum tolerable delay: 400 ms
- ❖ *loss tolerance*: depending on voice encoding, loss concealment, packet loss rates between 1% and 10% can be tolerated

# Delay jitter



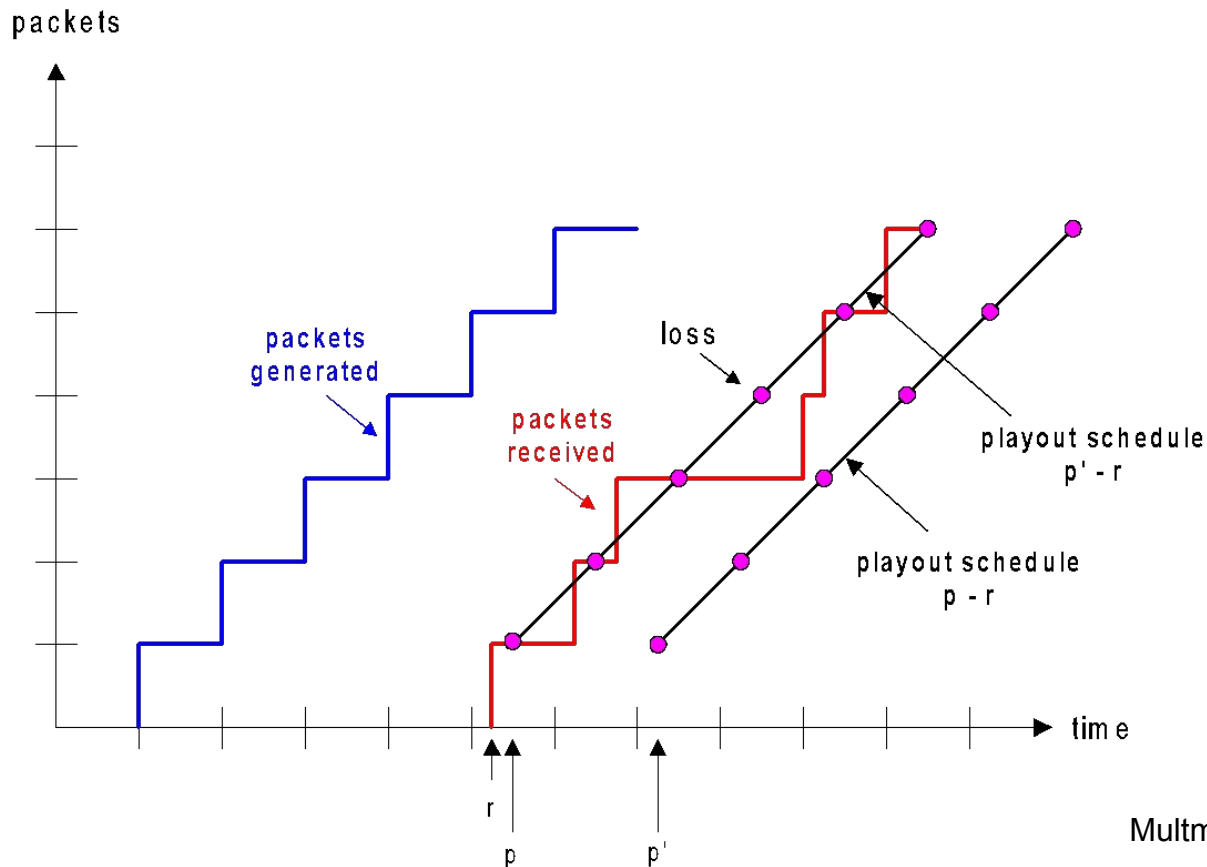
- ❖ end-to-end delays of two consecutive packets: difference can be more or less than 20 msec (transmission time difference)

# VoIP: fixed playout delay

- ❖ receiver attempts to playout each chunk exactly  $q$  msecs after chunk was generated.
  - chunk has time stamp  $t$ : play out chunk at  $t+q$
  - chunk arrives after  $t+q$ : data arrives too late for playout: data “lost”
- ❖ tradeoff in choosing  $q$ :
  - *large  $q$ : less packet loss*
  - *small  $q$ : better interactive experience*

# VoIP: fixed playout delay

- sender generates packets every 20 msec during talk spurt.
- first packet received at time  $r$
- first playout schedule: begins at  $p$
- second playout schedule: begins at  $p'$



# VoiP: recovery from packet loss (I)

*Challenge:* recover from packet loss given small tolerable delay between original transmission and playout

- ❖ each ACK/NAK takes  $\sim$  one RTT
- ❖ alternative: *Forward Error Correction (FEC)*
  - send enough bits to allow recovery without retransmission (recall two-dimensional parity in Ch. 5)

# Multimedia networking: outline

- 7.1 multimedia networking applications
- 7.2 streaming stored video
- 7.3 voice-over-IP
- 7.4 protocols for real-time conversational applications: RTP, SIP
- 7.5 network support for multimedia

# Multimedia networking: outline

7.1 multimedia networking applications

7.2 streaming stored video

7.3 voice-over-IP

7.4 protocols for real-time conversational applications

7.5 network support for multimedia



# Network support for multimedia

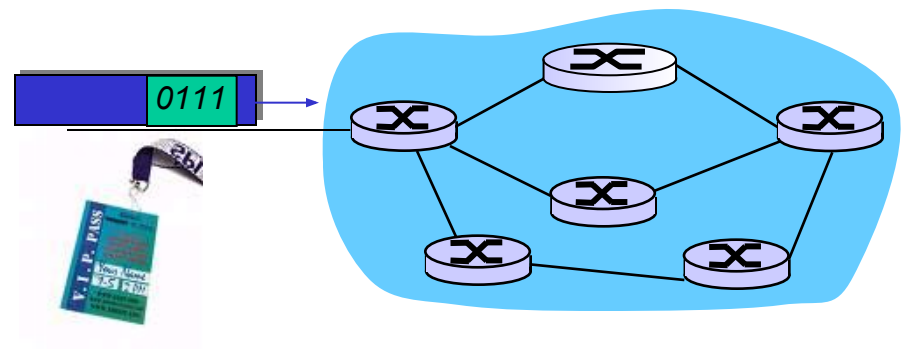
Approach	Granularity	Guarantee	Mechanisms	Complex	Deployed?
Making best of best effort service	All traffic treated equally	None or soft	No network support (all at application)	low	everywhere
Differentiated service	Traffic “class”	None or soft	Packet market, scheduling, policing.	med	some
Per-connection QoS	Per-connection flow	Soft or hard after flow admitted	Packet market, scheduling, policing, call admission	high	little to none

# Dimensioning best effort networks

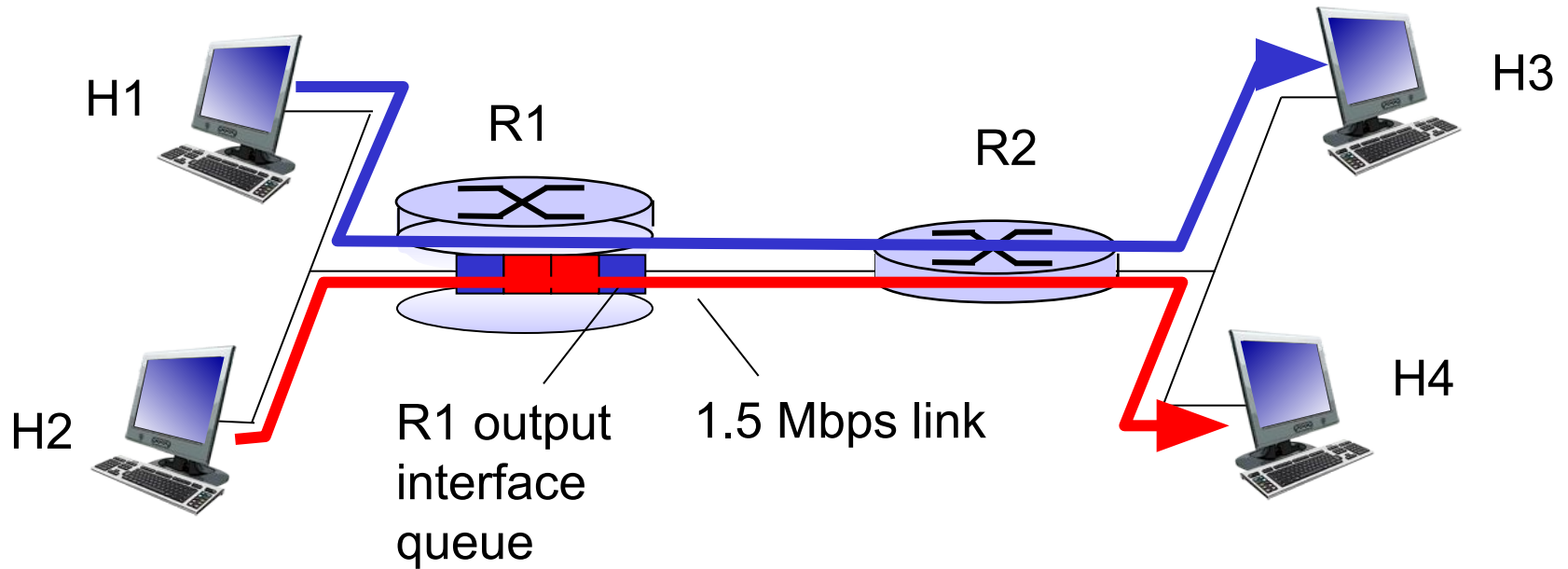
- ❖ *approach*: deploy enough link capacity so that congestion doesn't occur, multimedia traffic flows without delay or loss
  - low complexity of network mechanisms (use current “best effort” network)
  - high bandwidth costs
- ❖ challenges:
  - *network dimensioning*: how much bandwidth is “enough?”
  - *estimating network traffic demand*: needed to determine how much bandwidth is “enough” (for that much traffic)

# Providing multiple classes of service

- ❖ thus far: making the best of best effort service
  - one-size fits all service model
- ❖ alternative: multiple classes of service
  - partition traffic into classes
  - network treats different classes of traffic differently (analogy: VIP service versus regular service)
- ❖ granularity: differential service among multiple classes, *not among individual connections*
- ❖ history: ToS bits

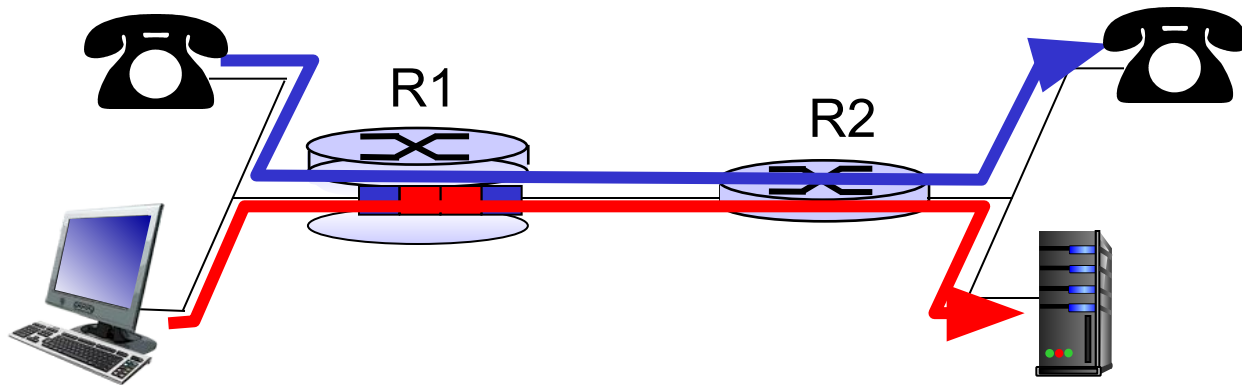


# Multiple classes of service: scenario



# Scenario 1: mixed HTTP and VoIP

- ❖ example: 1 Mbps VoIP, HTTP share 1.5 Mbps link.
  - HTTP bursts can congest router, cause audio loss
  - want to give priority to audio over HTTP

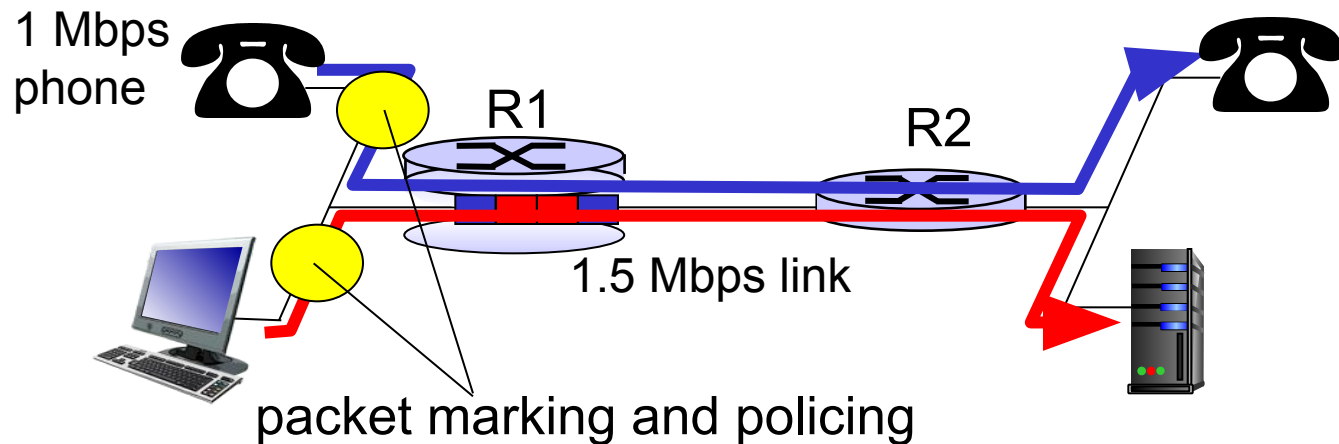


## *Principle*

packet marking needed for router to distinguish between different classes; and new router policy to treat packets accordingly

# Principles for QOS guarantees (more)

- ❖ what if applications misbehave (VoIP sends higher than declared rate)
  - policing: force source adherence to bandwidth allocations
- ❖ *marking, policing* at network edge

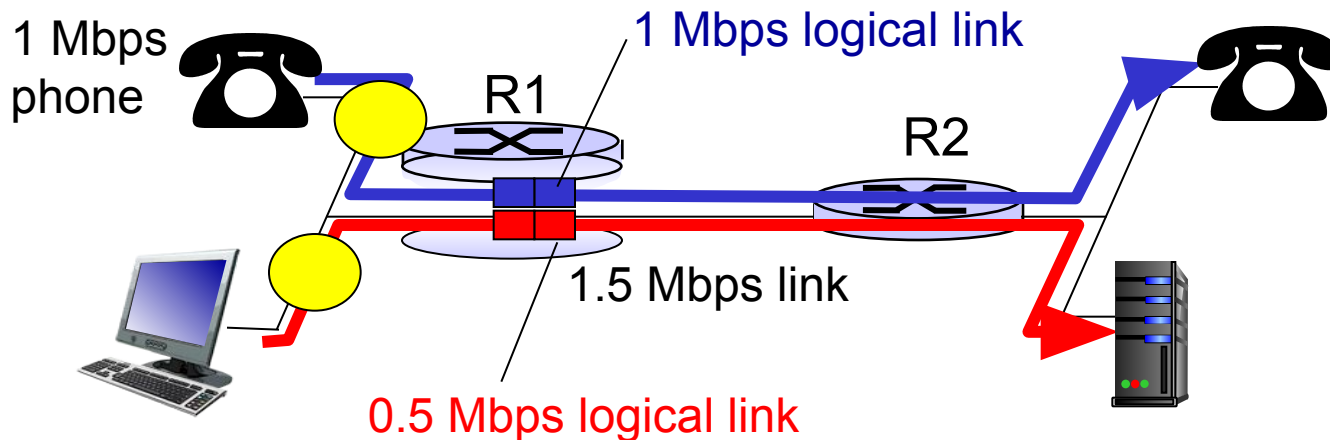


## *Principle*

**2** provide protection (isolation) for one class from others

# Principles for QOS guarantees (more)

- ❖ allocating *fixed* (non-sharable) bandwidth to flow: *inefficient* use of bandwidth if flows doesn't use its allocation

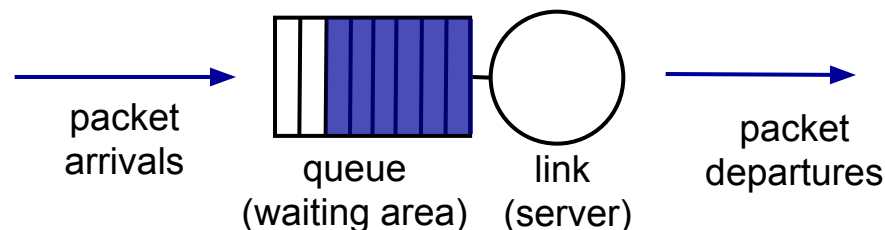


## Principle

while providing isolation, it is desirable to use resources as efficiently as possible

# Scheduling and policing mechanisms

- ❖ *scheduling*: choose next packet to send on link
- ❖ *FIFO (first in first out) scheduling*: send in order of arrival to queue
  - real-world example?
  - *discard policy*: if packet arrives to full queue: who to discard?
    - *tail drop*: drop arriving packet
    - *priority*: drop/remove on priority basis
    - *random*: drop/remove randomly



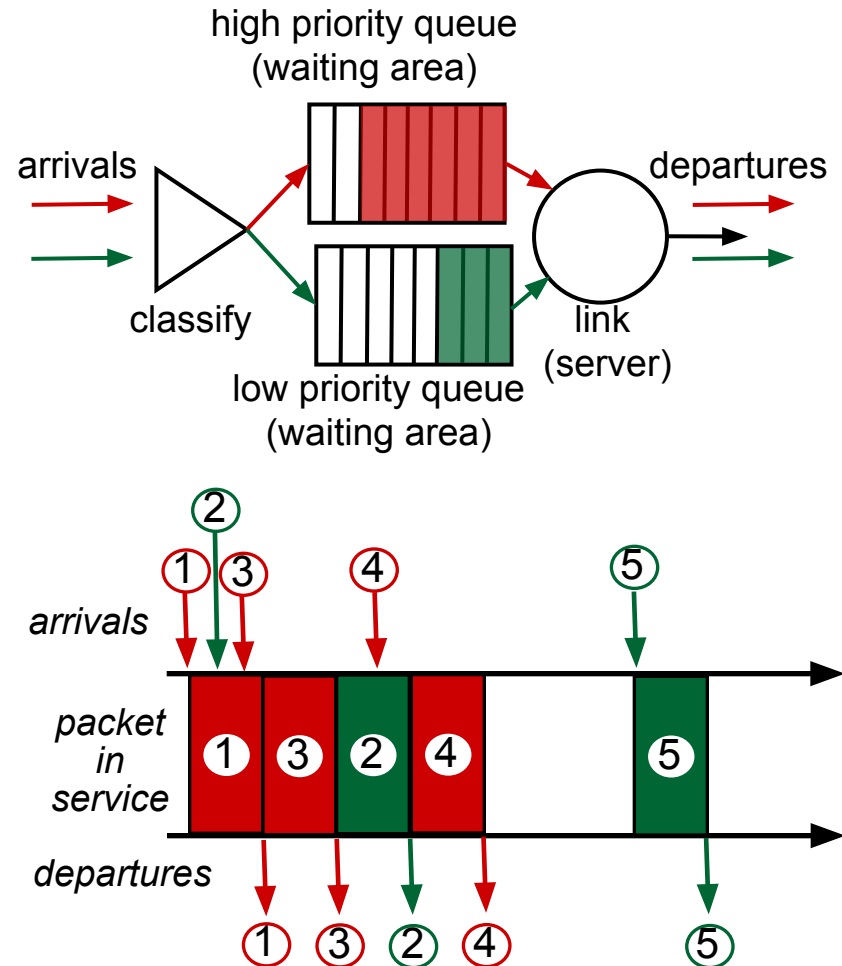


# Scheduling policies: priority

*priority scheduling*: send highest priority queued packet

❖ multiple *classes*, with different priorities

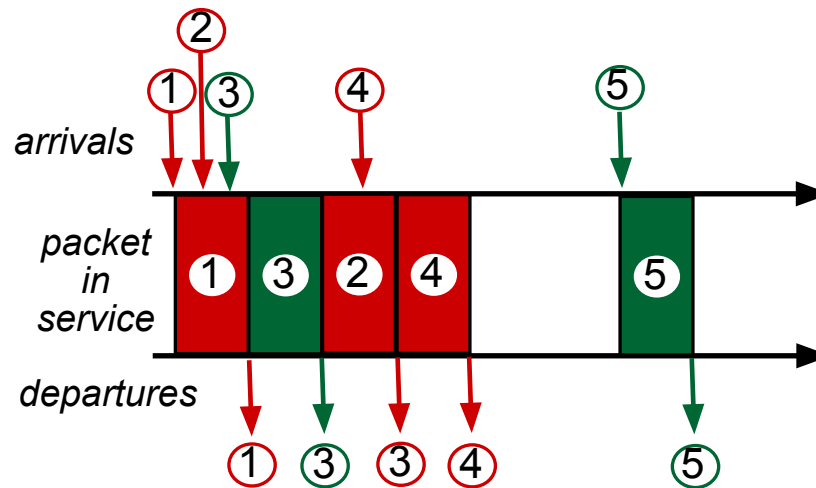
- class may depend on marking or other header info, e.g. IP source/dest, port numbers, etc



# Scheduling policies: still more

## *Round Robin (RR) scheduling:*

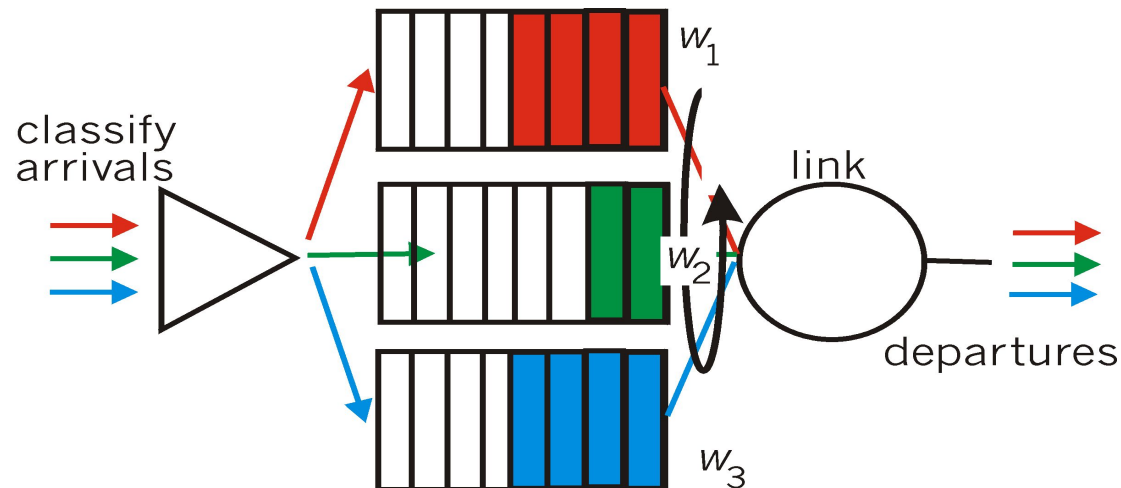
- ❖ multiple classes
- ❖ cyclically scan class queues, sending one complete packet from each class (if available)



# Scheduling policies: still more

## *Weighted Fair Queuing (WFQ):*

- ❖ generalized Round Robin
- ❖ each class gets weighted amount of service in each cycle



# Policing mechanisms

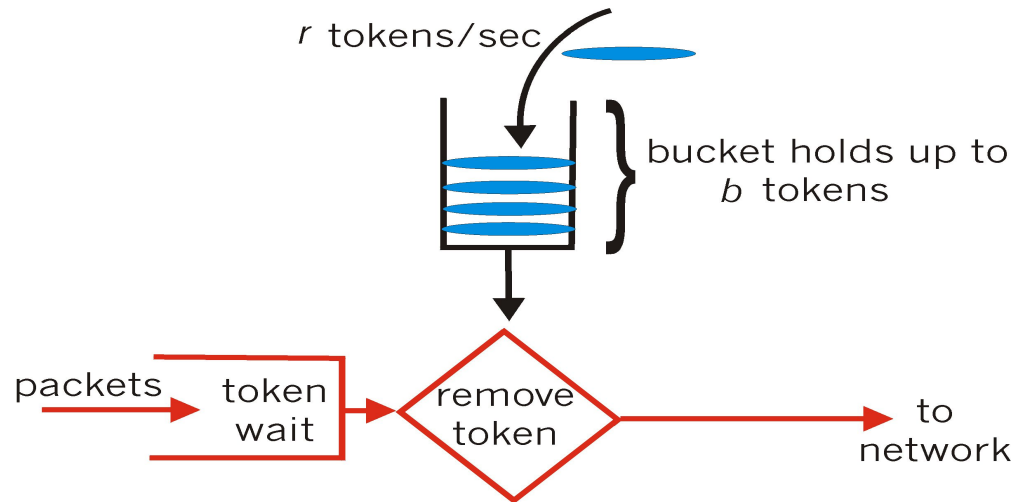
*goal:* limit traffic to not exceed declared parameters

Three common-used criteria:

- ❖ *(long term) average rate:* how many pkts can be sent per unit time (in the long run)
  - crucial question: what is the interval length: 100 packets per sec or 6000 packets per min have same average!
- ❖ *peak rate:* e.g., 6000 pkts per min (ppm) avg.; 1500 ppm peak rate
- ❖ *(max.) burst size:* max number of pkts sent consecutively (with no intervening idle)

# Policing mechanisms: implementation

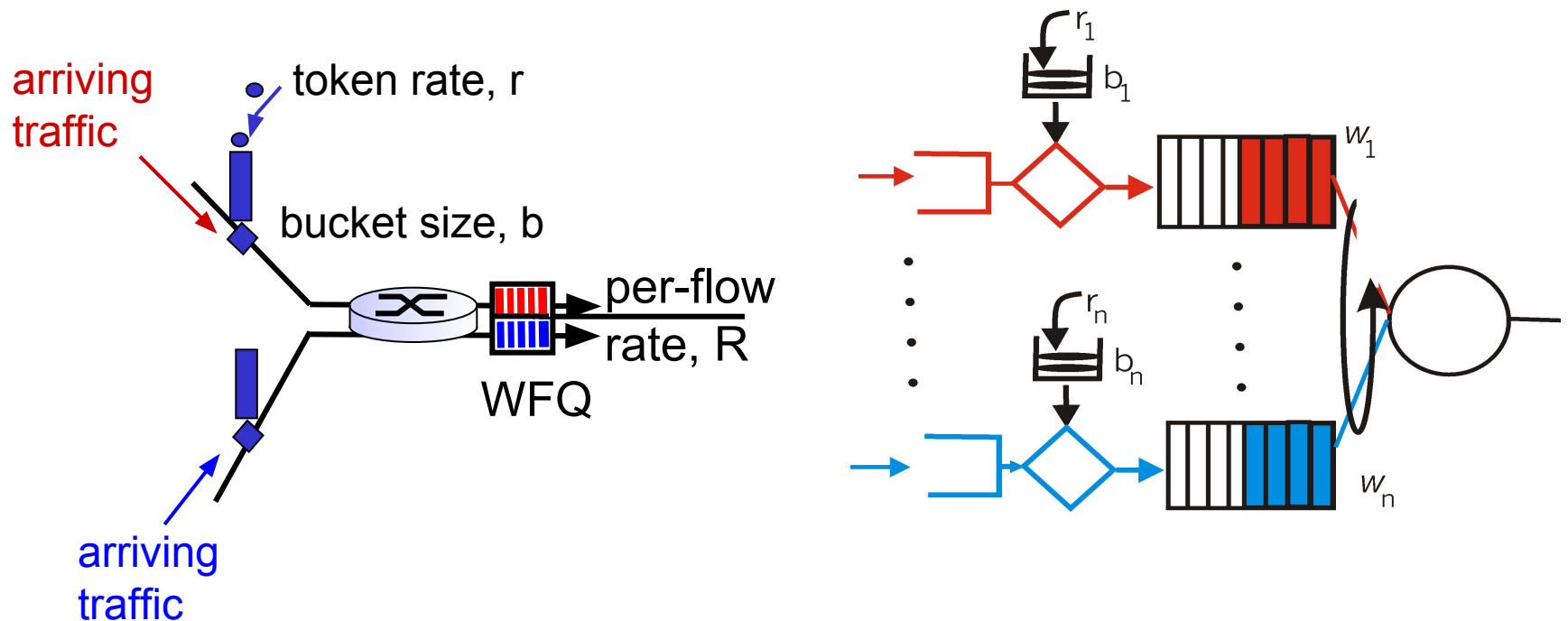
*token bucket*: limit input to specified *burst size* and *average rate*



- ❖ bucket can hold  $b$  tokens
- ❖ tokens generated at rate  $r$  token/sec unless bucket full
- ❖ *over interval of length  $t$ : number of packets admitted less than or equal to  $(r t + b)$*

# Policing and QoS guarantees

- ❖ token bucket, WFQ combine to provide guaranteed upper bound on delay, i.e., *QoS guarantee!*

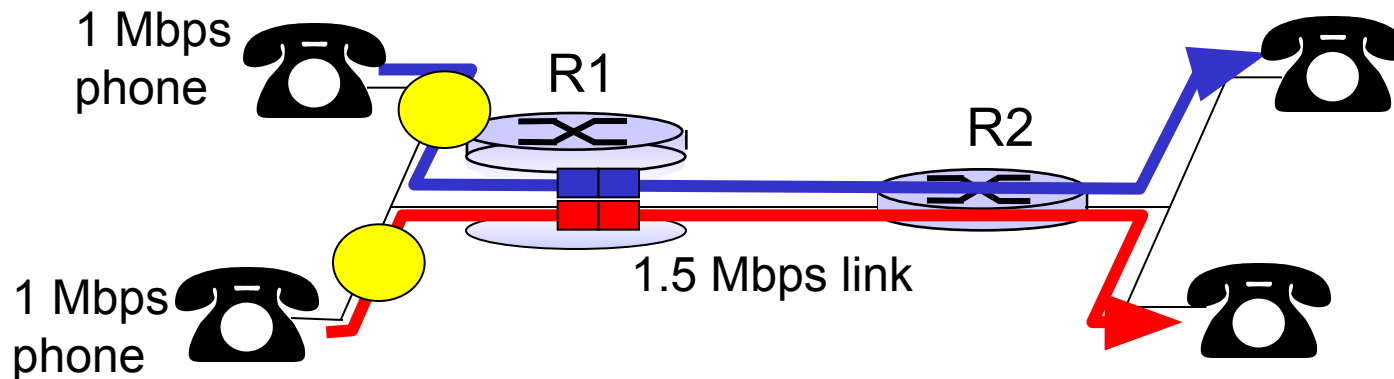


# Differentiated services

- ❖ want “qualitative” service classes
  - “behaves like a wire”
  - relative service distinction: Platinum, Gold, Silver
- ❖ *scalability*: simple functions in network core, relatively complex functions at edge routers (or hosts)
  - signaling, maintaining per-flow router state difficult with large number of flows

# Per-connection QOS guarantees

- ❖ *basic fact of life*: can not support traffic demands beyond link capacity



## *Principle*

*call admission*: flow declares its needs, network may

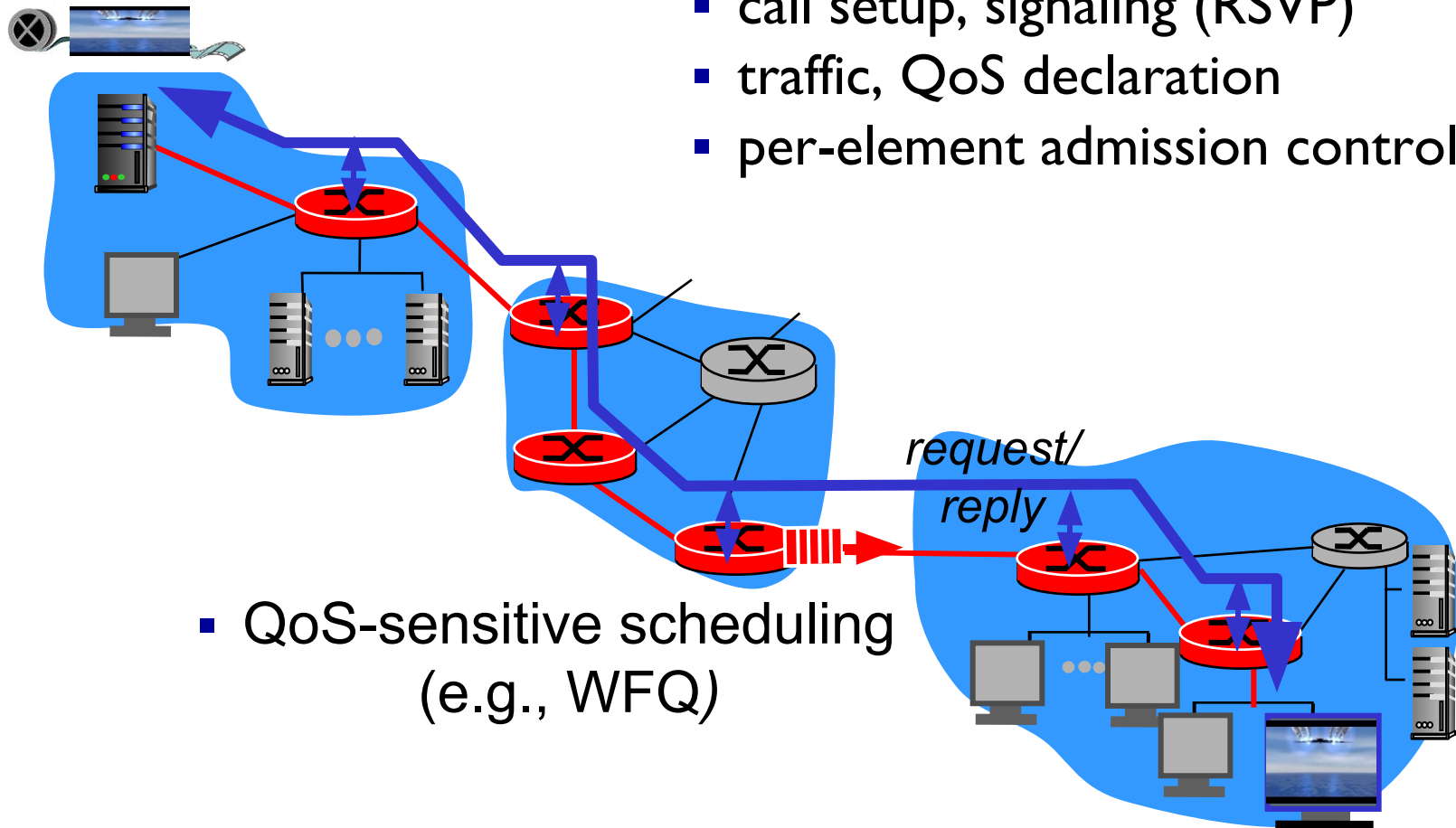
block call (e.g., busy signal) if it cannot meet needs



# QoS guarantee scenario

## ❖ *resource reservation*

- call setup, signaling (RSVP)
- traffic, QoS declaration
- per-element admission control



- QoS-sensitive scheduling (e.g., WFQ)