# A Survey of Scheduling Methods

**Raj Jain**

Raj Jain is now at
Washington University in Saint Louis
Jain@cse.wustl.edu
http://www.cse.wustl.edu/~jain/

# Overview

q Goals

q Metrics

q Classification

q Scheduling Methods

# Scheduling: Goals

q Sharing bandwidth

q Sharing bandwidth fairly

q Meeting bandwidth guarantees (min and max)
$\Rightarrow$ Provide isolation between users

q Meeting loss guarantees (multiple levels)

q Meeting delay guarantees (multiple levels)

q Reducing delay variation

# Goals (Cont)

q  Guarantees require call admission control,
    policing, shaping, drop policies,
    buffer allocation, and scheduling

q  These issues are, therefore, related.

q  For example, zero-loss can be obtained by
    allocating PCR (but no multiplexing gain).

# Scheduling: Methods

q FCFS

q Round Robin

q Priority Queueing

q Priority Queueing with Windows

q Generalized Processor Sharing (GPS)

q VirtualClock

q Weighted Fair Queueing (WFQ), WF2Q, WF2Q+

q Self-Clocked Fair Queueing (SCFQ)

q Stop and Go

q Rate Controlled Service Descipline (RCSD)

# Scheduling Metrics

q Complexity of enqueue+dequeue processes

q Fairness: If two flows are backlogged, difference between their weighted throughputs is bounded

q Complexity of adding and releasing connections (or changing quotas in ABR)

q Delay bounds should not depend upon behavior of other flows, number of other flows, reservations of other flows
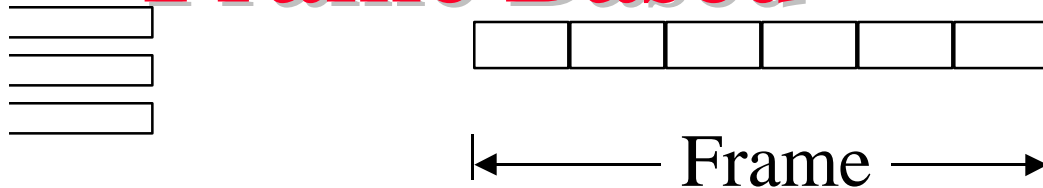
# Scheduling Classification

q Work conserving vs nonconserving

q Sorted priority vs frame based

q Control vs accomodate distortion

# Work Conserving vs Nonconserving

q Conserving: Server not idle if there is work.

   q Produces lower average delay, higher delay var.

   q Produces high total throughput

   q Examples: GPS, WFQ,  VirtualClock

q Nonconserving: Better for multiple hops

   q May produce lower worst case end-to-end delay

   q May produce higher network throughput

   q Reshaping at every hop $\Rightarrow$ additive hop delays

   q Examples: Stop & Go

# Sorted Priority vs Frame Based



Frame

- q  Sorted Priority: Virtual time for each flow/packet
  - q  Generally has a O(log(V)) complexity
  - q  Examples: VirtualClock, WFQ
- q  Frame Based: Time split into fixed/variable frames
  - q  Each flow reserves the time per frame
  - q  Delay and bandwidth allocations are dependent
    - q  Examples: Stop and Go uses constant frame size. DRR, WRR allow variable frame size
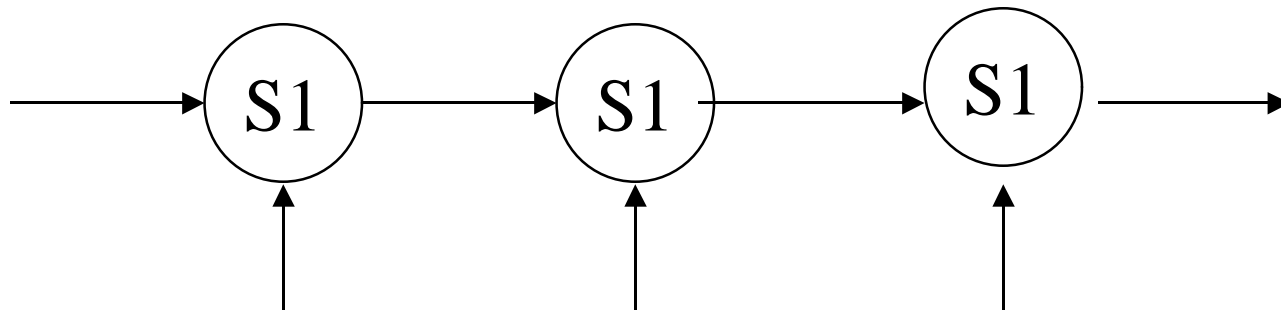
# Control Distortion vs Accomodate distortion

q Burstiness of the traffic increases along the path $\Rightarrow$ Need more resources

q Control Distortion: Reshape at each hop $\Rightarrow$ Non-work conserving

  q Example: Stop and Go, HRR, Jitter EDD

q Accomodate Distortion: Do not reshape

  q Example: VirtualClock, Fair Queueing, GPS, Delay EDD

# FCFS

q Unfair

q No isolation among users

# Round Robin

```
    ──►(S1)──►(S1)──►(S1)──►
        ▲      ▲      ▲
        │      │      │
```

q  Parking lot problem: Distance sources get lower

q  Classify incoming traffic into flows (Src-Dest pairs)

q  Round-robin among flows

q  Known Problems:
   Ignores packet length $\Rightarrow$ Fair Queueing

q  Ref: Nagle

# Priority Queueing

q  Also known as head of line (HOL)

q  Priority 0 through n-1

q  Priority 0 is always serviced first.

q  Priority i is serviced only if 0 through i-1 are empty

q  Highest priority has the lowest delay, highest throughput, lowest loss

q  Lower priority classes may be starved if higher priority are overloaded

# Priority Queueing with Windows

q Maximum $n_i$ packets from ith priority during a single round

q Come back to higher priority unless $n_i$ packets have been served

q Guarantees non-starvation but increases the delay for higher priorities

q Large $n_i \Rightarrow$ Priority queueing.
Small $n_i \Rightarrow$ Round robin

# Priority with Windows (Cont)

q $n_i$'s determine min bandwidth allocations and delays

q Quantitative relationships between $n_i$ and delays or loss not provided.

q VLSI design implemented

q Refs: El-Gebaly et al and Sabaa et al.

# VirtualClock

- q Goals: Provide average reserved throughput $R_i$ b/s
- q Provides isolation between users
- q Upon packet arrival:
  - q $VirtualClock_i = Max\{$wall clock time, $VirtualClock_i\}$
  - q Timestamp the packet with $VirtualClock_i$
  - q $VirtualClock_i = VirtualClock_i + $ packet size$/R_i$
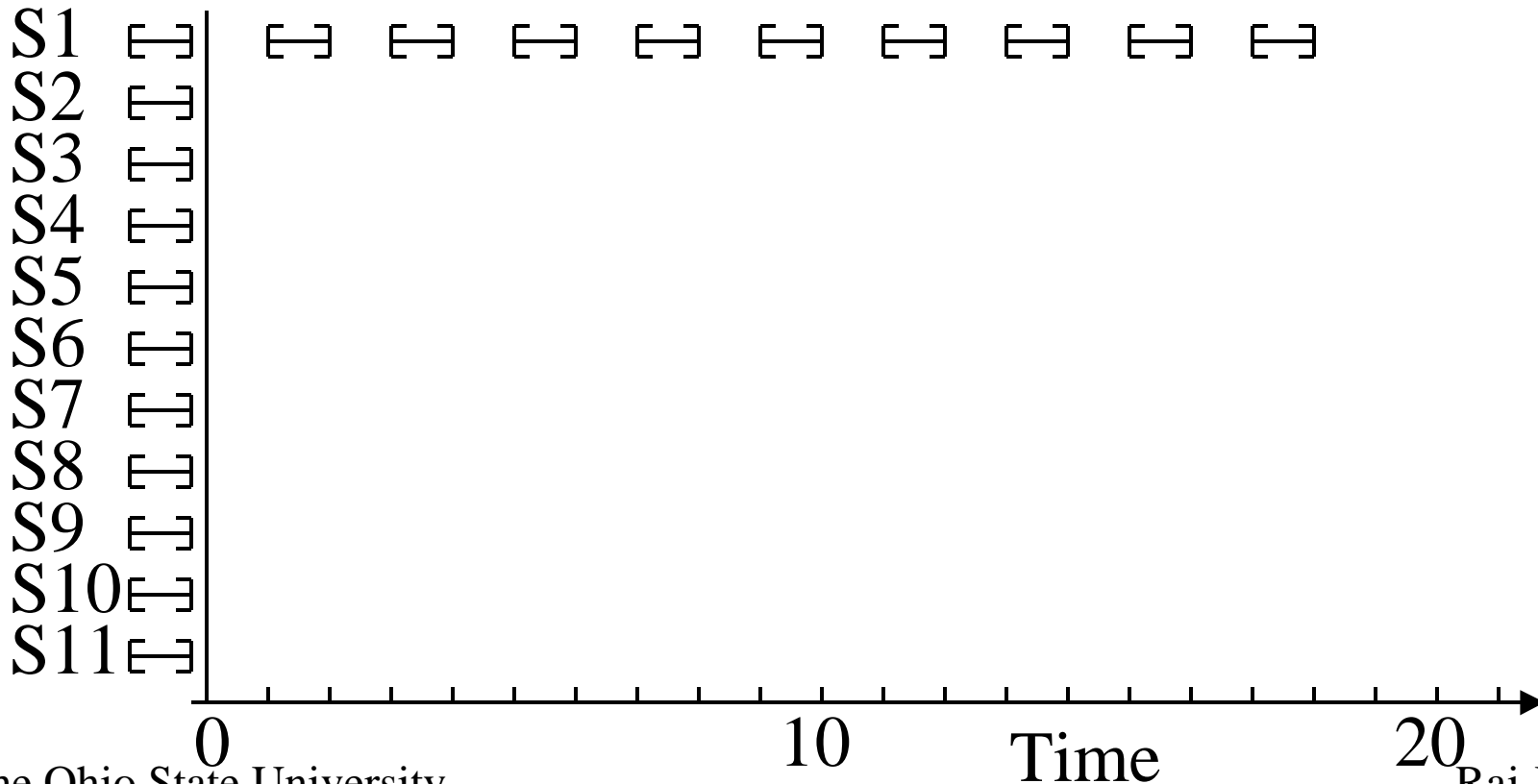- q Transmit packets in order of increasing timestamps

# VirtualClock (Cont)

q Possible to implement with one timestamp per flow rather than one per packet

q Known Problems: Flows do not accumulate credits

  q Flows using idle bandwidth are penalized later $\Rightarrow$ Virtual Clock is Unfair $\Rightarrow$ Several proposed fixes, e.g., Time-shift scheduling

  q No CAC policy $\Rightarrow$ no delay bounds

  q Need to implement priority queues $\Rightarrow$ O(log V) complexity, V=# of VCs

q Refs: Zhang, Srinivasan et al, Stilidias and Varma, Cobb et al

# Generalized Processor Sharing

q Idealized policy to split bandwidth

q Each user has a fraction si of the bandwidth

q All unused bandwidth is allocated in proportion to the fraction $\phi_i$

q At time t, ith active user gets a fraction ri

$$r_i = \phi_i / \Sigma_{\text{active } j} \phi_j$$

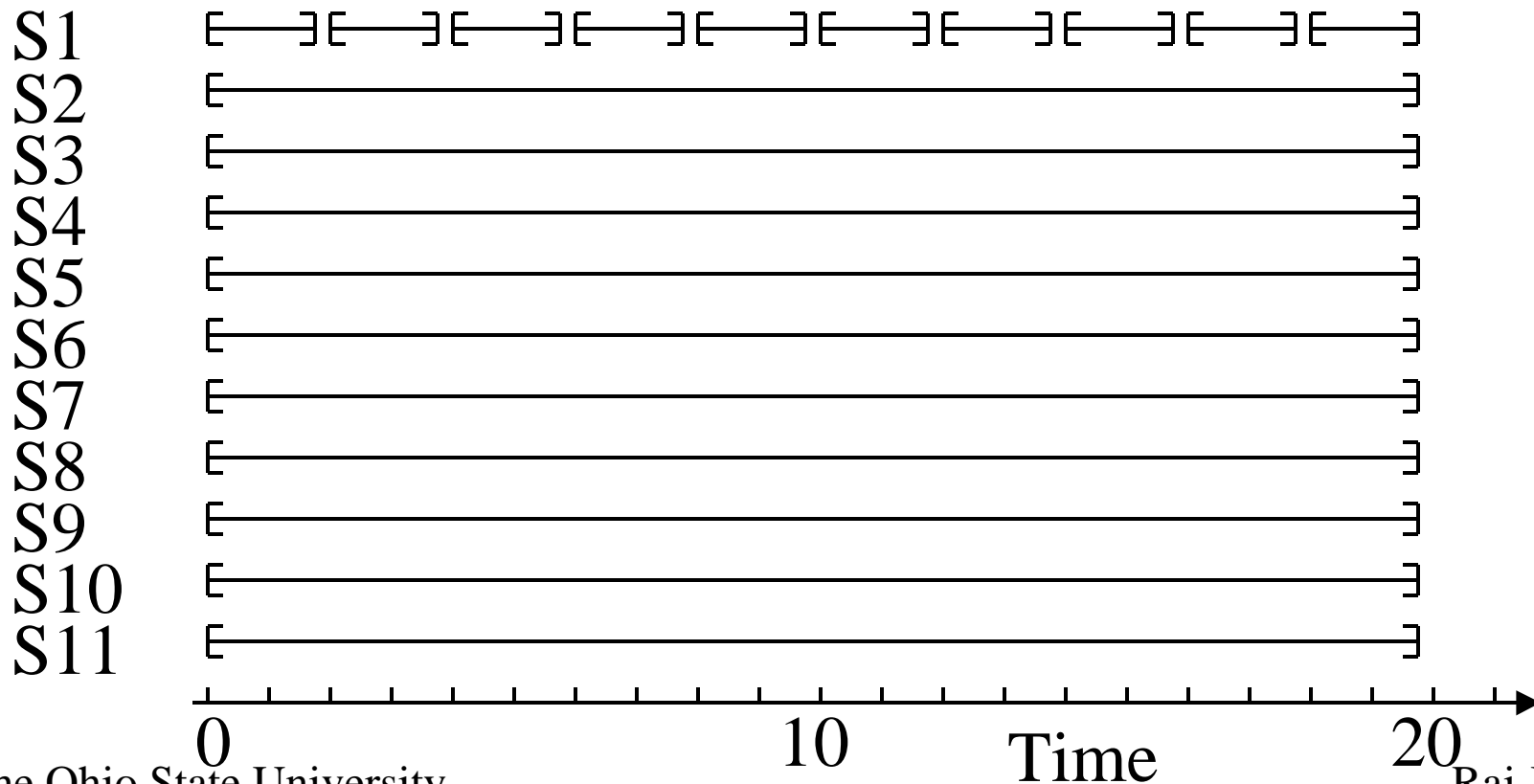q Weighted round-robin with infinitely small service quantum

# GPS Example: Arrivals

q Eleven Sources. First source gets 0.5. Other 10 sources get 0.05 each. First source sends 10 cells. 2-11 send one each at t=0.
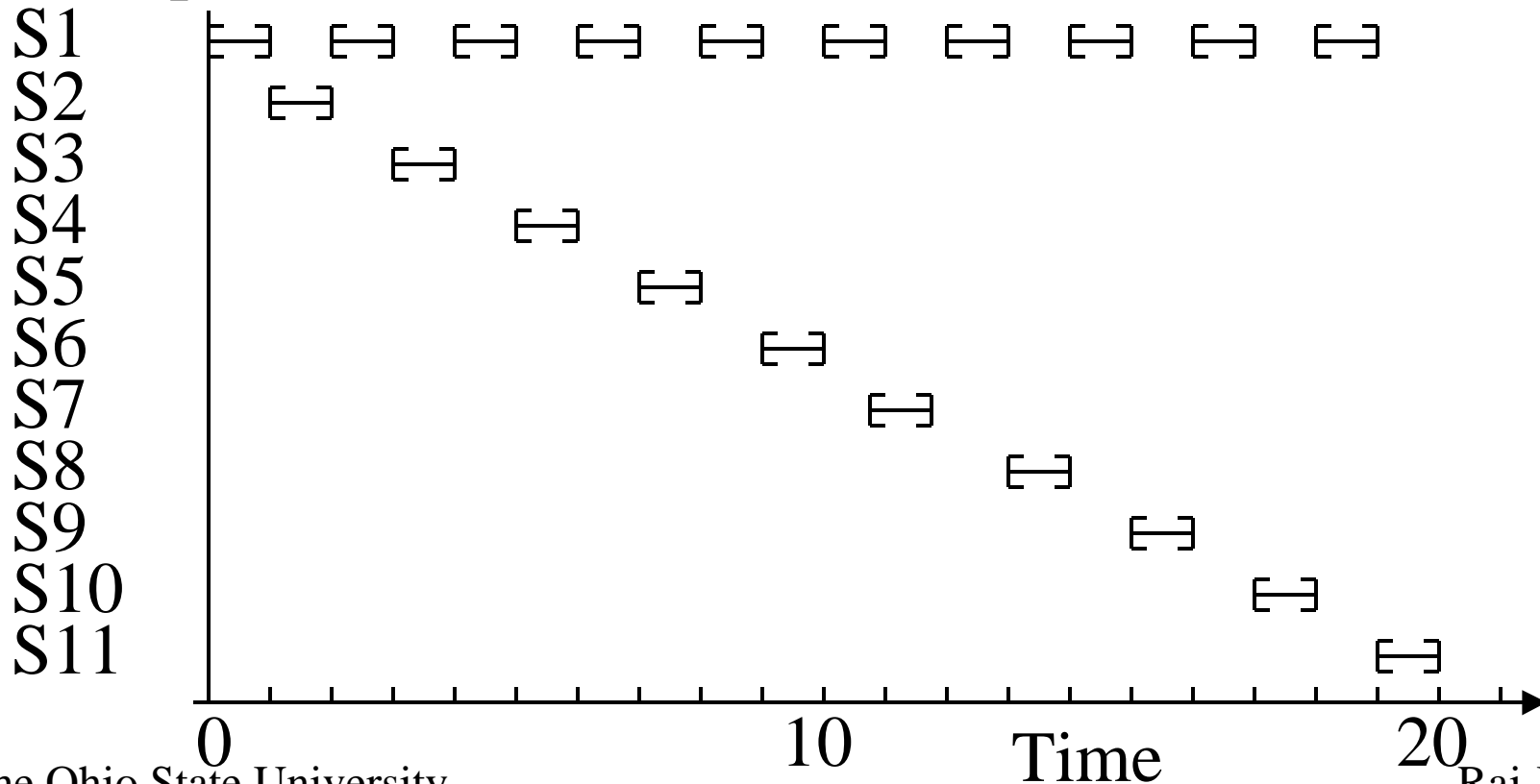
# GPS Example: Service

q  Each cell of the first source takes 2 units of time. Sources 2-11 take 20 units each.

Raj Jain

# Weighted Fair Queueing

q  Approximates bit-by-bit round robin.
   Compute GPS finish time and schedule
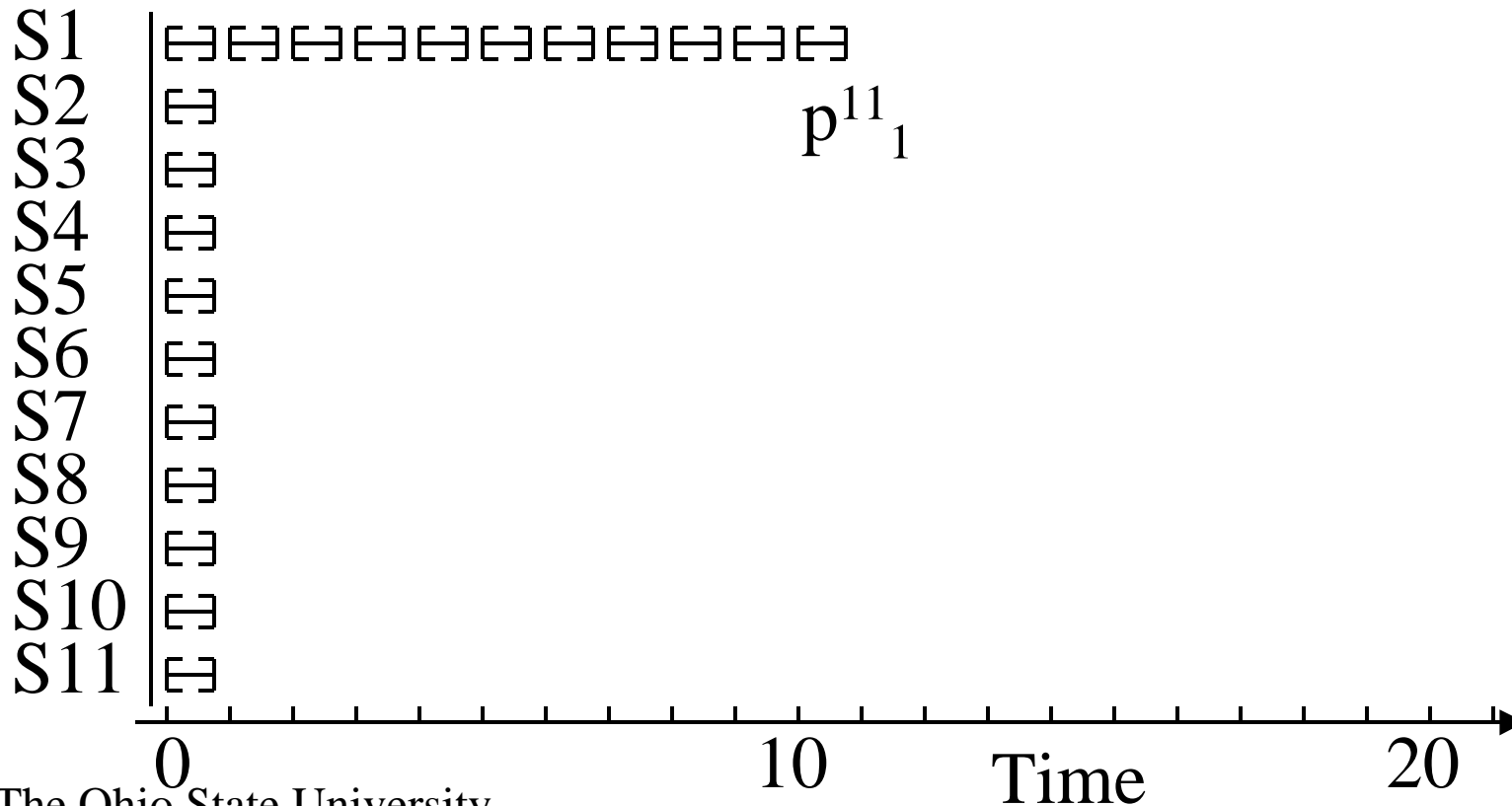   the packet with the smallest finish time.

# Weighted Fair Queueing (WFQ)

q Basis of IETF's integrated services

q Naive implementation requires O(log(m)), m=# of packets

q Keshav's implementation requires O(log(V)), V=# of flows

q **Known Properties**: CAC and End-to-end delay bounds have been derived for leaky-bucket shaped sources

q Parekh and Gallagher showed that leaky bucket +FQ $\Rightarrow$ delay guarantees

# WFQ (Cont)

q **Known Problems**:

 q Need large bandwidth reservation to get small delay bound.

 q Complex to implement.

 q Packets can be serviced much earlier than GPS. Can introduce significant unfairness over GPS.

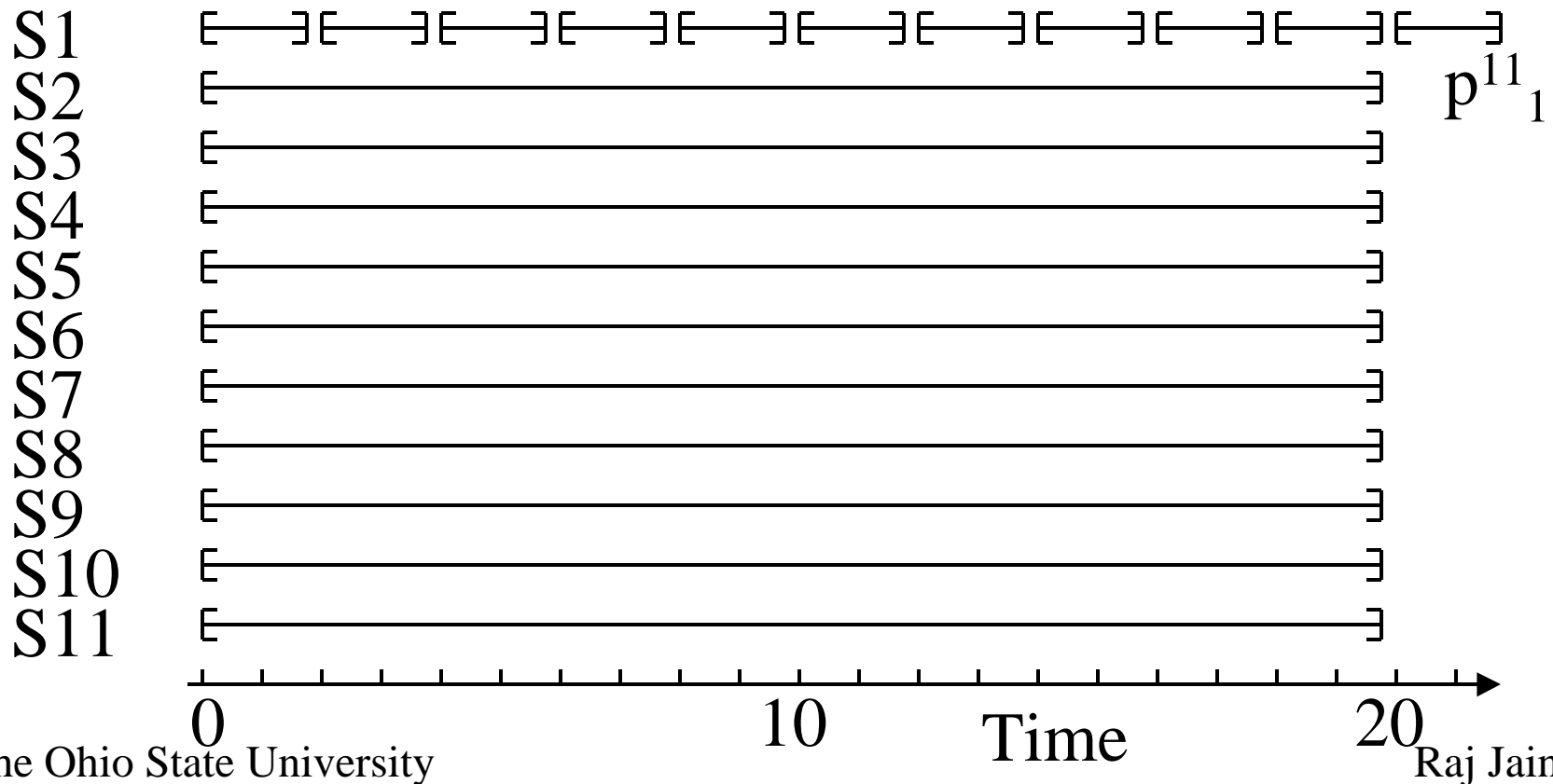q **Refs**: Demers et al, Keshav, Srinivasan et al

# GPS Example 2: Arrivals

q  Eleven Sources. First source gets 0.5.
   Other 10 sources get 0.05 each. First
   source sends 11 cells. 2-11 send one each at t=0.
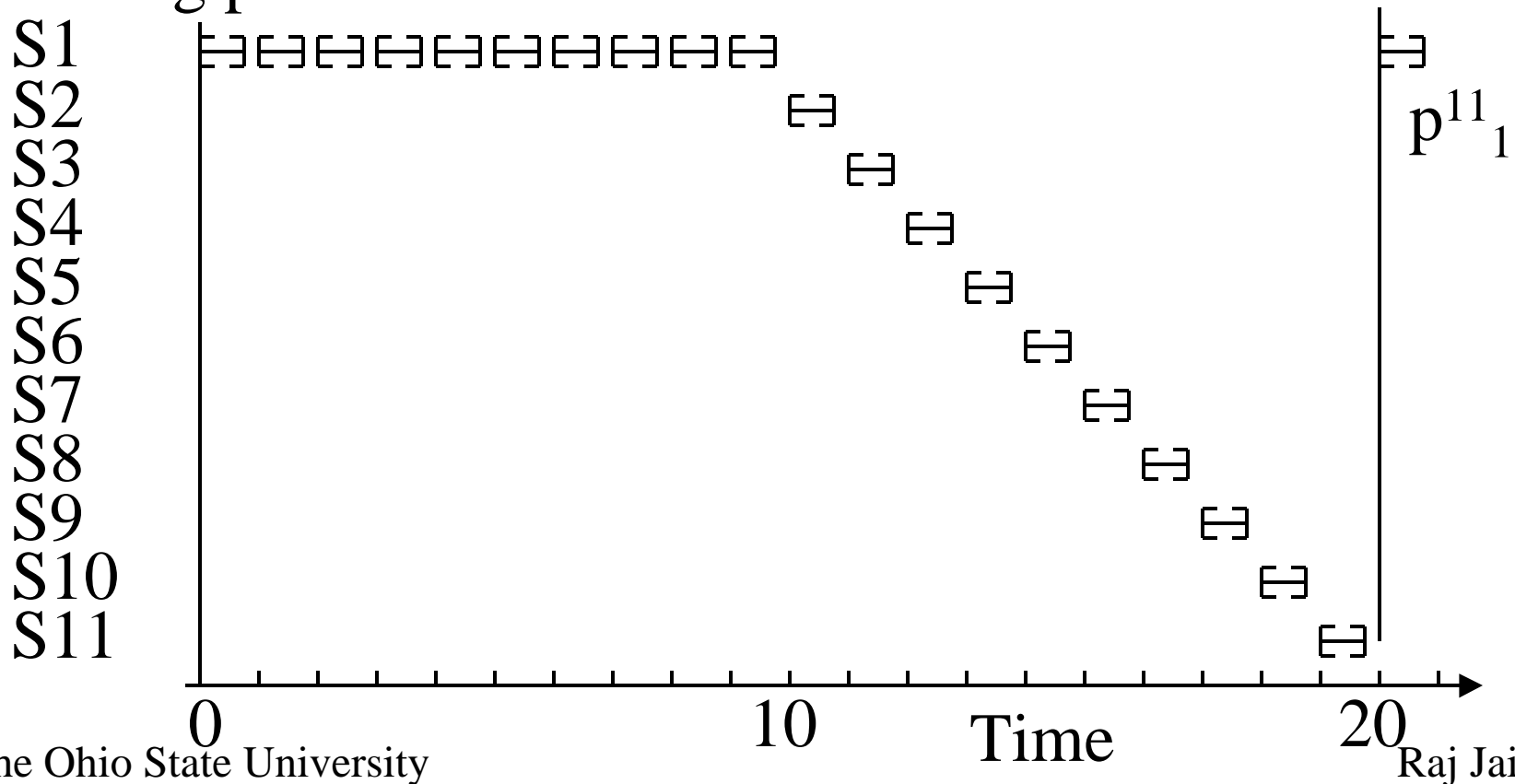
S1, S2, S3, S4, S5, S6, S7, S8, S9, S10, S11

$p^{11}_1$

Time: 0, 10, 20

# GPS Example 2: Service

q Each cell of the first source takes 2 units of time. Sources 2-11 take 20 units each.
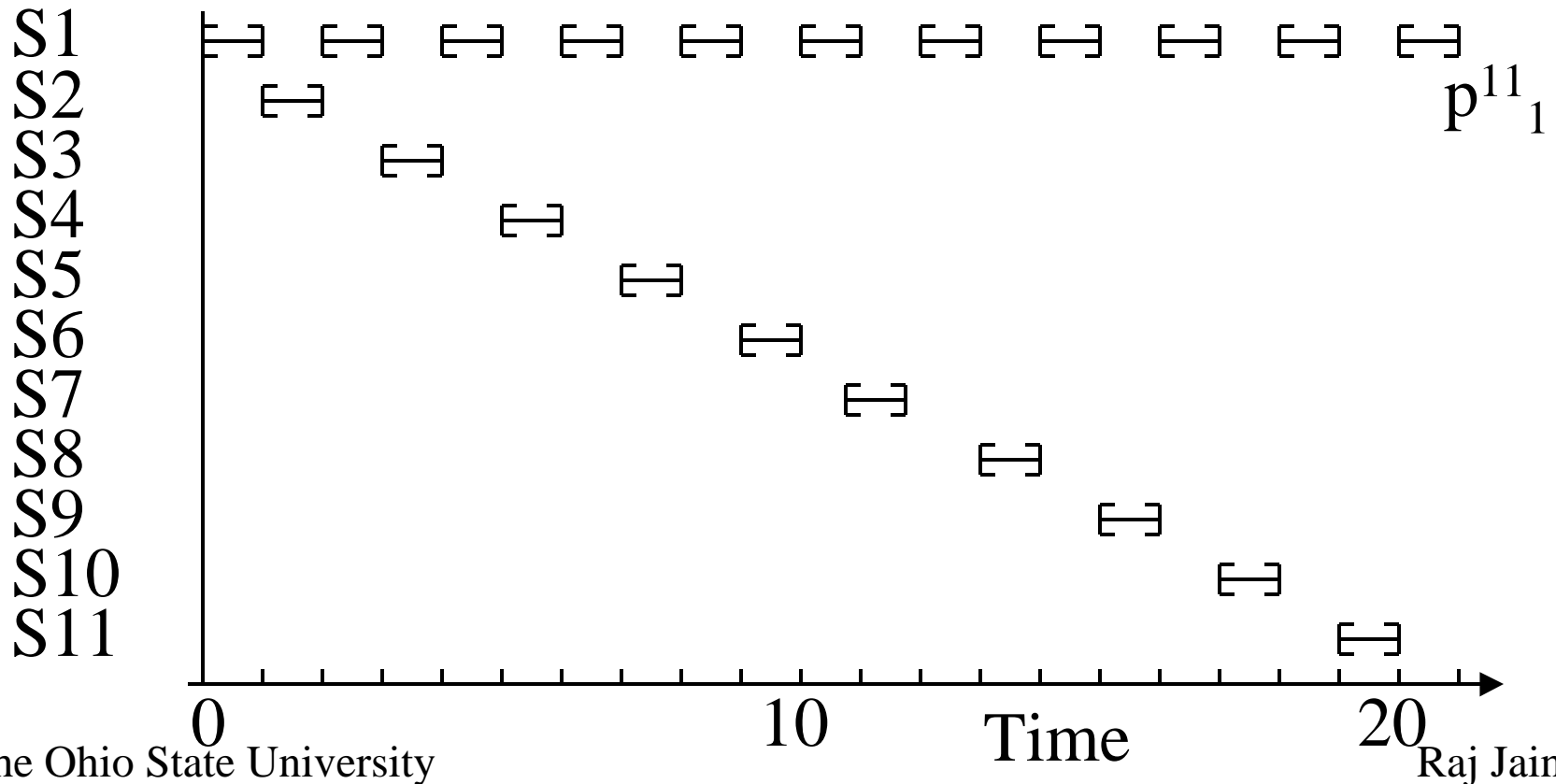


$p^{11}{}_1$

Time

Raj Jain

# WFQ: Service

q Packets finish at the same time or earlier than GPS. Some packets finish much earlier. Long period of no service $\Rightarrow$ Unfair.

S1
S2
S3
S4
S5
S6
S7
S8
S9
S10
S11

$p^{11}_1$

0          10          20
                Time

# Worst Case Fair Weighted Fair Queuing (WF2Q )

q WF2Q fixes the unfairness problem in WFQ.

  q WFQ: Among packets <u>waiting</u> in the system, pick one that will finish service first under GPS.

  q WF2Q: Among packets waiting in the system that have <u>started</u> service under GPS, select one that will finish first under GPS.

q WF2Q provides service close to GPS (difference in packet service time bounded by max. packet size).

q WF2Q+ is an simpler implementation of WF2Q

q Refs: Jon Bennett, Hui Zhang.

# WF2Q: Service

S1
S2
S3
S4
S5
S6
S7
S8
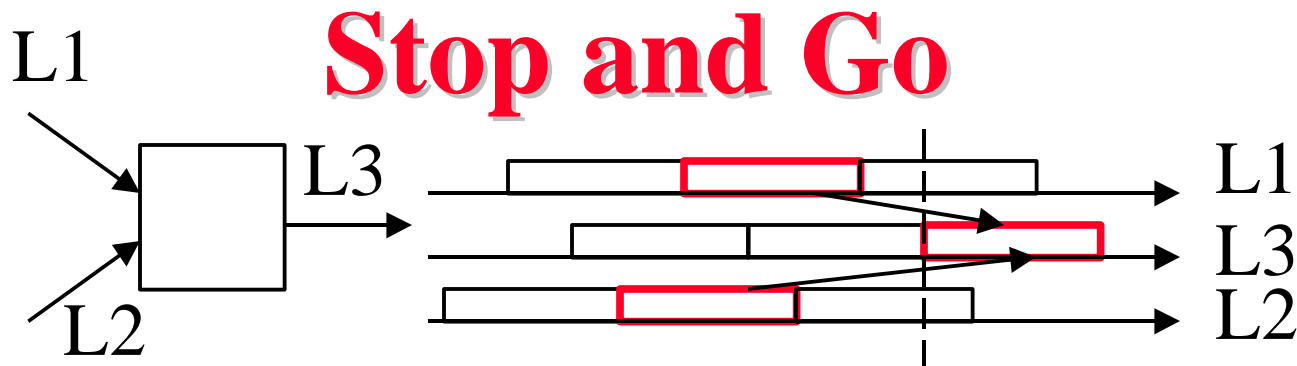S9
S10
S11

$p^{11}{}_1$

0          10          20

Time

# Self-Clocked Fair Queueing (SCFQ)

q Computational complexity of computing virtual finishing time in WFQ depends upon the number of times flows change from idle to busy and vice versa. SCFQ reduces it to $O(1)$. Dequeue and enqueue is still $O(\log(V))$

q Uses system clock instead of wall clock (as in Virtual Clock)

q A packet's tag = Length/rate + Max{tag of previous packet in that flow, tag of packet in service at the time of arrival}
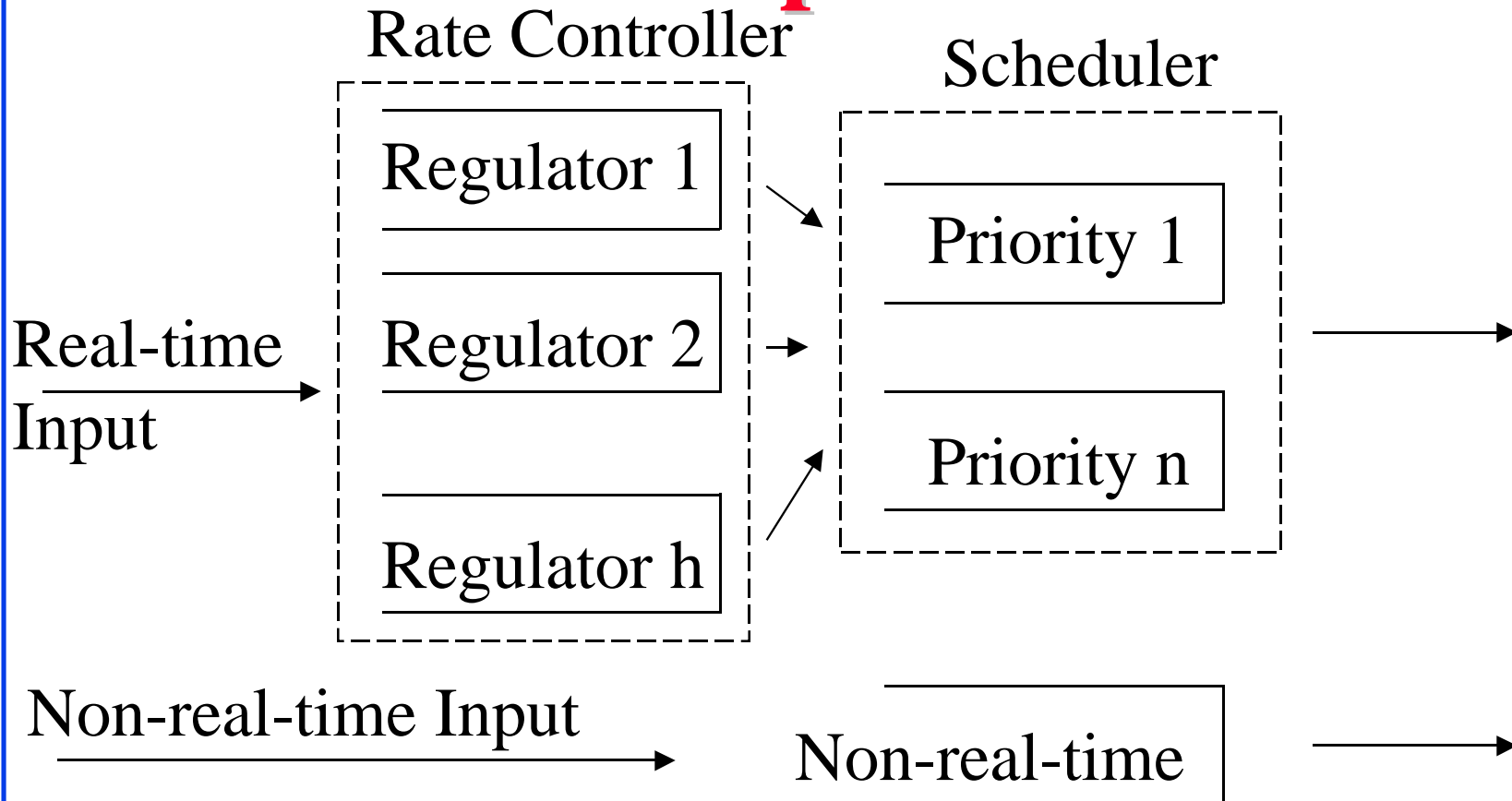
q Ref: Golestani

# Stop and Go

L1

L3

L1
L3
L2

L2

q Time is divided into constant size frames

q An arriving packet leaves on the next departing frame

q Server cannot idle if there are eligible packets

q Known Properties: Shaping maintained throughout $\Rightarrow$ Delay bound possible $\Rightarrow$ Tight delay jitter

q Non-Work conserving. Link idle if no packets for current frame but there are the next frame

q Known Problems: Allocation = PCR $\Rightarrow$ Inefficient

# Rate Controlled Service Disciplines

Rate Controller

Scheduler

```
Real-time    ┌ ─ ─ ─ ─ ─ ─ ─ ┐      ┌ ─ ─ ─ ─ ─ ─ ─ ┐
Input        │  [Regulator 1] │      │  [Priority 1]  │
─────────────│  [Regulator 2] │  →   │               │  ────────→
             │       ⋮        │      │  [Priority n]  │
             │  [Regulator h] │      └ ─ ─ ─ ─ ─ ─ ─ ┘
             └ ─ ─ ─ ─ ─ ─ ─ ┘
```

Non-real-time Input

[Non-real-time]  ────────→

q  Rate-Controller (shaper) + Packet Scheduler

Decouples bandwidth allocation and delay bounds

# RCSD (Cont)

q Rate Controller:

    q Each packet is assigned an eligibility time on arrival

    q Packet is held in rate controller and released to scheduler at its eligibility time

q Many different schedulers and rate controllers can be combined to produce different algorithms.
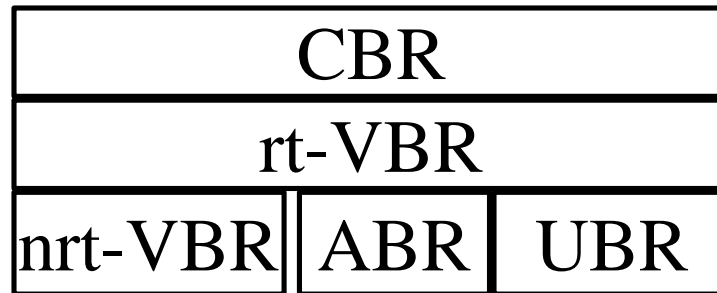
Rate-Jitter          Earliest due date

                    $\times$     Static priority

Delay-Jitter             FCFS

# Multi-class Scheduling For ATM

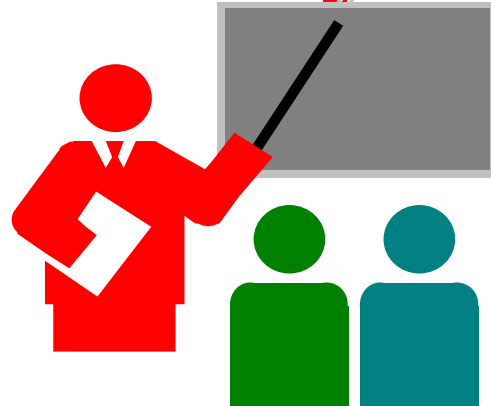| CBR | | |
|-----|-----|-----|
| rt-VBR | | |
| nrt-VBR | ABR | UBR |

q Each class has an allocation = Guaranteed under overload

q Some classes need minimum delay $\Rightarrow$ have priority.

q Some classes are greedy: They will send more than allocated and will want to use all left-over. No left-over capacity.

q Left-over capacity must be fairly allocated.

# Scheduling Methods: Comparison

| Algorithm | Unfairness | Complexity |
|---|---|---|
| Round Robin | ∞ | O(1) |
| Fair Queueing | Max | O(log(V)) |
| SCFQ | 2Max | O(log(V)) |
| DRR | 3 Max | O(1) |
| Virtual Clock | ∞ | O(log(V)) |
| WRR | Max | O(1) |

# Summary

q Schedulers can provide bandwidth, delay, loss guarantees

q Large # of VCs $\Rightarrow$ Need O(1) complexity

q Frequent rate changes $\Rightarrow$ Allocation to a VC should depend upon that VC's demands and not on others

q Non-work conserving schedulers provide end-to-end delay guarantees.

# References

q   For detailed list of references, see
    http://www.cis.ohio-state.edu/
    ~jain/refs/ref_schd.htm

# Thank You!