# Adaptive Real-Time Resource Management
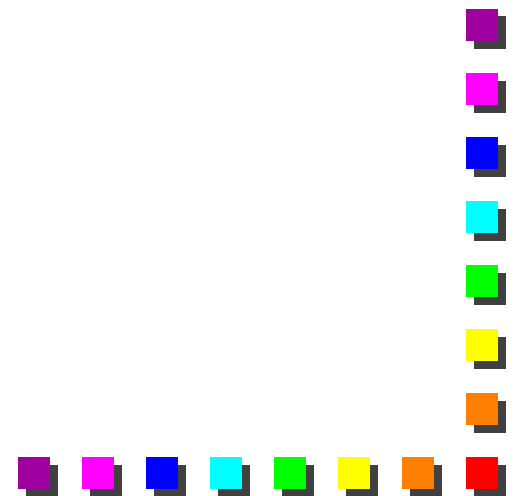
## Richard West

Boston University

Computer Science Department
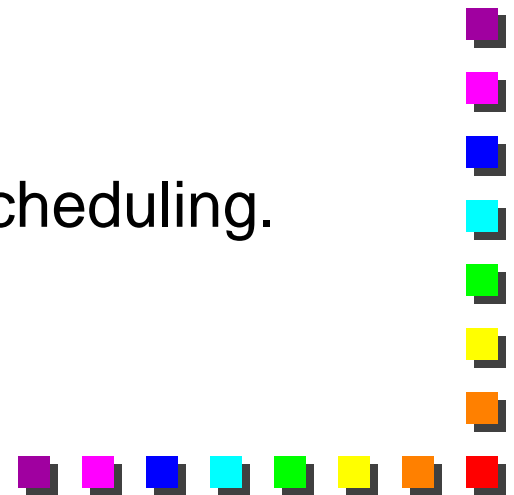
# Outline of Talk

- Problem Statement.
  - How to guarantee QoS to applications?
  - Variable resource demands / availability.
- Approach.
  - System mechanisms.
    - Dionisys.
  - System policies.
    - Dynamic Window-Constrained Scheduling.
- Conclusions.
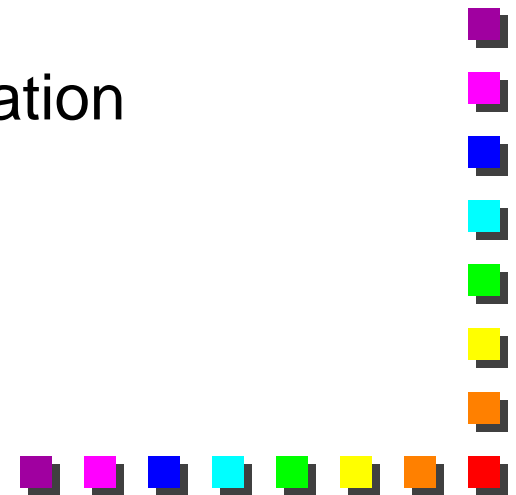
BOSTON UNIVERSITY Rich West (2001)

# Problem Statement

- Distributed, real-time (RT) applications:
  - e.g., VE, RT multimedia, tele-medicine, ATR.
  - Require QoS guarantees on end-to-end transfer of information.
- How do we guarantee QoS?
- Need system support to **maintain / maximize QoS**:
  - Policies & mechanisms.
  - **Adaptive** / **coordinated** resource management.

Rich West (2001)

# Application Characteristics

- Dynamic exchanges between processes.
  - The information (content & type) to be exchanged changes with time.
- Variable rates (bursts) of exchanges.
- Variable resource demands.
  - Bandwidth, CPU cycles, memory.
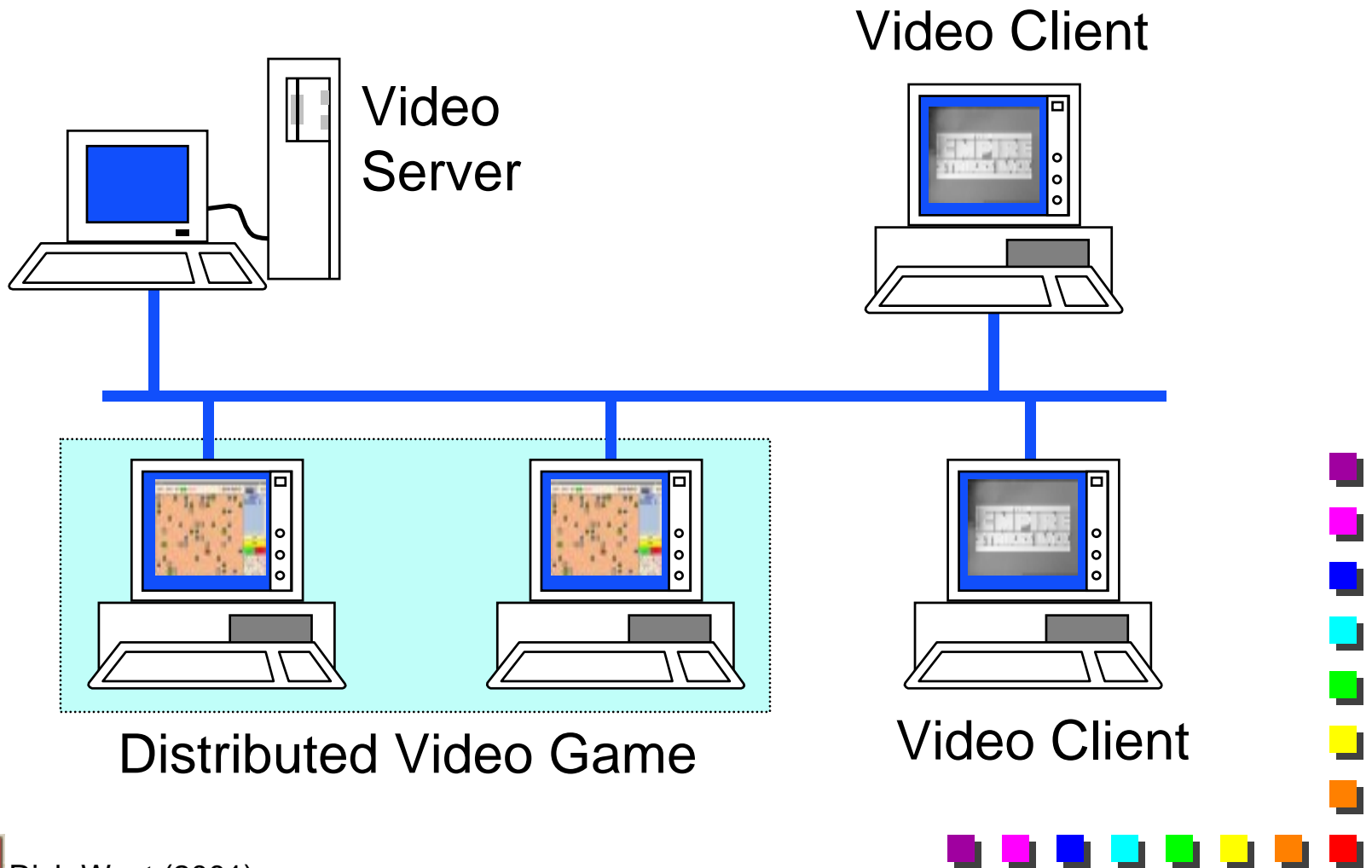- Variable **QoS requirements** on information exchanged.

Rich West (2001)
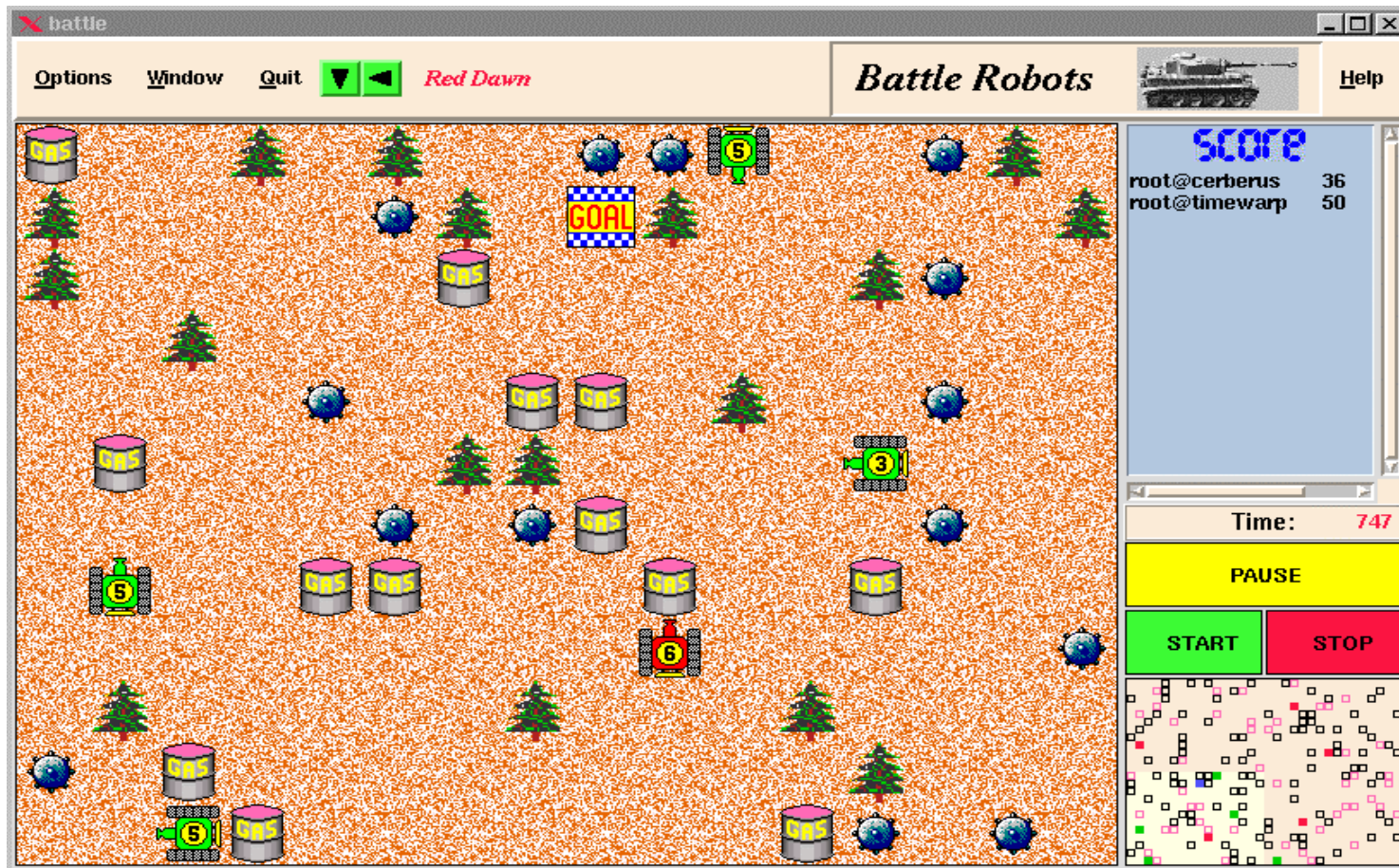
# QoS Requirements

- **Delay**: e.g., max end-to-end delay, delay variation.

- **Loss-tolerance, fidelity, resolution**:

  - Minimum degree of detail.

- **Throughput, rate**:

  - e.g., 30 fps video.

  - e.g., min/max updates per second to shared data.

- **Consistency constraints**:

  - *When, with whom* semantics.

Rich West (2001)

# Example Scenario



Video Server

Video Client

Distributed Video Game

Video Client

Rich West (2001)

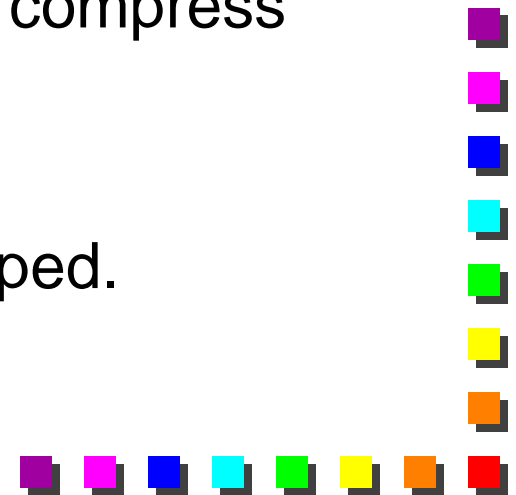# Example: Distributed Video Game



Rich West (2001)

# Distributed Video Game

- Requires consistency of shared (tank) objects.
- Here **QoS (and, hence, resource) requirements vary** with time based on current state of application.

- Application-level **spatial** & **temporal** semantics.
  - Exchange state info only when two objects less than distance $d$ apart.
  - Exchange position, orientation and (varying amounts of) graphical info about shared objects based on their distance apart.

Rich West (2001)

# Example: Video Server

- QoS requirements: Loss-tolerance and frame rate.

- Suppose a client requires at least 15fps playback rate but prefers 30fps.

- If network bandwidth is limited:

  - **Adapt CPU service.**

    - e.g. allocate more CPU cycles to compress video info.

  - **Adapt network service.**

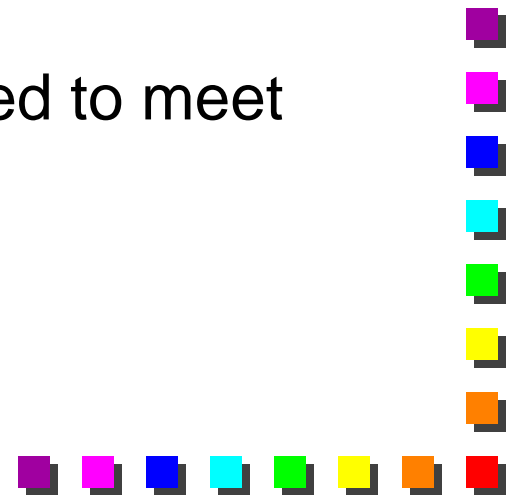    - e.g. allow 1 frame in 2 to be dropped.

Rich West (2001)

# Video Server (continued)

- If CPU cycles are limited:
  - **Adapt CPU service**.
    - If possible, reduce frame generation rate.
  - **Adapt network service**.
    - e.g. ensure no frames are now dropped.
- If CPU and network resources are limited:
  - Adapt to new QoS region / requirements if possible! Re-negotiation?

Rich West (2001)

# Summary of Problem

- Need to maintain / maximize QoS on end-to-end transfer of information.

- **Varying** resource requirements & availability.

- **Static** resource allocation too expensive.
    - Poor resource utilization & scalability.

- Suppose enough resources are reserved to meet the minimum needs of all applications.
    - How can we do better?

Rich West (2001)

# Approach

- **Dionisys** QoS mechanisms.
  - Allow real-time applications to specify:
    - <u>How</u> actual service should be adapted to meet required / improved QoS.
    - <u>When</u> and <u>where</u> adaptations should occur.
- **Coordinated** CPU and network management.
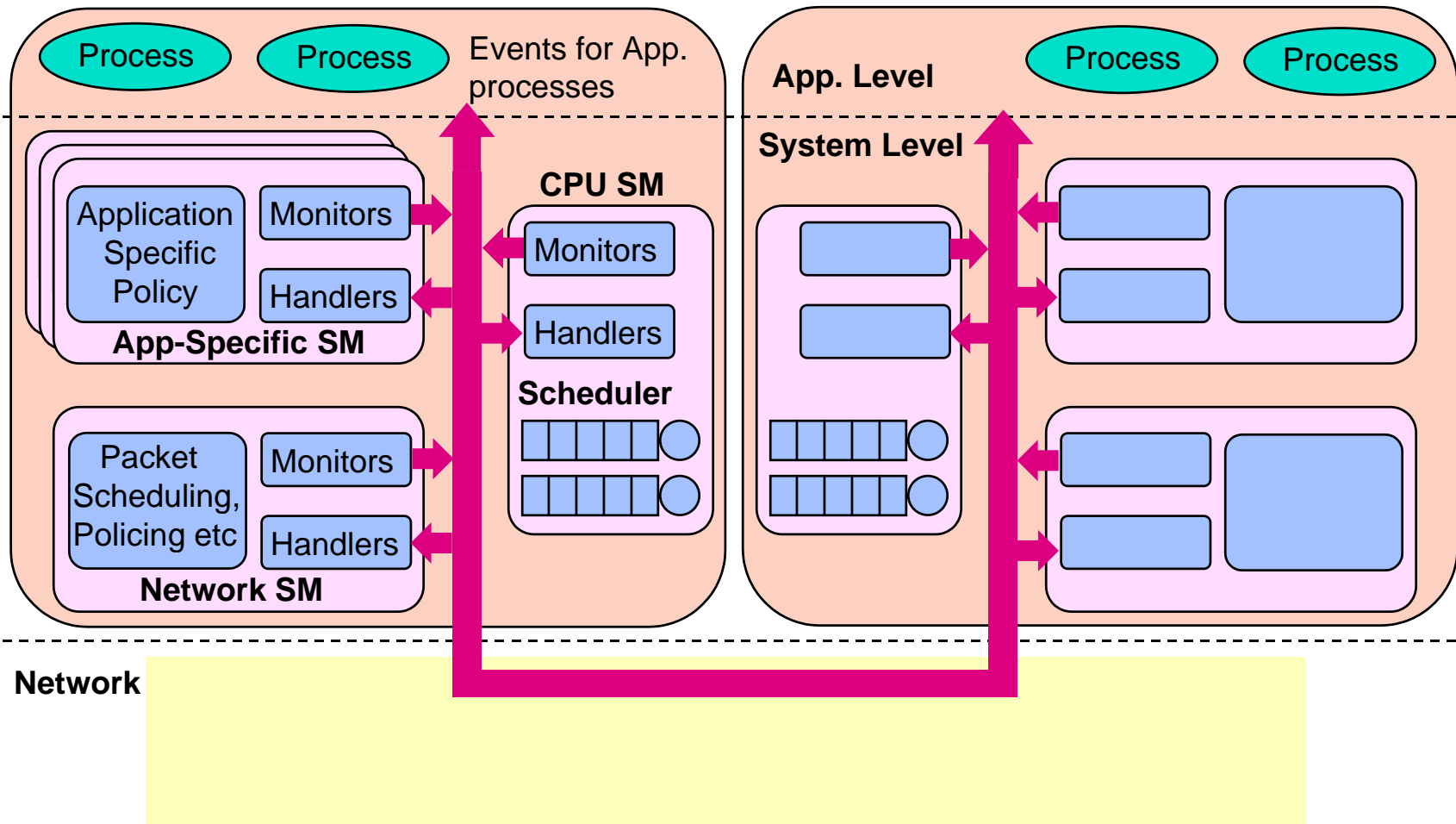  - **Dynamic Window-Constrained Scheduling.**

BOSTON UNIVERSITY  Rich West (2001)

# Dionisys

- Key components:
    - Service managers (SMs).
    - Monitors - influence <u>when</u> to adapt.
    - Handlers - influence <u>how</u> to adapt.
    - Events.
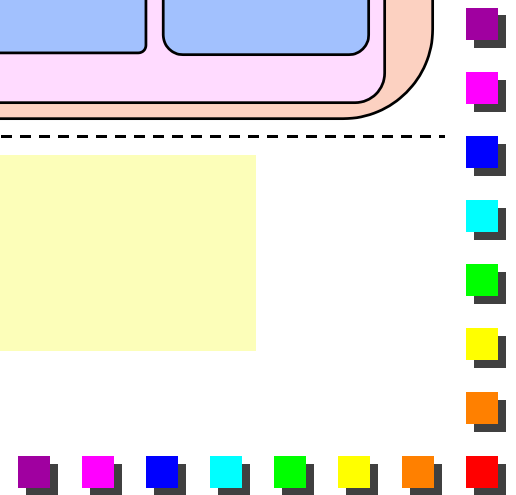        - Delivered to SMs, <u>where</u> adaptation is needed.
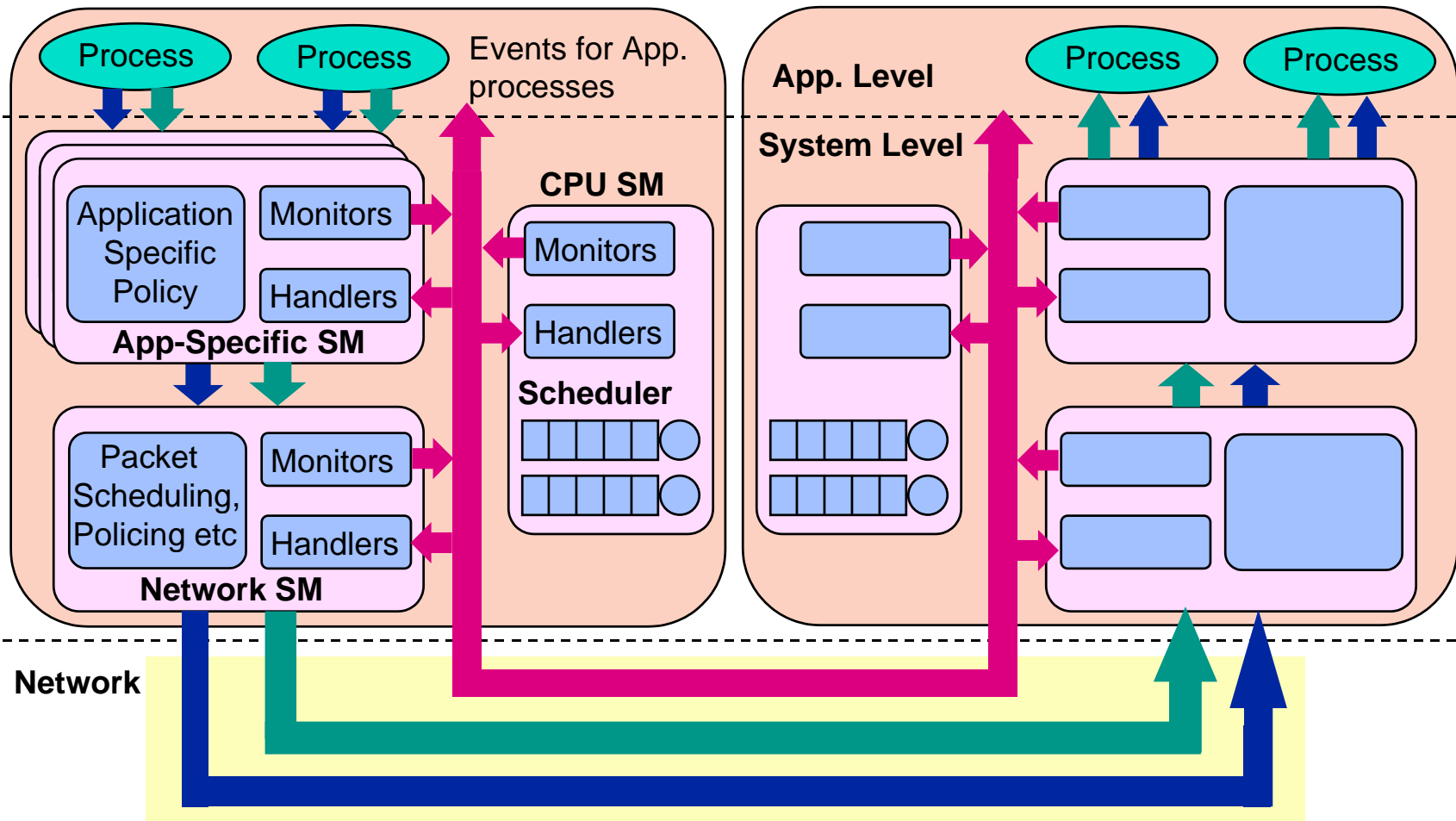    - Event channels.

Rich West (2001)

**SOURCE HOST**

**DESTINATION HOST**

Process   Process   Events for App. processes

App. Level   Process   Process

System Level

Application Specific Policy   Monitors
Handlers
**App-Specific SM**

**CPU SM**
Monitors
Handlers
**Scheduler**

Packet Scheduling, Policing etc   Monitors
Handlers
**Network SM**

**Network**

➡ Control path

Rich West (2001)

**SOURCE HOST**

**DESTINATION HOST**

Process   Process   Events for App. processes

App. Level

Process   Process

System Level

Application Specific Policy   Monitors   Handlers

**App-Specific SM**

**CPU SM**

Monitors   Handlers

**Scheduler**

Packet Scheduling, Policing etc   Monitors   Handlers

**Network SM**

Network

Control path   QoS attribute path   Data path

BOSTON UNIVERSITY   Rich West (2001)

SOURCE HOST

DESTINATION HOST

Process  Process  Events for App. processes

App. Level

Process  Process

System Level

Application Specific Policy  Monitors  Handlers

App-Specific SM

CPU SM

Monitors

Handlers

Scheduler

Monitors  Handlers

Scheduler

Monitors  Application Specific Policy

Handlers

Packet Scheduling, Policing etc  Monitors

Handlers

Network SM

Monitors

Handlers

Buffer Mgmt. etc

Network

Control path          QoS attribute path          Data path

BOSTON UNIVERSITY  Rich West (2001)

# Dionisys Key

**Process** Application process.

➡ Event channel.

➡ QoS attribute channel (shared memory on a single host).

➡ Data channel.

Service Manager (SM) e.g., CPU SM.

SM functions: App-specific monitors, handlers and service policy.
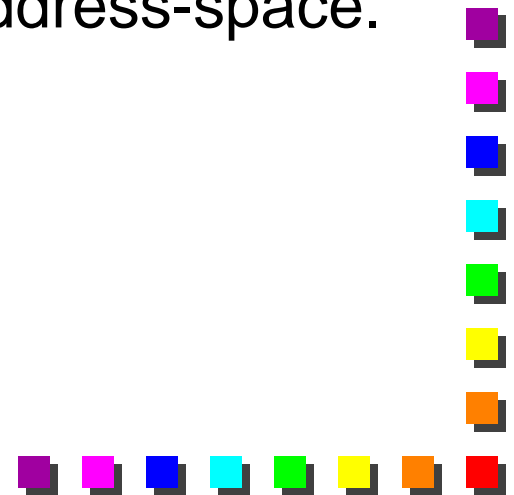
Host machine.

# Service Managers

- Responsible for:
  - Monitoring application-specific service.
  - Handling events for service adaptation.
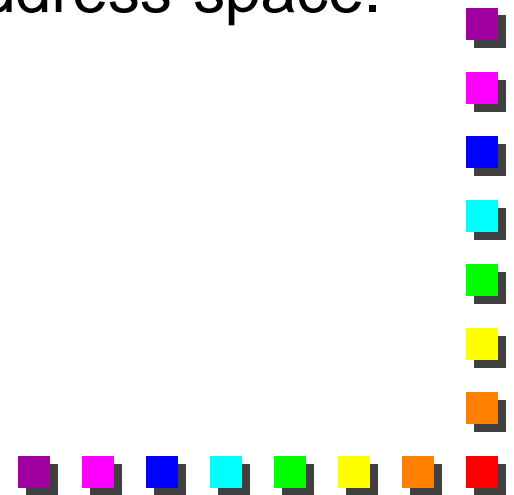  - Providing service to applications.
    - Resource allocation.
- Kernel level threads.

Rich West (2001)

# Monitors

- Functions that monitor a specific service.

- Influence <u>when</u> to adapt service provided to an application.

    - e.g., QoS below desired level, or unacceptable.
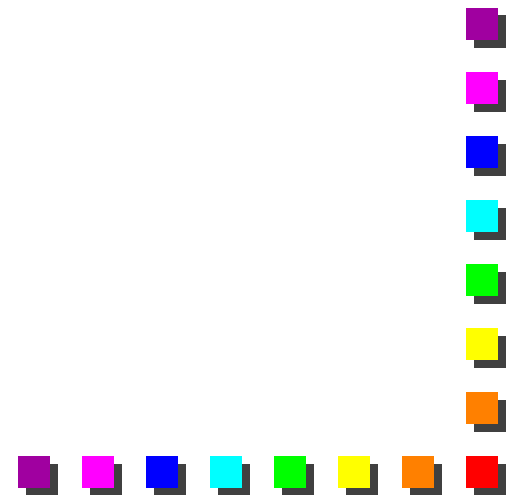
- Compiled into objects.

    - Dynamically-linked into target SM address-space.

Rich West (2001)

# Handlers

- Functions executed in SMs to decide <u>how</u> to adapt service provided to an application.

    - e.g., increase / decrease CPU cycles, or network bandwidth.

- Compiled into objects.

    - Dynamically-linked into target SM address-space.
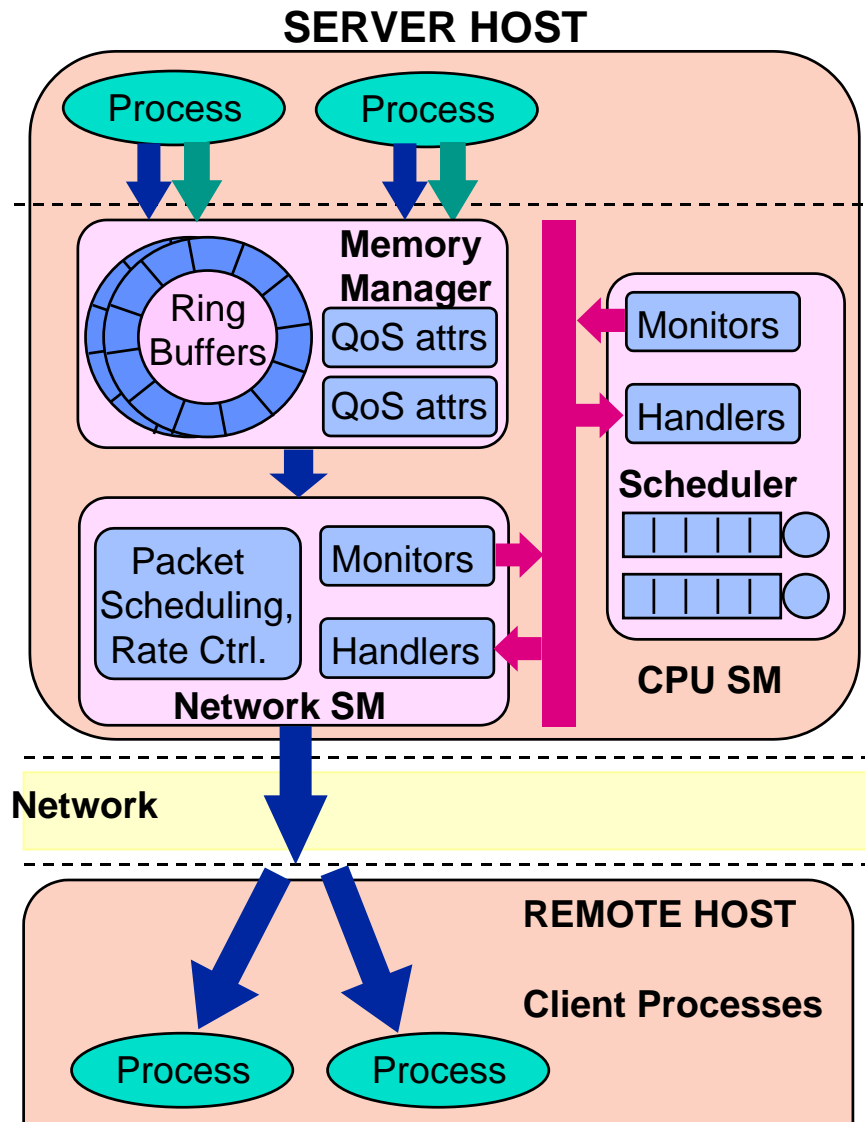
Rich West (2001)

# Events

- Generated <u>when</u> service adaptation is necessary.

- Delivered to handlers <u>where</u> service needs adapting.

- Have attributes that influence extent to which service is adapted.
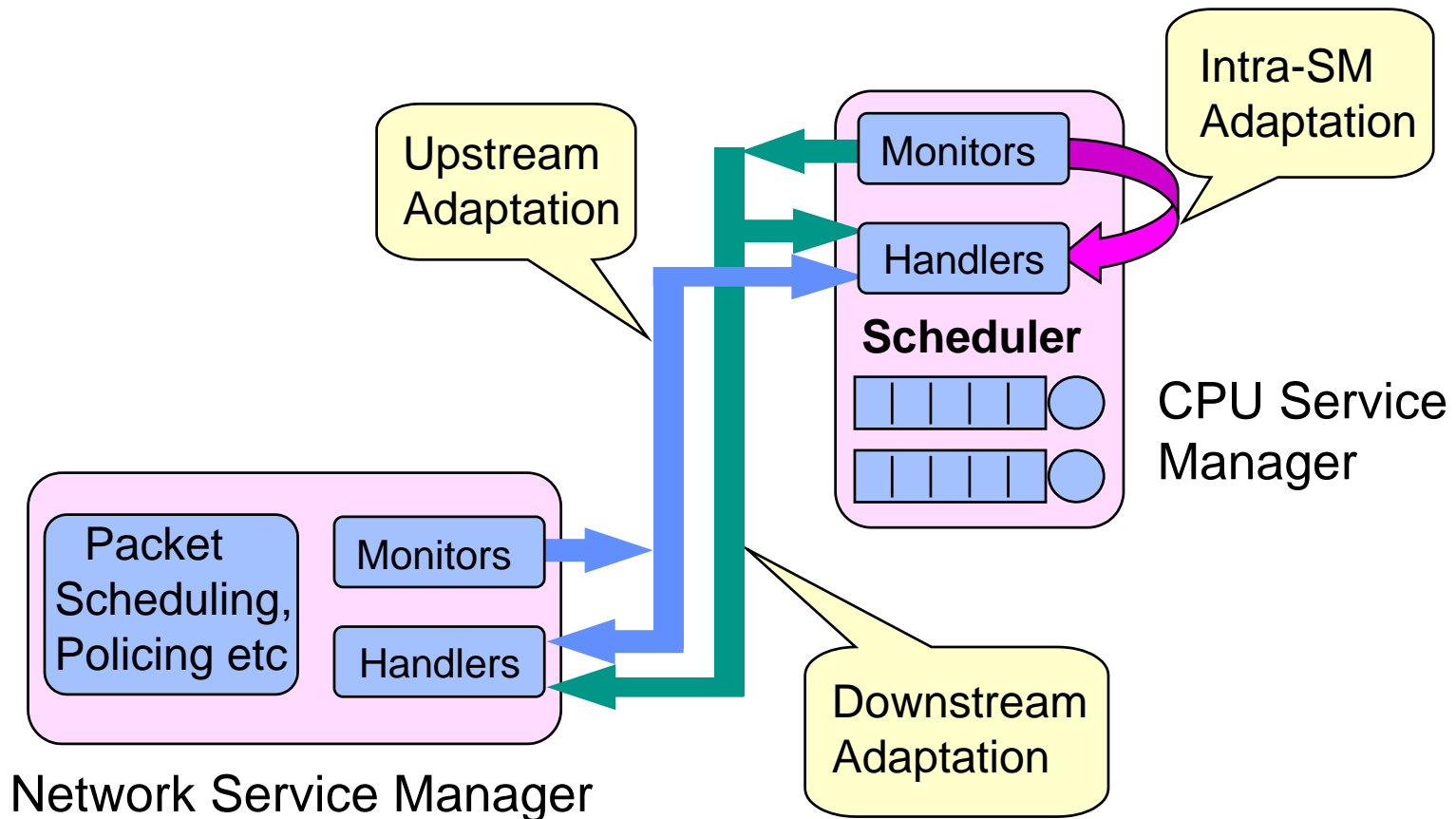
  - "Quality Events".

# Event Channels



Monitors       Handlers

SM 1

M(1)   M(2)    M(m1)    H(1)   H(2)    H(h1)

. . .

SM 2

M(1)   M(2)    M(m2)    H(1)   H(2)    H(h2)

. . .

Rich West (2001)

**SERVER HOST**

Process  Process

Memory Manager

Ring Buffers

QoS attrs

QoS attrs

Monitors

Handlers

Scheduler

Packet Scheduling, Rate Ctrl.

Monitors

Handlers

**Network SM**

**CPU SM**

**Network**

**REMOTE HOST**

**Client Processes**

Process  Process

**Example: Video Server**

Rich West (2001)

# Example Adaptation Strategies



Intra-SM Adaptation

Upstream Adaptation

Monitors

Handlers

**Scheduler**

CPU Service Manager

Packet Scheduling, Policing etc

Monitors

Handlers

Downstream Adaptation

Network Service Manager

Rich West (2001)

# Adaptation Strategies (continued)

- **Upstream adaptation**:
  - Applied in direction opposing flow of data.
    - e.g. feedback congestion control.
- **Downstream adaptation**:
  - Applied in direction corresponding to flow of data.
    - e.g. forward error correction.
- **Intra-SM adaptation**:
  - Applied to current service manager.
  - Lacks coordination between SMs.

BOSTON UNIVERSITY Rich West (2001)

# Adaptation Example: Video Server

- QoS requirements: Loss-tolerance and frame rate.
- If network bandwidth is limited:
  - Apply **upstream adaptation** to increase CPU cycles to e.g. compress video information.
  - Apply **intra-SM adaptation** in the network SM to increase loss-tolerance.

Rich West (2001)

# Adaptation Example (continued)

- If CPU cycles are limited:
    - Apply **intra-SM adaptation** in the CPU-SM to reduce, for example, frame (generation) rate.
    - Apply **downstream adaptation** to reduce loss-tolerance.

Rich West (2001)

# Experimental Scenario - Part 1

- Server-side processes (one per stream):
  - Generate data for streaming to remote clients.
    - Stream of MPEG-1 I-frames (160x120 pixels) per generator process.
    - Data placed in circular queues in shared memory.
- QoS attributes associated with each data stream:
  - Min / Max / Target frame rate.
- "Quality" event channels between Network and CPU service managers.

Rich West (2001)

# Experimental Scenario - Part 2

- Client-side processes (one per stream):
  - Decode and playback incoming frames.
- SparcStation Ultra-2 170Mhz dual processor server, running Solaris 2.6 connected via switched 100Mbps Ethernet to one client (w/ UDP connection).
- 3 Streams:
  - Stream 1: Target 30fps +/- 10% (3000 frames)
  - Stream 2: Target 20fps +/- 10% (2000 frames)
  - Stream 3: Target 10fps +/- 20% (1000 frames)
  - 3 second exponential idle time every 1000 frames.

BOSTON UNIVERSITY Rich West (2001)

# Adaptation in Video Server

- (**Downstream**) CPU SM monitors frame generation rate.

- (**Upstream**) Net SM monitors frame transmission rate.

- Apply adaptation if (monitored rate != target rate).

- All monitors / SMs run at 10mS intervals.

Rich West (2001)

# Adaptation Handlers

- CPU-Level:
  - Adjust priorities & time-slices of generator processes by a function of target and monitored service rates.

- Network-Level:
  - Invoke rate control if monitored rate exceeds maximum rate.
  - Raise priority of packet stream $S_i$ if its service falls below minimum service rate.
    - i.e., alter bandwidth allocation $(y_i - x_i) / y_i$.

Rich West (2001)

# Adaptive Rate Control Block Diagram



Rich West (2001)

# Quality Functions - Example

Q

180
160
140
120
100
80
60
40
20
0

8  10  12        17  20  23    26    30    34

**S (Rate, fps)**

- Can embed quality functions into handlers.
- Service adaptation is a function of actual and required service of all applications.

Rich West (2001)
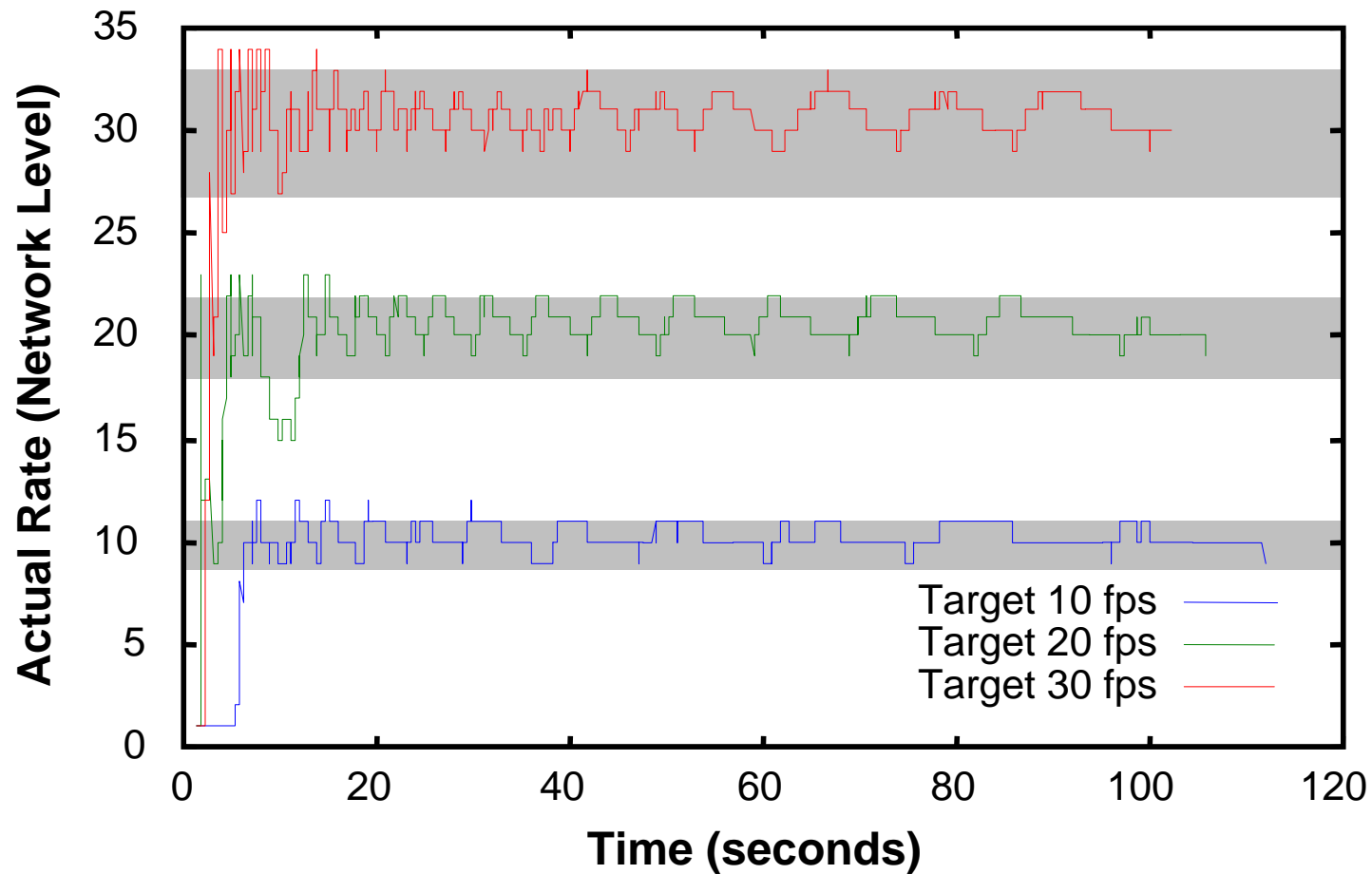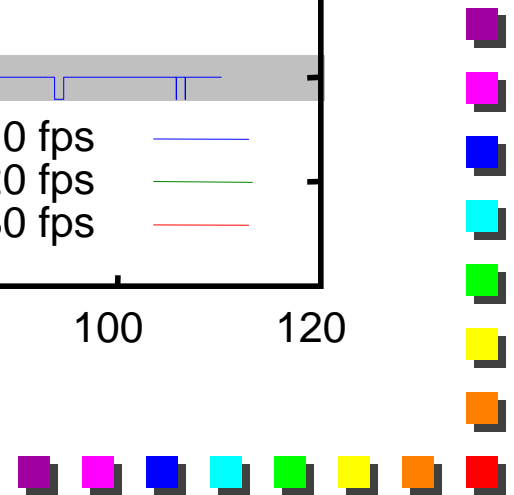
# Non-Adaptive Rate Allocating Service



Rich West (2001)

# Non-Adaptive Rate Controlled Service



Rich West (2001)

# Network Rate - Upstream Adaptation



Target 10 fps
Target 20 fps
Target 30 fps

Rich West (2001)

# Network Rate - Downstream Adaptation
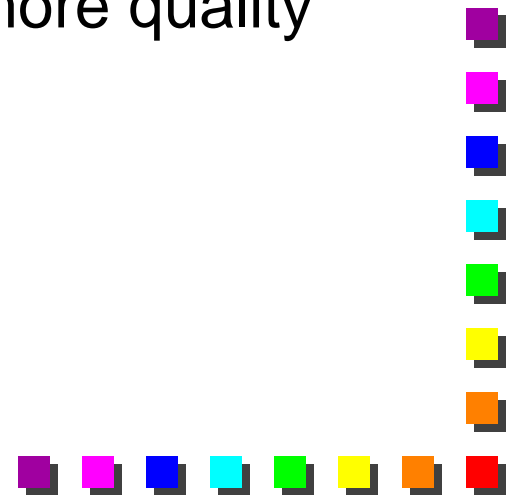


Rich West (2001)

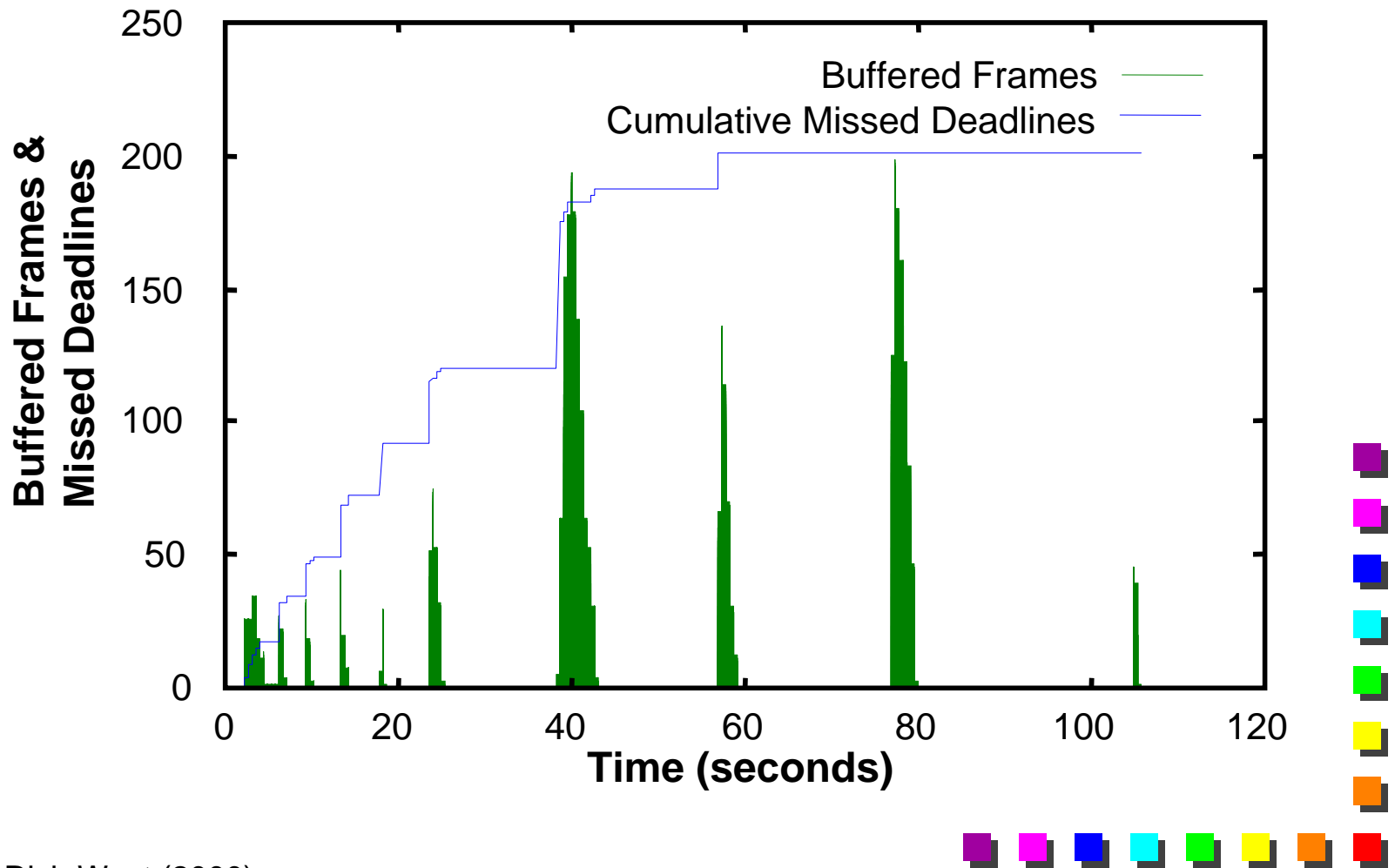# Comparison of Rate Control Methods



Rich West (2001)

# Rate Control

- Upstream adaptation leads to poorer rate control.
    - Longer time to reach steady state.
    - More prominent "sawtooth" effect as target rate is tracked.
    - Larger fluctuations of actual rate from target.
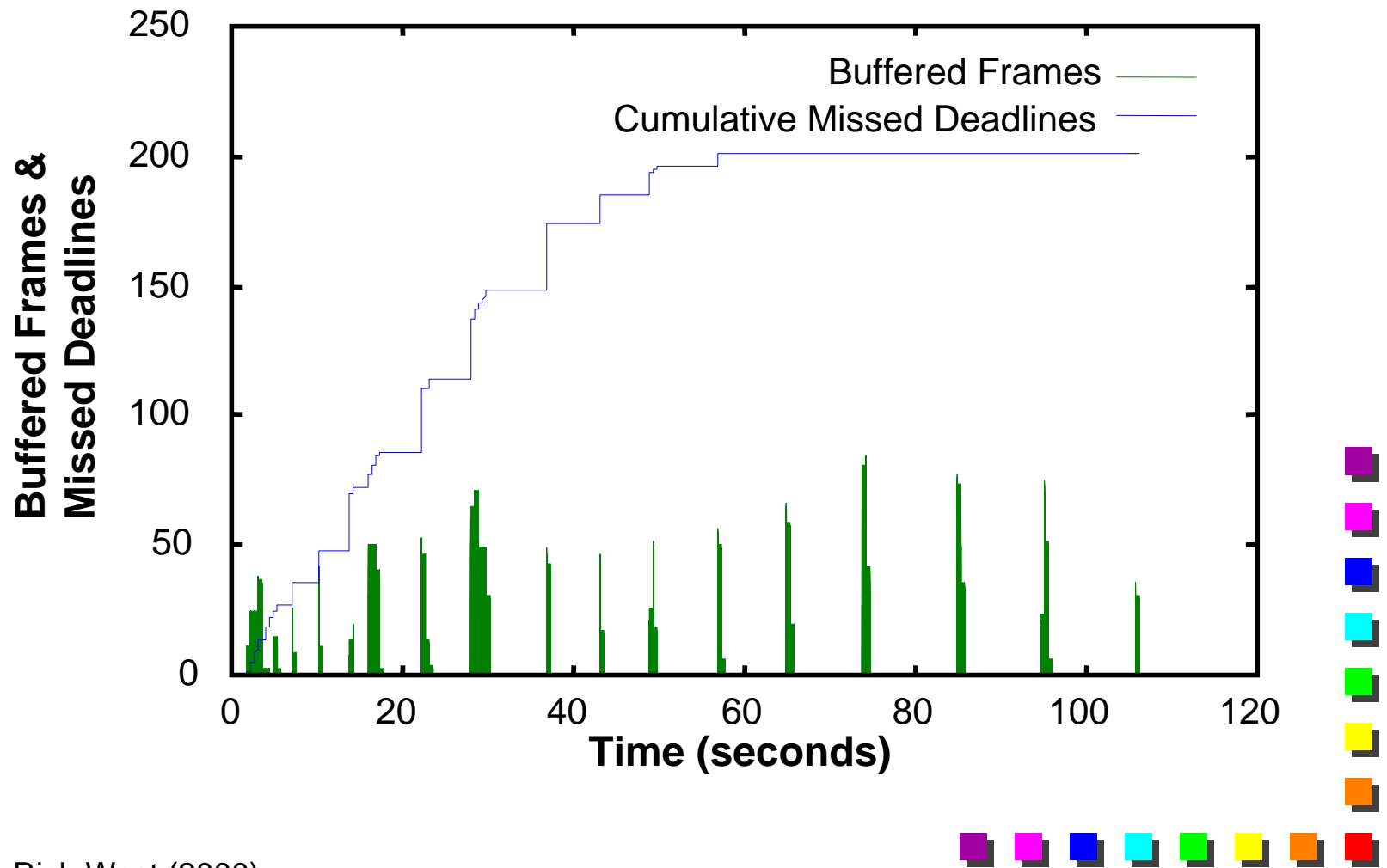        - Better tracking of target rate for more quality critical streams.

Rich West (2001)

# Upstream Adaptation - 10fps
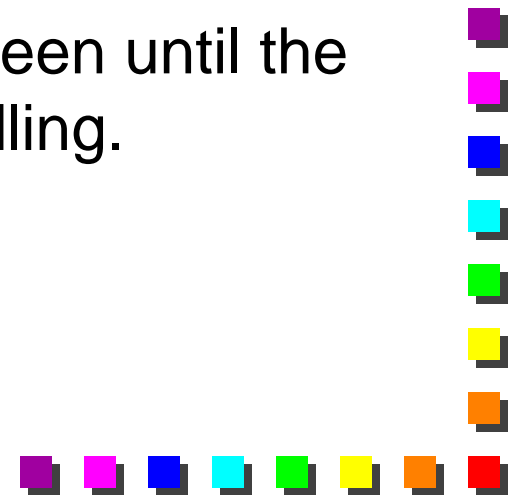


Rich West (2000)

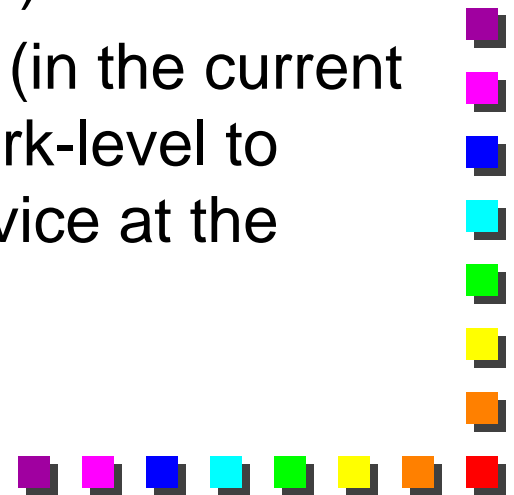# Downstream Adaptation - 10fps



Rich West (2000)

# Buffering

- Upstream adaptation leads to greater variance in buffer usage, compared to downstream / intra SM adaptation.

  - Network monitor triggers "request" for generation of frames "too late". That is, after buffer has emptied.

  - Effect of an event being raised not seen until the next "phase" of monitoring and handling.
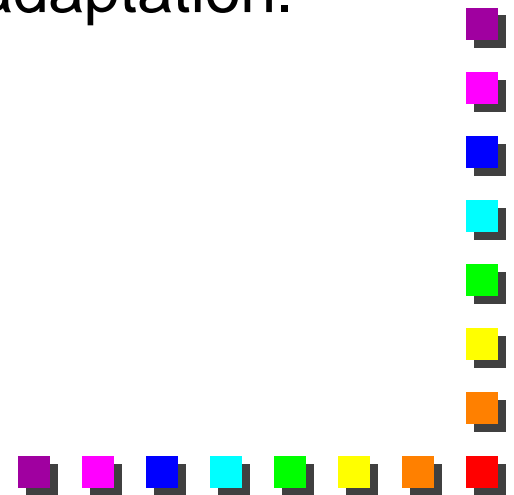
Rich West (2001)

# Missed Deadlines

- Higher buffering variance and, consequently, higher queueing delays, imply potentially higher consecutive numbers ("bursts") of missed deadlines.

- Downstream adaptation can reduce the number of consecutive deadlines missed at any time by:
  - Providing more accurate (responsive) service.
  - By effecting changes "more quickly" (in the current event/monitoring cycle) at the network-level to compensate for inadequacies in service at the CPU-level.
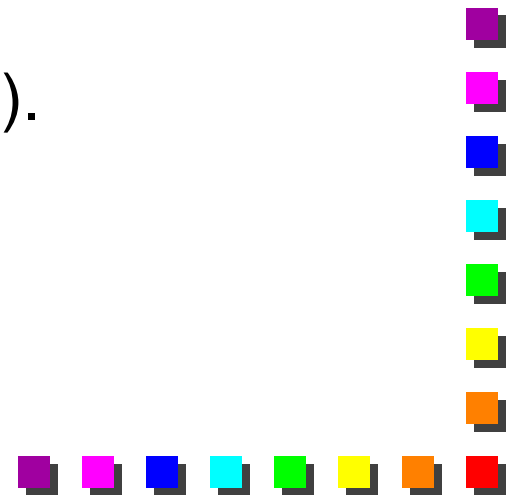
# Summary

- **Dionisys** QoS mechanisms allow real-time applications to specify:
    - <u>How</u> actual service should be adapted to meet required / improved QoS.
    - <u>When</u> and <u>where</u> adaptations should occur.
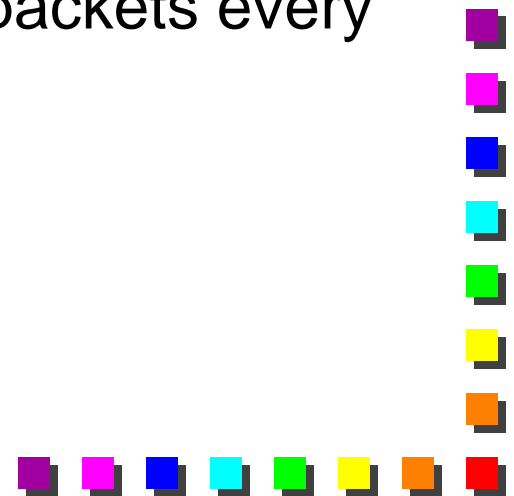- **Flexible** approach to run-time service adaptation.

# What About Service Policies?

- Certain applications can tolerate lost / late information.

- Restrictions on:
    - when losses of info can occur.
    - when info must be generated.

- Need real-time scheduling of:
    - threads / processes (info generators).
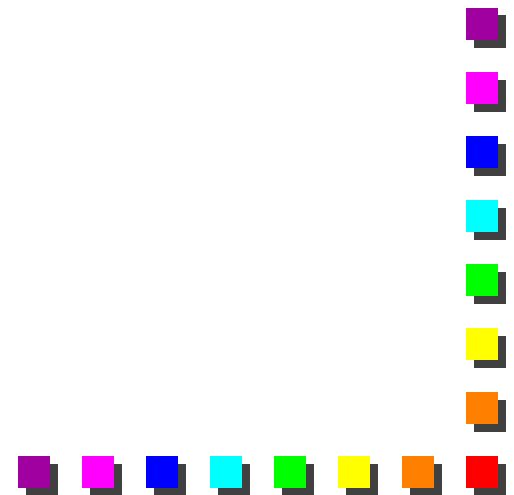    - packets (info carriers).

Rich West (2001)

# DWCS

- Dynamic Window-Constrained Scheduling of:
  - Threads
    - "Guarantee" **minimum** quantum of service every fixed window of service time.
  - Packets
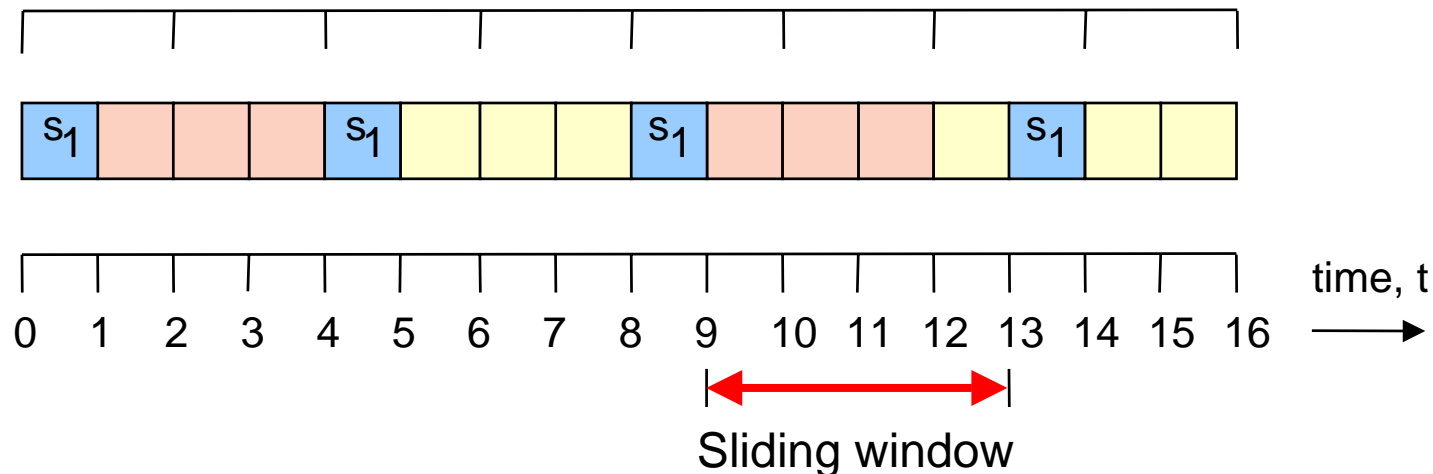    - "Guarantee" at most **x** late / lost packets every window of **y** packets.

Rich West (2001)

# DWCS Packet Scheduling

- Two attributes per packet stream, $S_i$:
    - Request period, $T_i$.
        - Defines interval between deadlines of consecutive pairs of packets in $S_i$.
    - Window-constraint, $W_i = x_i/y_i$.
        - Essentially, a "loss-tolerance".
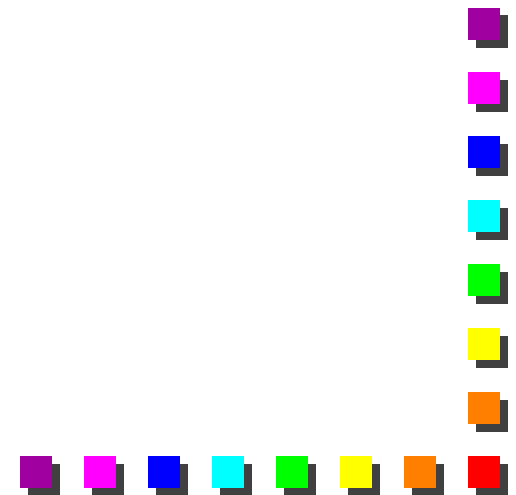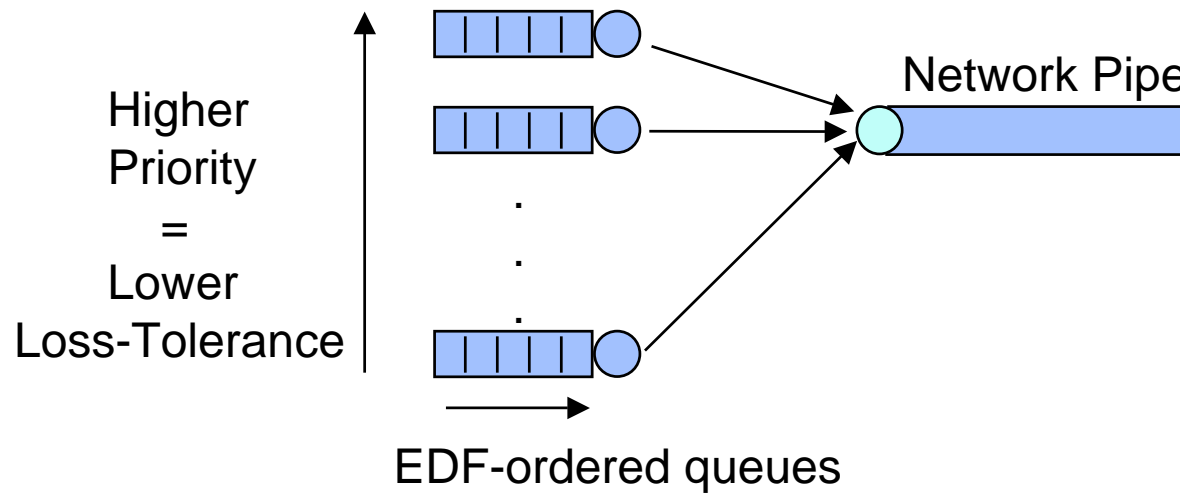
Rich West (2001)

# "x out of y" Guarantees

- e.g., Stream $S_1$ with $C_1=1$, $T_1=2$ and $W_1=1/2$



Sliding window

time, t

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16

- Feasible schedule if "x out of y" guarantees are met.

# DWCS - Original Conceptual View



Higher
Priority
=
Lower
Loss-Tolerance
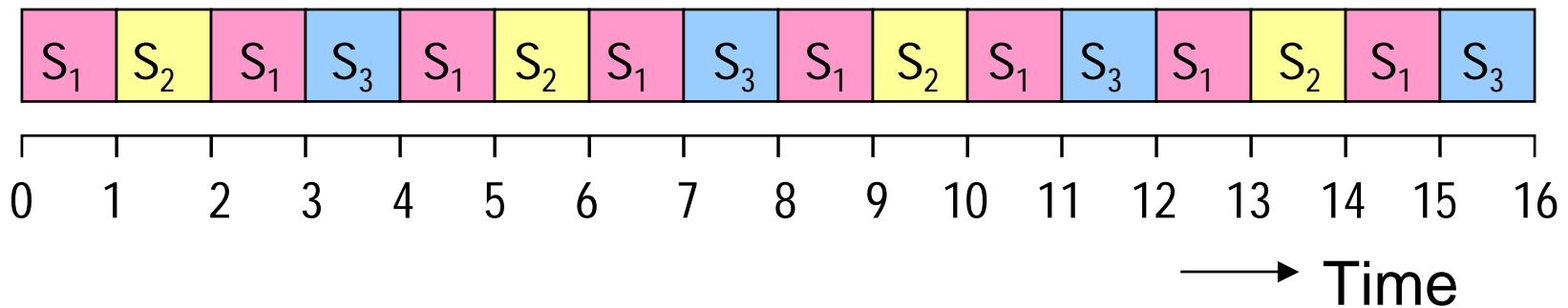
Network Pipe

EDF-ordered queues

# (x,y)-firm DWCS: Pairwise Packet Ordering Table

| Precedence amongst pairs of packets |
|---|
| • Lowest window-constraint first |
| • Same non-zero window-constraints, order EDF |
| • Same non-zero window-constraints & deadlines, order lowest window-numerator first |
| • Zero window-constraints and denominators, order EDF |
| • Zero window-constraints, order highest window-denominator first |
| • All other cases: first-come-first-serve |

Rich West (2001)

# Example: "Fair" Scheduling

| $S_1$ | $S_2$ | $S_1$ | $S_3$ | $S_1$ | $S_2$ | $S_1$ | $S_3$ | $S_1$ | $S_2$ | $S_1$ | $S_3$ | $S_1$ | $S_2$ | $S_1$ | $S_3$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16
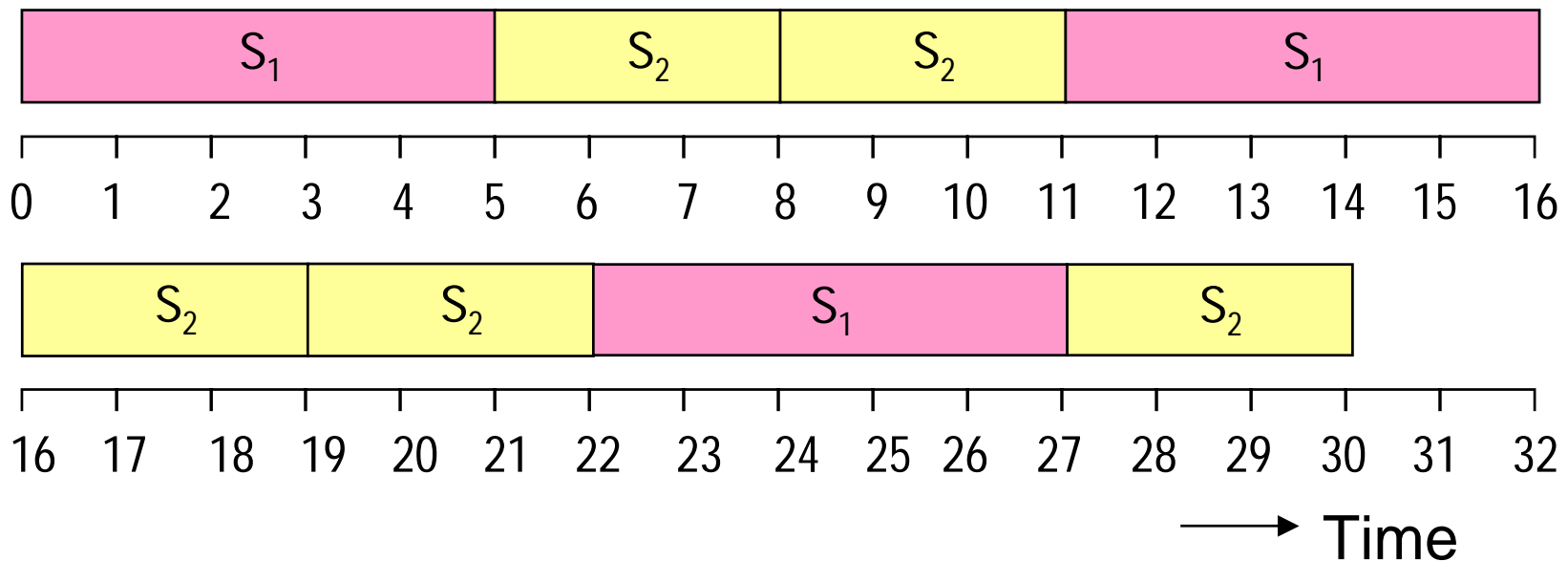
→ Time

$S_1$ 1/2(0) 1/1(1) 1/2(2) 1/1(3) 1/2(4)...

$S_2$ 3/4(0) 2/3(1) 2/2(2) 1/1(3) 3/4(4)...

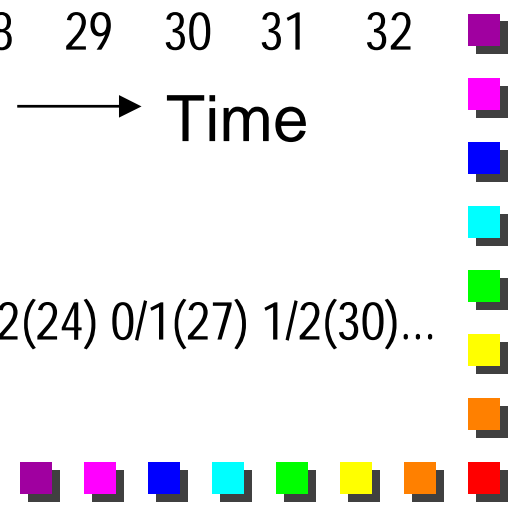$S_3$ 6/8(0) 5/7(1) 4/6(2) 3/5(3) 3/4(4) 2/3(5) 1/2(6) 0/1(7) 6/8(8)...

Rich West (2001)

# Example: Variable Length Packets

| S₁ | S₂ | S₂ | S₁ |
|---|---|---|---|

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16

| S₂ | S₂ | S₁ | S₂ |
|---|---|---|---|

16  17  18  19  20  21  22  23  24  25  26  27  28  29  30  31  32

→ Time

$S_1$ 1/2(0) 1/1(5) 1/2(10) 0/1(15) 1/2(20) 0/1(25) 1/2(30)…

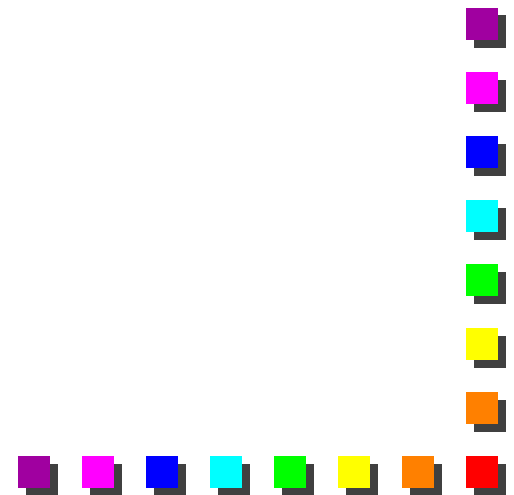$S_2$ 1/2(0) 0/1(3) 1/4(6) 1/3(9) 1/2(12) 0/1(15) 1/4(18) 1/3(21) 1/2(24) 0/1(27) 1/2(30)…

Rich West (2001)

# Window-Constraint Adjustment (A)

- For stream $S_i$ whose head packet is serviced **before** its deadline:
    - if ($y_i' > x_i'$) then $y_i'=y_i'-1$;
    - else if ($y_i' = x_i'$) **and** ($x_i' > 0$) then
        - $x_i'=x_i'-1$; $y_i'=y_i'-1$;
    - if ($x_i'=y_i'=0$) **or** ($S_i$ **is tagged**) then
        - $x_i'=x_i$; $y_i'=y_i$;
    - if ($S_i$ **is tagged**) then **reset tag**;

Rich West (2001)

# Window-Constraint Adjustment (B)

- For stream $S_j$ whose head packet **misses** its deadline:
    - if ($x_j' > 0$) then
        - $x_j' = x_j' - 1$; $y_j' = y_j' - 1$;
        - if ($x_j' = y_j' = 0$) then $x_j' = x_j$; $y_j' = y_j$;
    - else if ($x_j' = 0$) **and** ($y_j > 0$) then
        - violation! One solution…
        - $y_j' = y_j' + \varepsilon$;
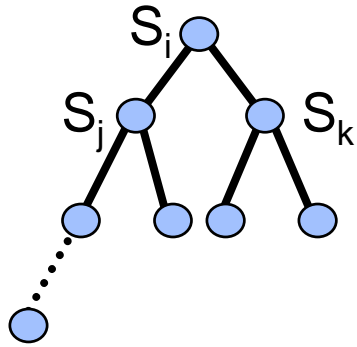        - **Tag $S_j$** with a violation;

Rich West (2001)

# DWCS Algorithm Outline

- Find stream $S_i$ with highest priority **(see Table)**

- Service head packet of stream $S_i$

- Adjust $W_i$' according to **(A)**

- **$Deadline_i = Deadline_i + T_i$**

- For each stream $S_j$ missing its deadline:

  - While deadline is missed:

    - Adjust $W_j$' according to **(B)**
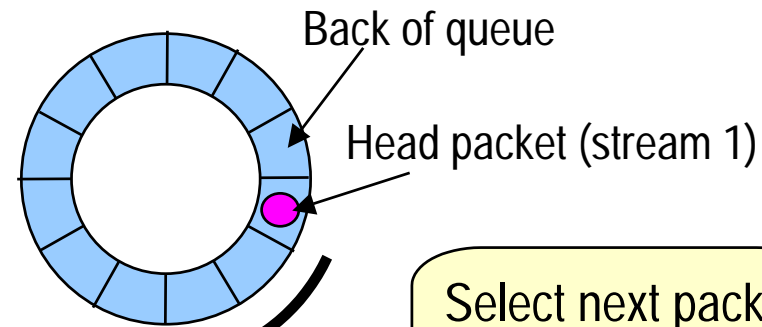    - Drop head packet of stream $S_j$ if droppable
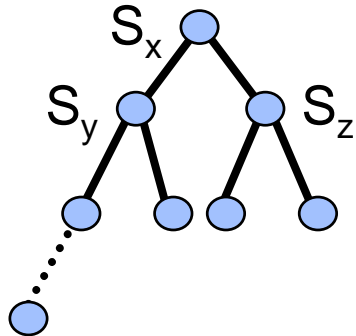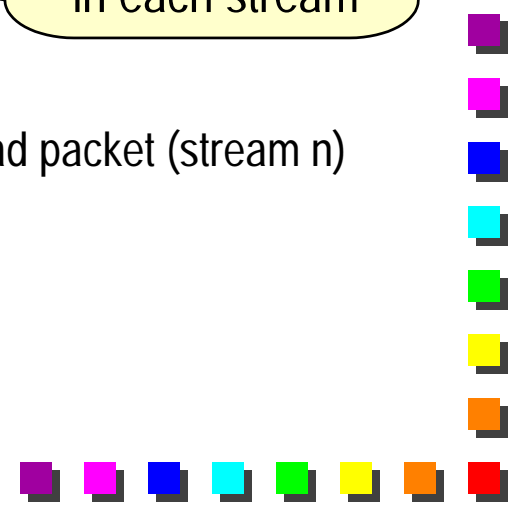    - **$Deadline_j = Deadline_j + T_j$**

Rich West (2001)

# DWCS Implementation

Deadline Heap

$S_i$

$S_j$     $S_k$

Loss-Tolerance Heap

$S_x$

$S_y$     $S_z$

Back of queue

Head packet (stream 1)

To back

Select next packet from head packets in each stream
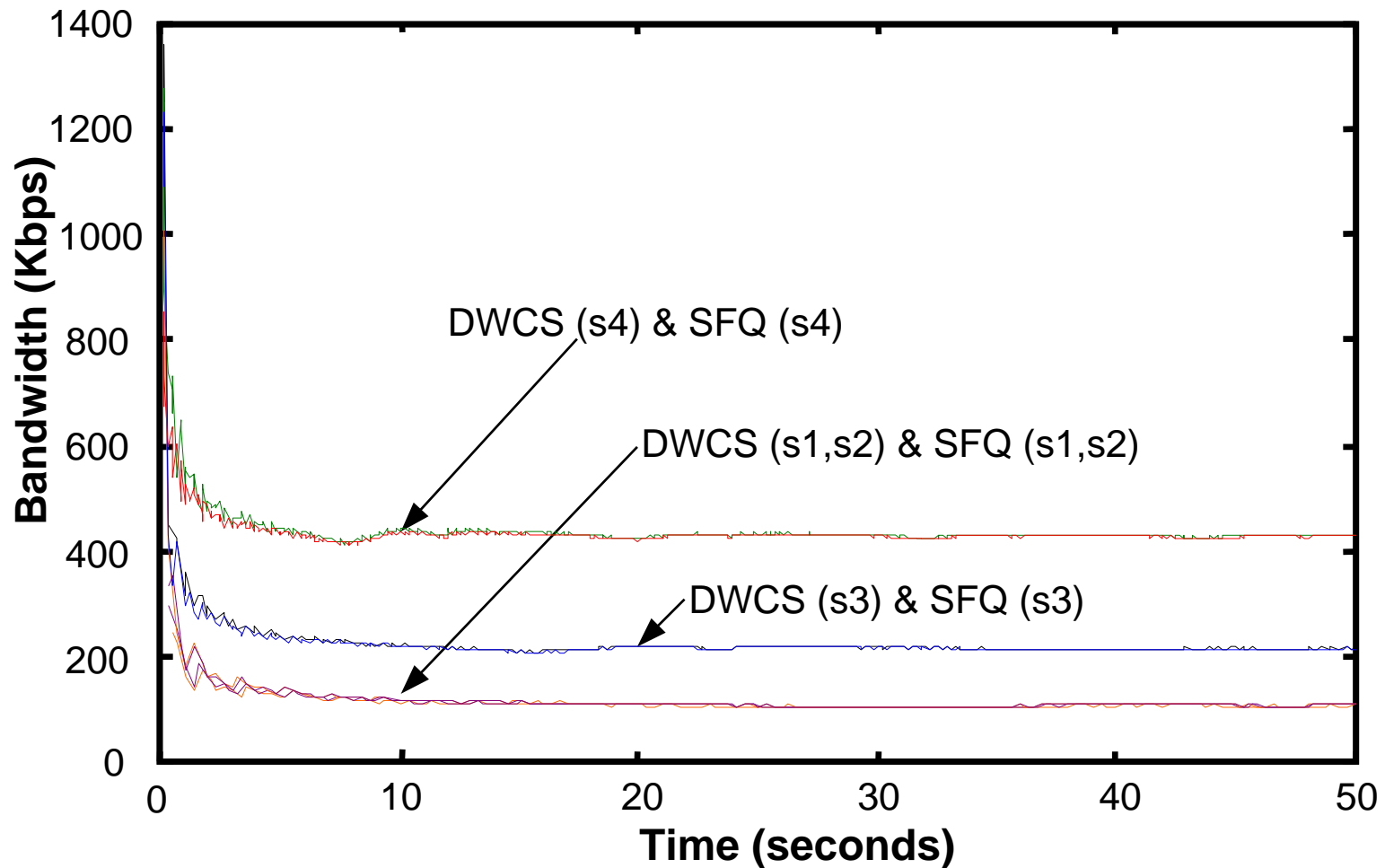
Head packet (stream n)

To back

Rich West (2001)

# Scheduling Overhead

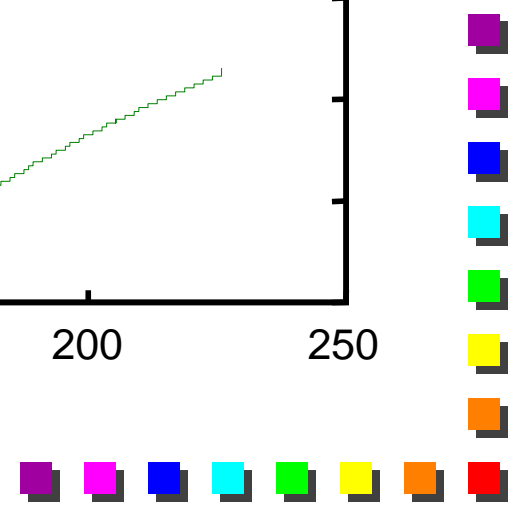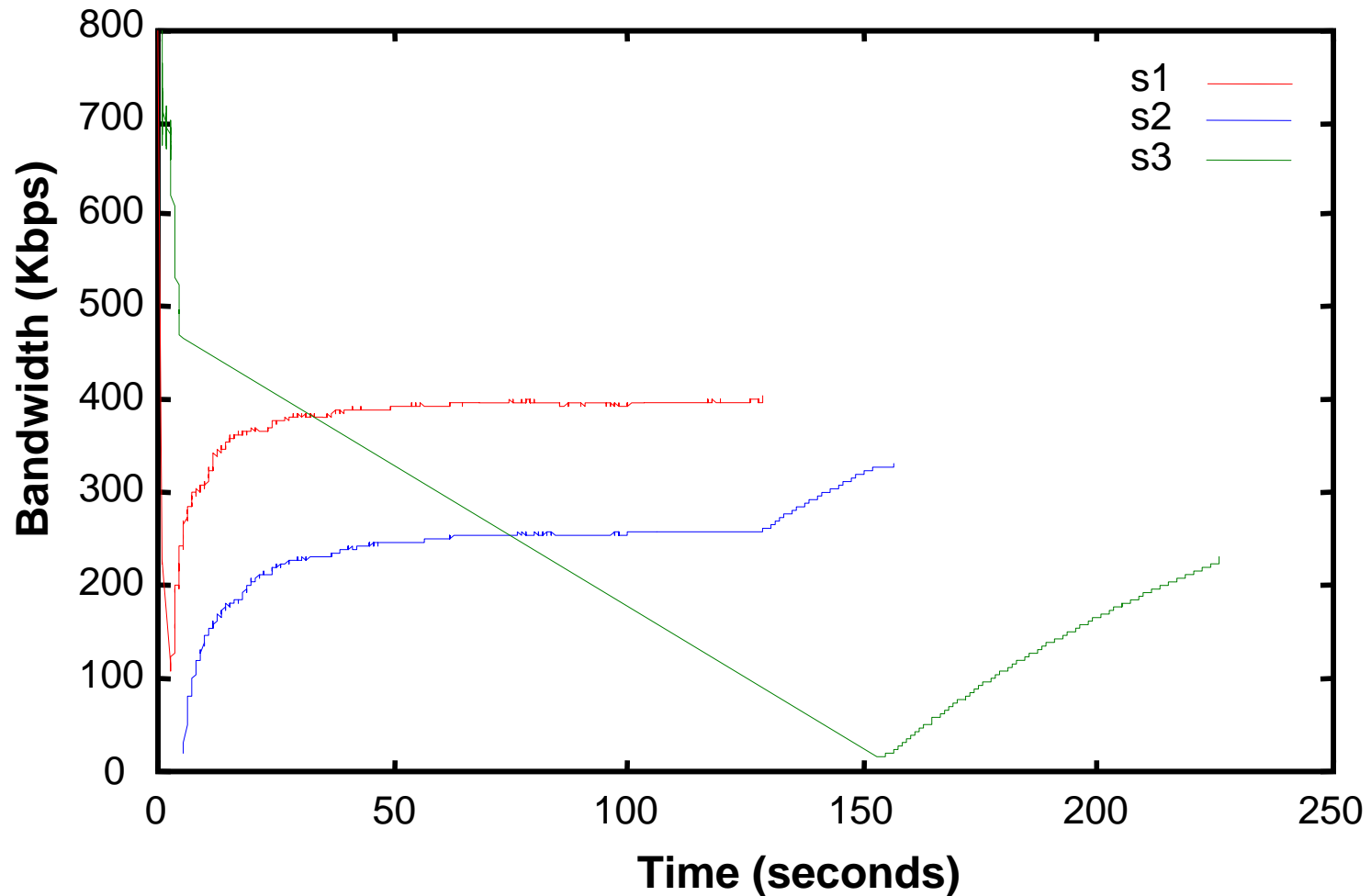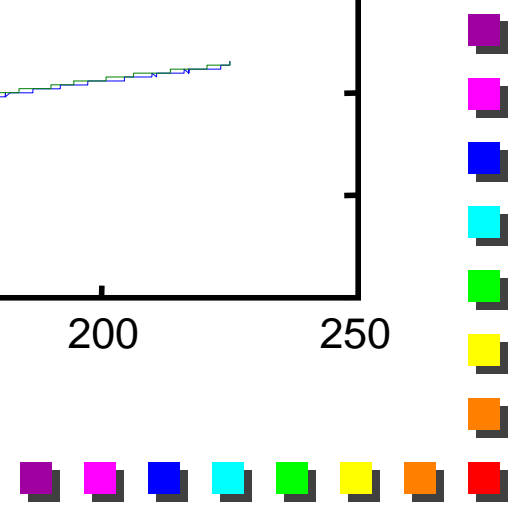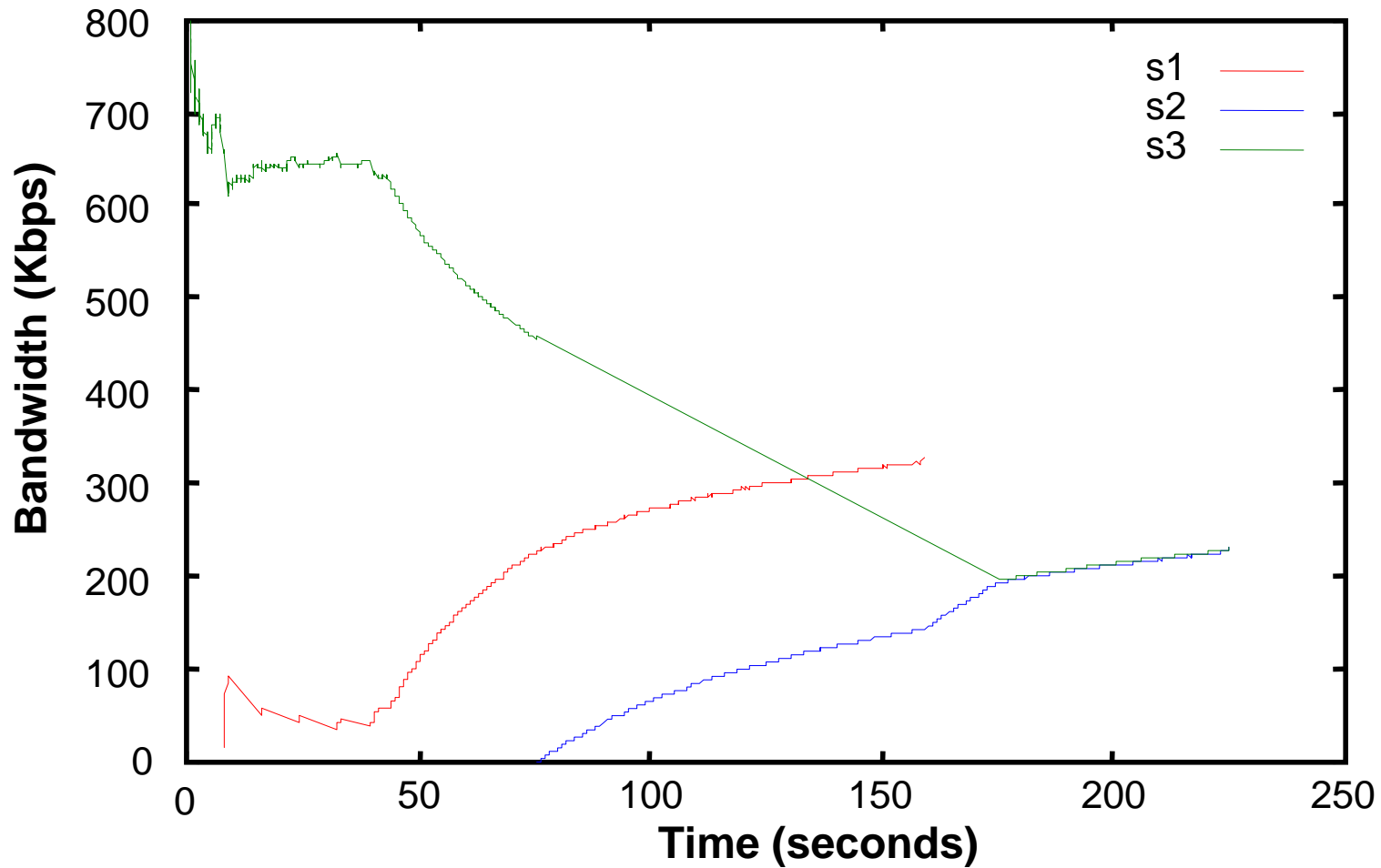# Fair Scheduling: b/w ratios:1,1,2,4
## W's=7/8,14/16,6/8,4/8



Rich West (2001)

# Mixed Traffic: W1=1/3,W2=2/3, W3=0/100,T1=1,T2=1,T3=∞

Rich West (2001)

# Mixed Traffic: W1=1/3,W2=2/3, W3=0/1500,T1=1,T2=1,T3=∞

Rich West (2001)

# Loss-Tolerance Violations (T=500, C=1)



Rich West (2001)

# DWCS Spreads Losses

DWCS

| X | X |  | X | X |  | X | X |  | X | X |  | X | X |  |

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15

⟶ Time

FIFO

| X |  |  | X | X | X | X | X | X |  |  |  | X | X | X |

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15

⟶ Time

■ Here, loss tolerance of 1/3 is violated more times with DWCS than FIFO, but losses are spread evenly.

# Approximation Overheads (T=500)



Rich West (2001)

# Approximation Overheads (T=200)



Rich West (2001)

# Deadlines Missed (T=500)

# Deadlines Missed (T=200)



Rich West (2001)

# Loss-Tolerance Violations (T=500)



Rich West (2001)

# Loss-Tolerance Violations (T=200)



Rich West (2001)

# DWCS - Recent Developments

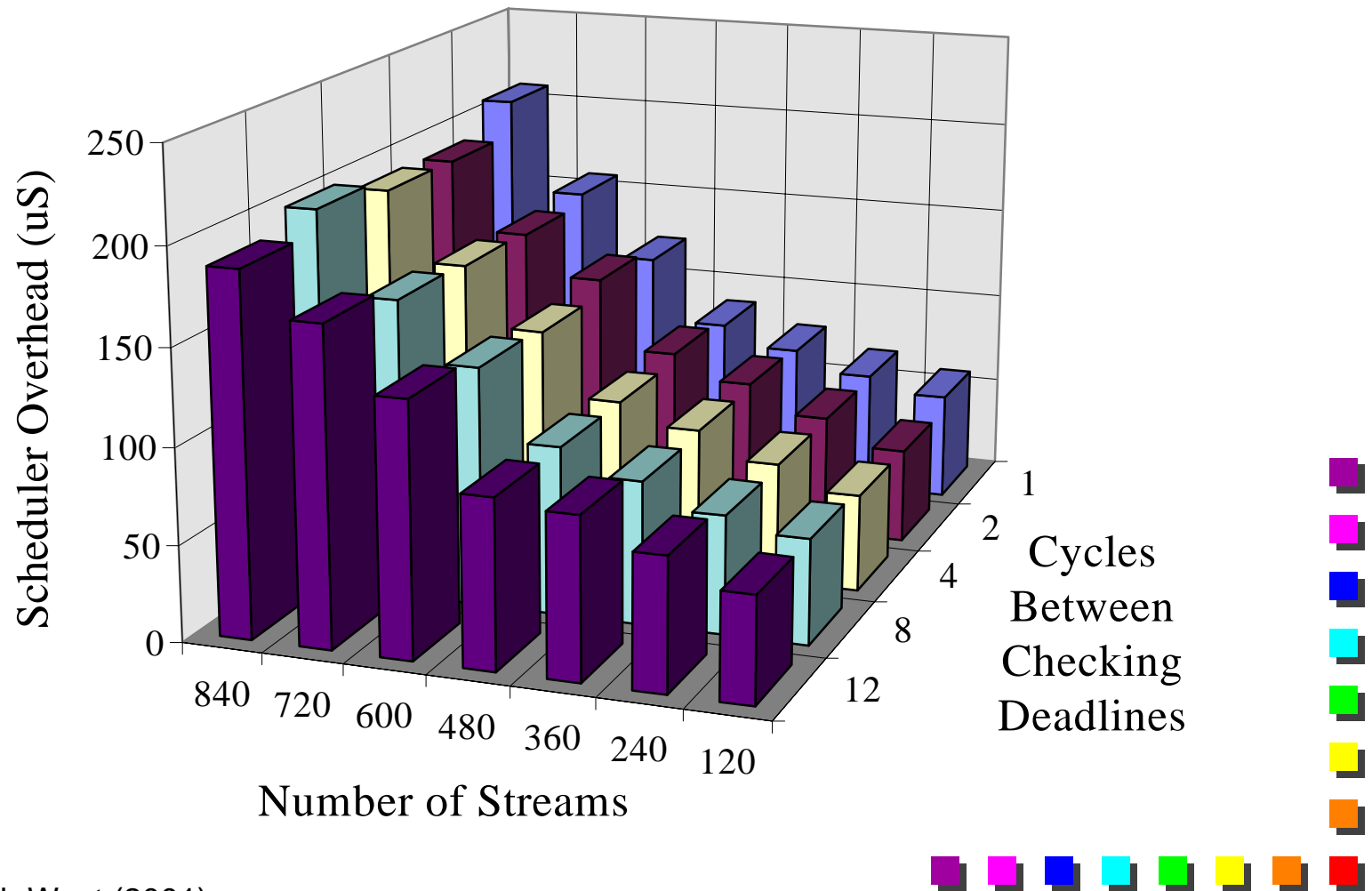- Support for **(x,y)-hard** deadlines as opposed to **(x,y)-firm** deadlines.
    - Bounded service delay.
    - Guaranteed service in a finite window of time.
    - Optimal (100%) utilization bound for fixed-length packets or (variable-length preemptive) threads.
- Replacement CPU scheduler in Linux kernel.
    - **www.cc.gatech.edu/~west/dwcs.html**

# (x,y)-Hard DWCS: Pairwise Packet Ordering Table

| Precedence amongst pairs of packets |
|---|
| • Earliest deadline first (EDF) |
| • Same deadlines, order lowest window-constraint first |
| • Equal deadlines and zero window-constraints, order highest window-denominator first |
| • Equal deadlines and equal non-zero window-constraints, order lowest window-numerator first |
| • All other cases: first-come-first-serve |

Rich West (2001)

# EDF versus DWCS

| $s_1$ | $s_2$ | $s_3$ | $s_1$ | $s_2$ | $s_3$ | $s_1$ | $s_2$ | $s_3$ | $s_1$ | $s_2$ | $s_3$ | $s_1$ | $s_2$ | $s_3$ | $s_1$ | EDF |

| $s_1$ | $s_2$ | $s_1$ | $s_3$ | $s_1$ | $s_2$ | $s_1$ | $s_3$ | $s_1$ | $s_2$ | $s_1$ | $s_3$ | $s_1$ | $s_2$ | $s_1$ | $s_3$ | DWCS |

time, t

0  1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16

$s_1$    1/2(1),1/1(2),1/2(3),1/1(4),1/2(5)...

$s_2$    3/4(1),2/3(2),2/2(3),1/1(4),3/4(5),2/3(6),2/2(7),1/1(8),3/4(9)...

$s_3$    6/8(1),5/7(2),4/6(3),3/5(4),3/4(5),2/3(6),1/2(7),0/1(8),6/8(9)...

Rich West (2001)

# DWCS Delay Characteristics

- If feasible schedule, max delay of service to $S_i$ is:
  - $(x_i + 1)T_i - C_i$
  - Note: Every time $S_i$ is not serviced for $T_i$ time units $x_i'$ is decremented by 1 until it reaches 0.

- If no feasible schedule, max delay of service to $S_i$ is still bounded.
- Function of time to have:
  - Earliest deadline, lowest window-constraint, highest window-denominator.

# Bandwidth Utilization

- Minimum utilization factor of stream $S_i$ is:

$$U_i = \frac{(y_i - x_i)C_i}{y_i T_i}$$

  - i.e., min req'rd fraction of bandwidth.
- **Least upper bound** on utilization is min of utilization factors for all streams that fully utilize bandwidth.
  - i.e., guarantees a feasible schedule.
  - L.U.B. is 100% in a slotted-time system.

Rich West (2001)

# Scheduling Test

- If:

$$\sum_{i=1}^{n} \frac{(1 - \frac{x_i}{y_i}).C_i}{T_i} \leq 1.0$$

and $C_i = K$, $T_i = qK$ for all $i$, where $q$ is 1,2,…etc, then a feasible schedule exists.

- For variable length packets:
  - let $C_i <= K$ for all $i$ or fragment/combine packets & translate service constraints.
    - e.g., ATM SAR layer.

Rich West (2001)

# Simulation Scenario

- 8 classes of packet streams:

| $W_i$ | 1/10 | 1/20 | 1/30 | 1/40 | 1/50 | 1/60 | 1/70 | 1/80 |
|-------|------|------|------|------|------|------|------|------|
| $T_i$ | 400  | 400  | 480  | 480  | 560  | 560  | 640  | 640  |

- Varied number of streams **n**, uniformly distributed amongst traffic classes.
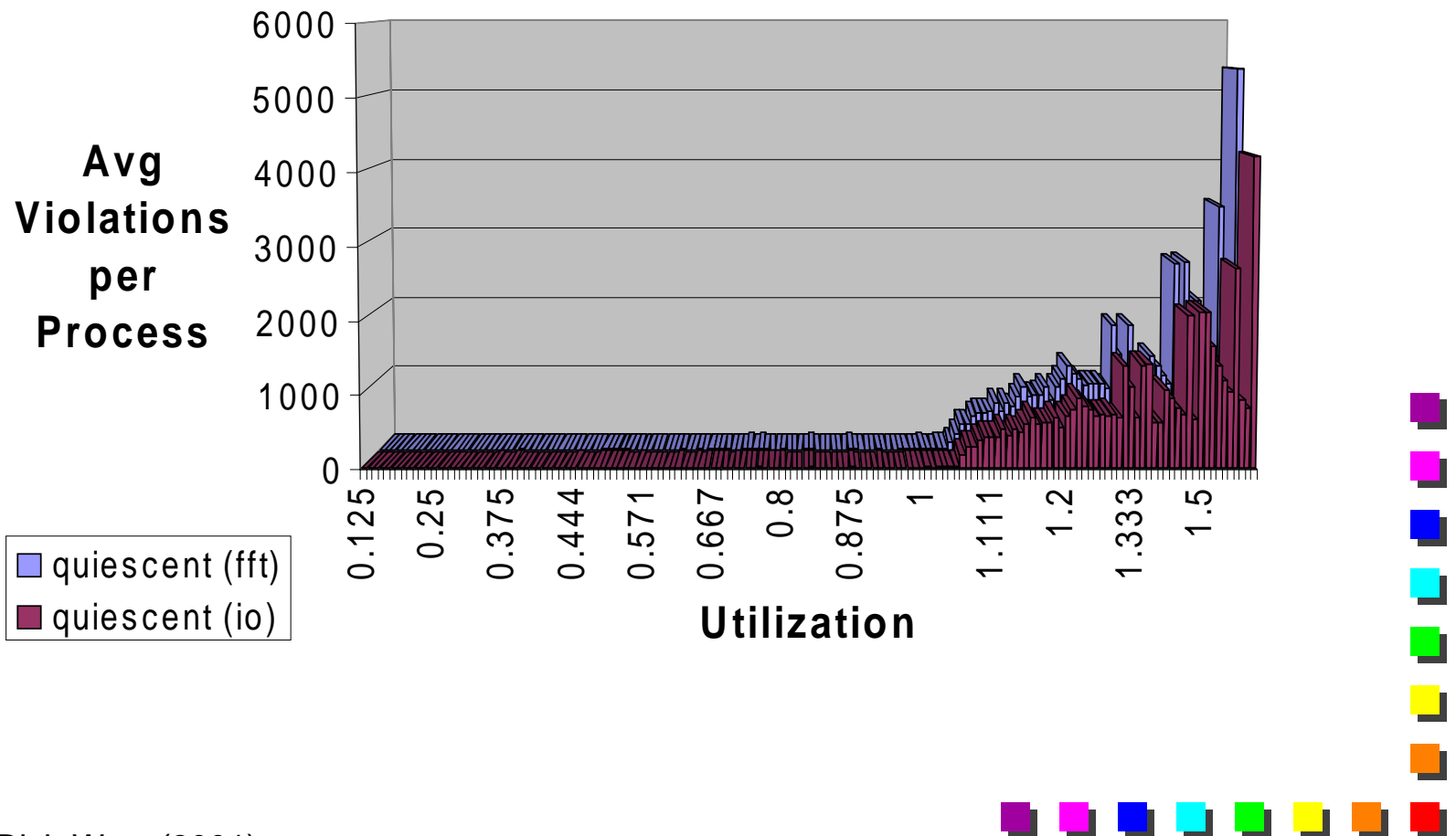- Total of a million packets serviced.
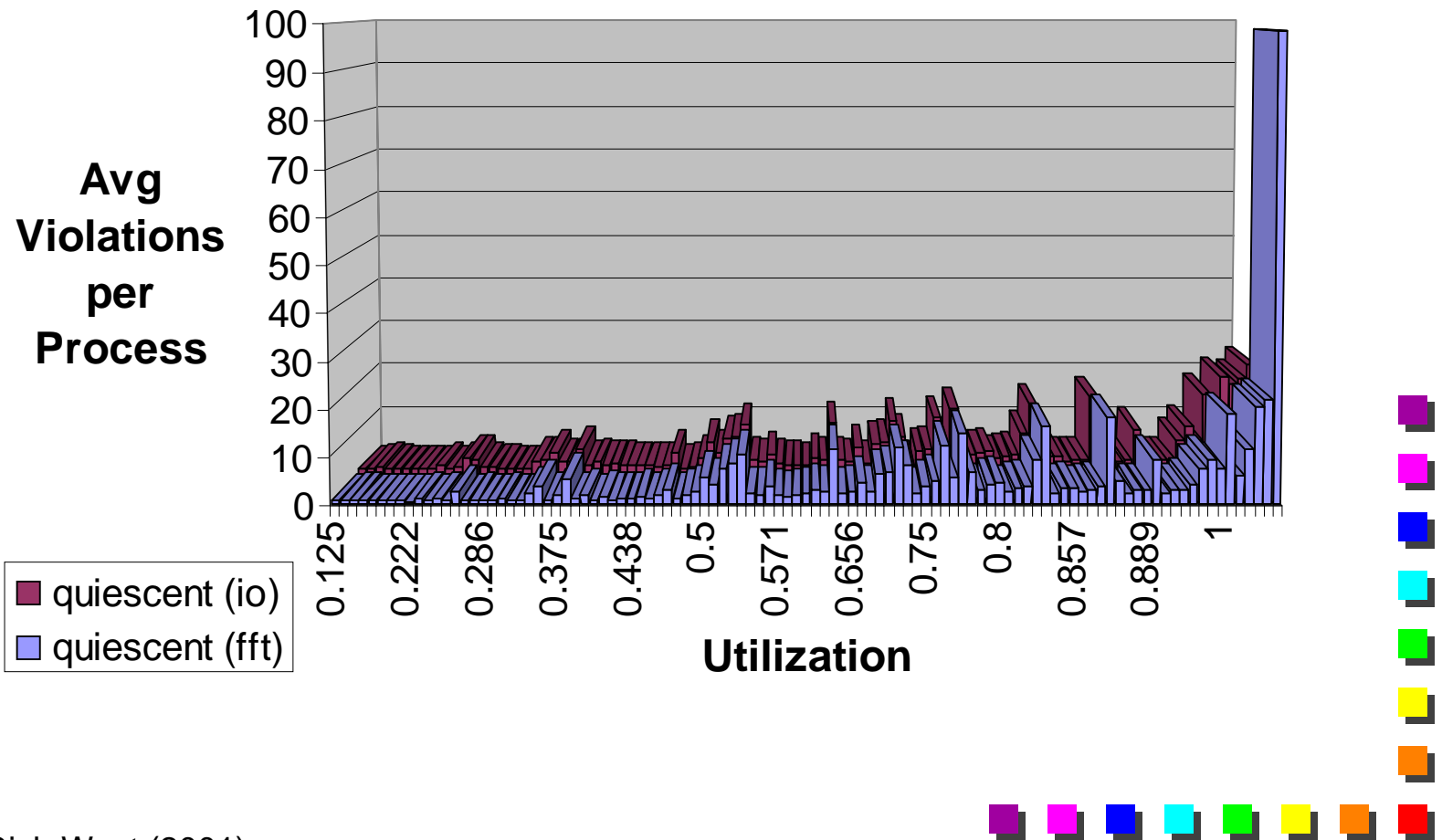
# Bandwidth Utilization Results

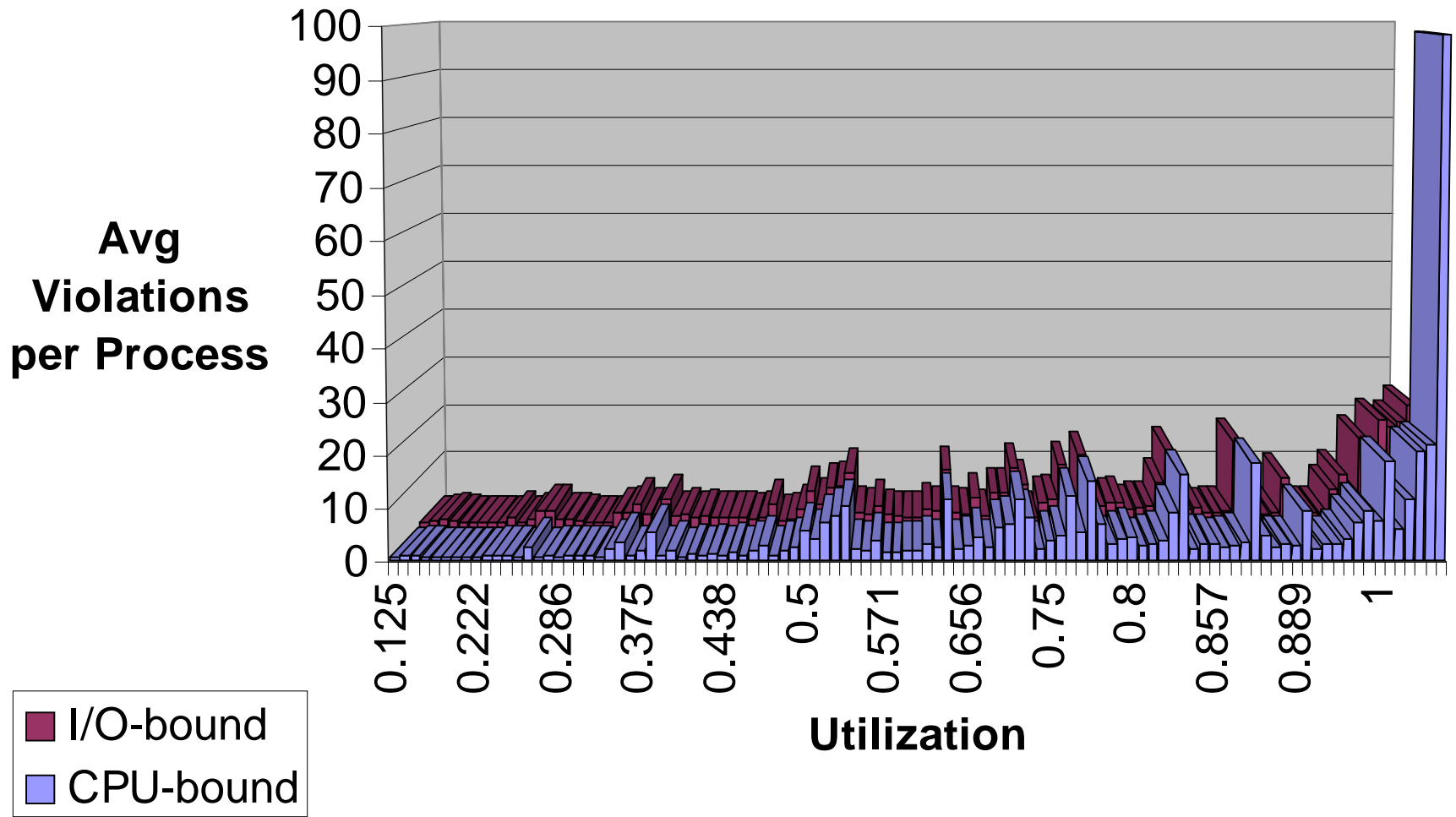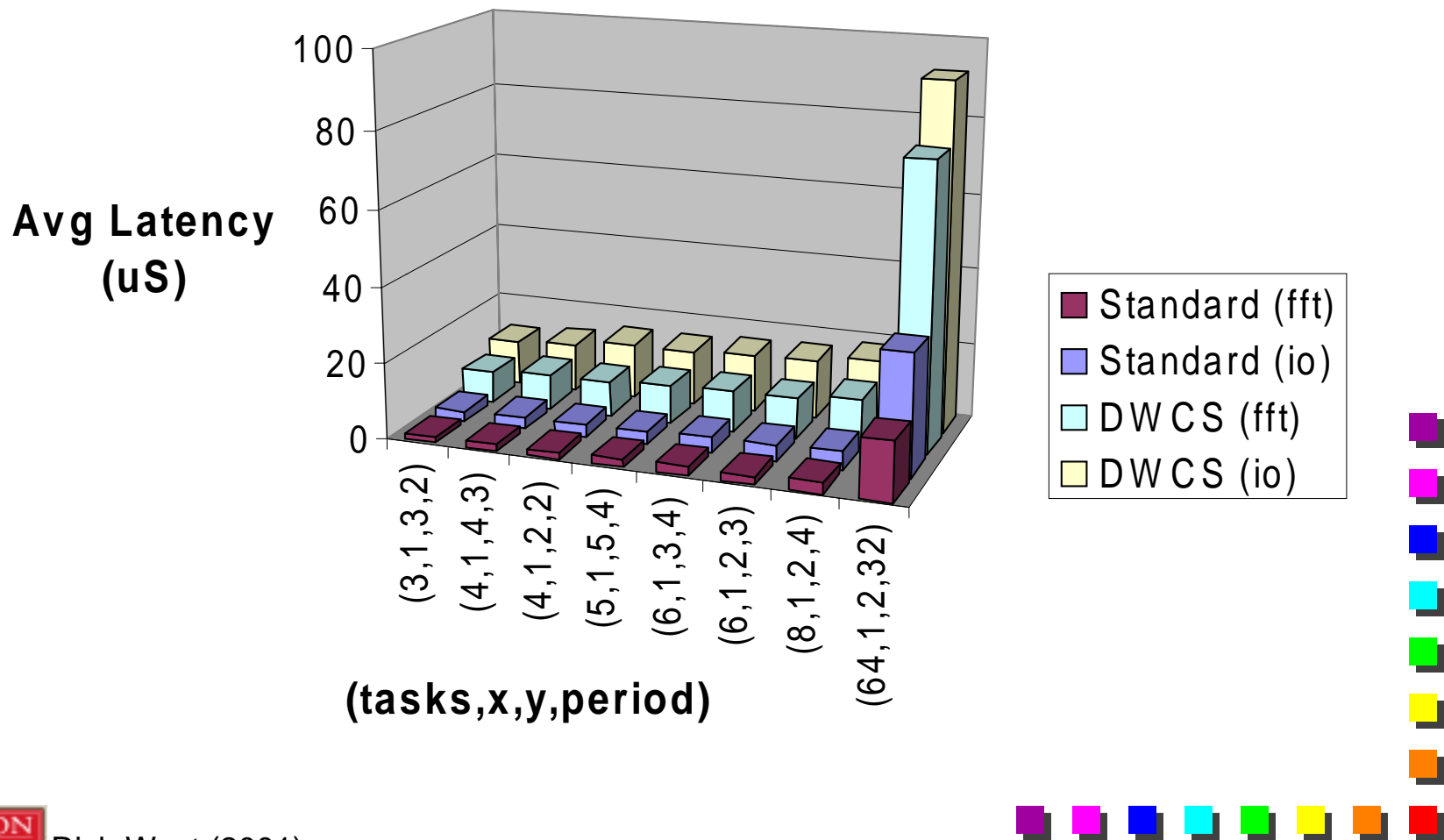| n | D | V | U | $\frac{n}{8} \cdot \sum_{i=1}^{8} \frac{c_i}{T_i}$ |
|---|---|---|---|---|
| 480 | 0 | 0 | 0.9156 | 0.9518 |
| 496 | 0 | 0 | 0.9461 | 0.9835 |
| 504 | 0 | 0 | 0.9613 | 0.9994 |
| 512 | 15152 | 0 | 0.9766 | 1.0152 |
| 520 | 30990 | 0 | 0.9919 | 1.0311 |
| 528 | 46828 | 7038 | 1.0071 | 1.047 |
| 544 | 78528 | 31873 | 1.0376 | 1.0787 |
| 560 | 110240 | 53455 | 1.0681 | 1.1104 |
| 640 | 268800 | 148143 | 1.2207 | 1.269 |

Rich West (2001)

# (x,y)-hard Linux CPU DWCS: Average Violations per Process

# (x,y)-hard Linux CPU DWCS: Average Violations per Process



Rich West (2001)

# (x,y)-hard Linux CPU DWCS: Scheduling Latency



Avg Latency (uS)

(tasks,x,y,period)

Legend:
- Standard (fft)
- Standard (io)
- DWCS (fft)
- DWCS (io)

x-axis labels: (3,1,3,2) (4,1,4,3) (4,1,2,2) (5,1,5,4) (6,1,3,4) (6,1,2,3) (8,1,2,4) (64,1,2,32)
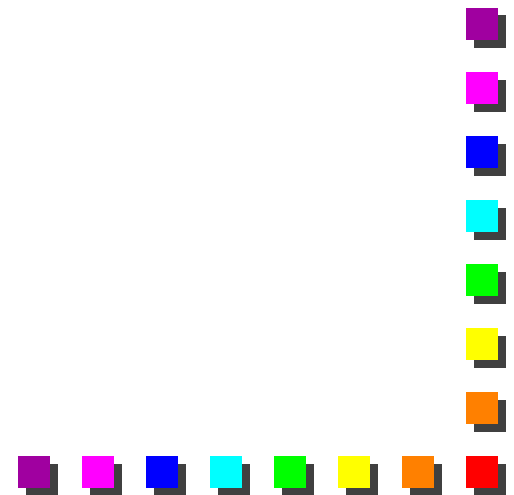
Rich West (2001)

# (x,y)-hard Linux CPU DWCS: % Execution Time in Violation


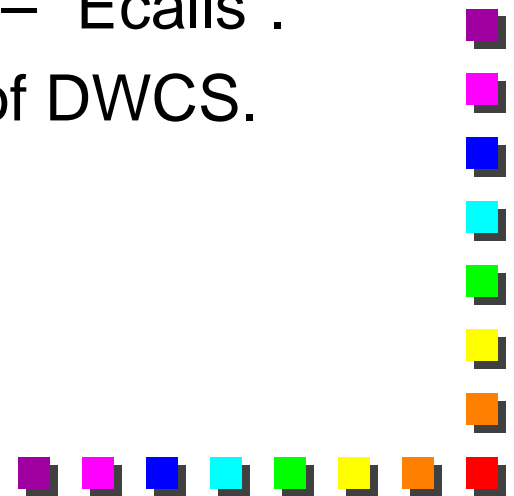
Rich West (2001)

# Conclusions

- **Flexible** approach to run-time service adaptation.
  - <u>When</u>, <u>where</u> and <u>how</u> to adapt.
- **Coordinated** resource management.
  - Dionisys "quality events", monitors, handlers etc.
- **DWCS** guarantees explicit loss and delay constraints for real-time / multimedia applications.

# Current & Future Work

- Linux kernel-level implementation of Dionisys mechanisms.
  - Cluster-wide coordination of resources.
  - Language support for "QoS safety".
    - Stability analysis.
  - Real-time "batched" events in Linux – "Ecalls".
- Switch / co-processor implementation of DWCS.
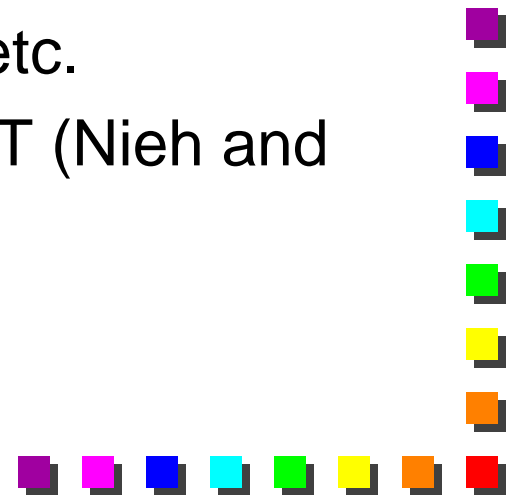  - Scheduling variable-length packets.

# Related Work

- **QoS Architectures**: QoS-A (Campbell), Washington Univ. (Gopalakrishna & Parulkar), QoS Broker (Nahrstedt et al), U. Michigan (Abdelzaher, Shin), QuO (BBN) + more…

- **QoS Specification/Translation**: Tenet (Ferrari), EPIQ (Illinois).

- **QoS Evaluation**: Rewards (Abdelzaher), Value fns (Jensen), Payoffs (Kravets).

- **System Service Extensions**: SPIN (U. Washington), Exokernel (MIT).

# Scheduling Related Work

- **Fair Scheduling**: WFQ/WF$^2$Q (Shenker, Keshav, Bennett, Zhang etc), SFQ (Goyal et al), EEVDF/Proportional Share (Stoica, Jeffay et al).

- **(m,k) Deadline Scheduling**: Distance-Based Priority (Hamdaoui & Ramanathan), Dual-Priority Scheduling (Bernat & Burns), Skip-Over (Koren & Shasha).

- **Pinwheel Scheduling**: Holte, Baruah etc.

- **Other multimedia scheduling**: SMART (Nieh and Lam).

Rich West (2001)

# Related Research Papers

- **Quality Events: A Flexible Mechanism for Quality of Service Management**, *RTAS 2001.*

- **Analysis of a Window-Constrained Scheduler for Real-Time and Best-Effort Traffic Streams**, *RTSS 2000.*

- **Dynamic Window-Constrained Scheduling for Multimedia Applications**, *ICMCS'99.*

- **Scalable Scheduling Support for Loss and Delay-Constrained Media Streams**, *RTAS'99.*

- **Exploiting Temporal and Spatial Constraints on Distributed Shared Objects**, *ICDCS'97.*

Rich West (2001)