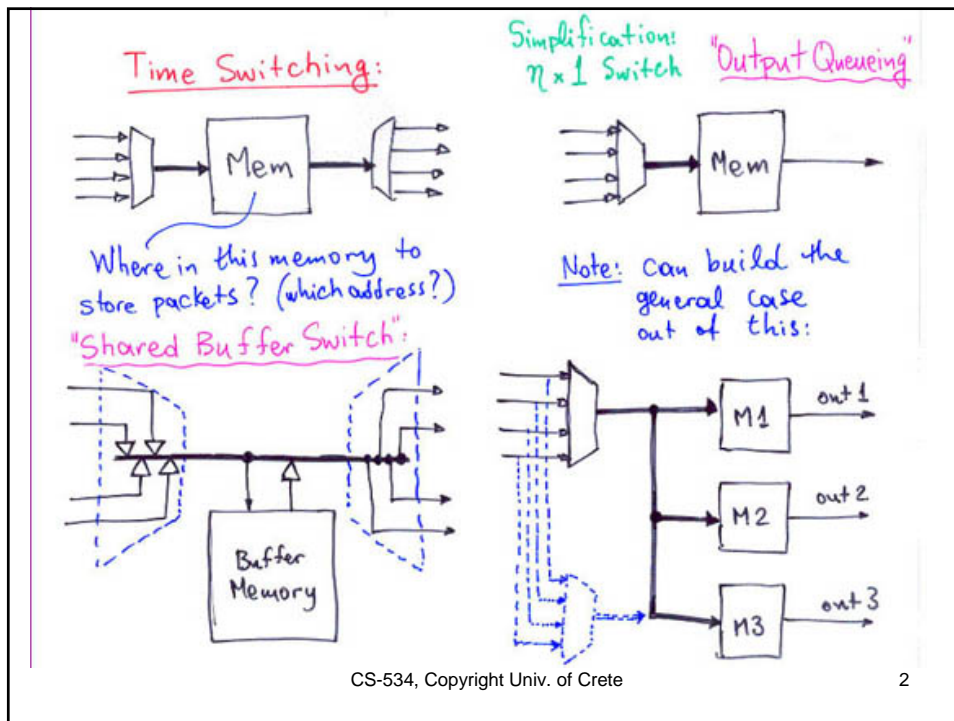


### 3.1 The need for multiple queues within a same buffer memory: Multi-Queue Data Structures

- Buffers memories for time switching: what data structures?
- Single queue feeding multiple destinations/classes
  - ⇒ Head-of-Line (HOL) Blocking ⇒ poor performance
- Multiple Queues in one buffer:
  - Partitioned Space (underutilized) ⇒ circular queues
  - Shared Space (efficient) ⇒ linked-list queues

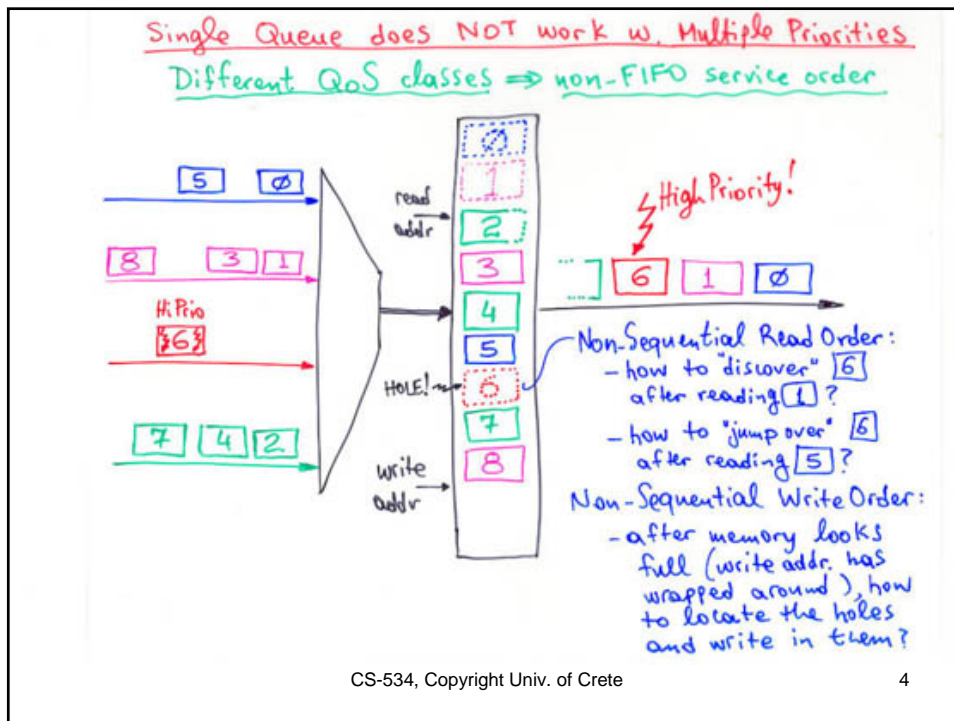
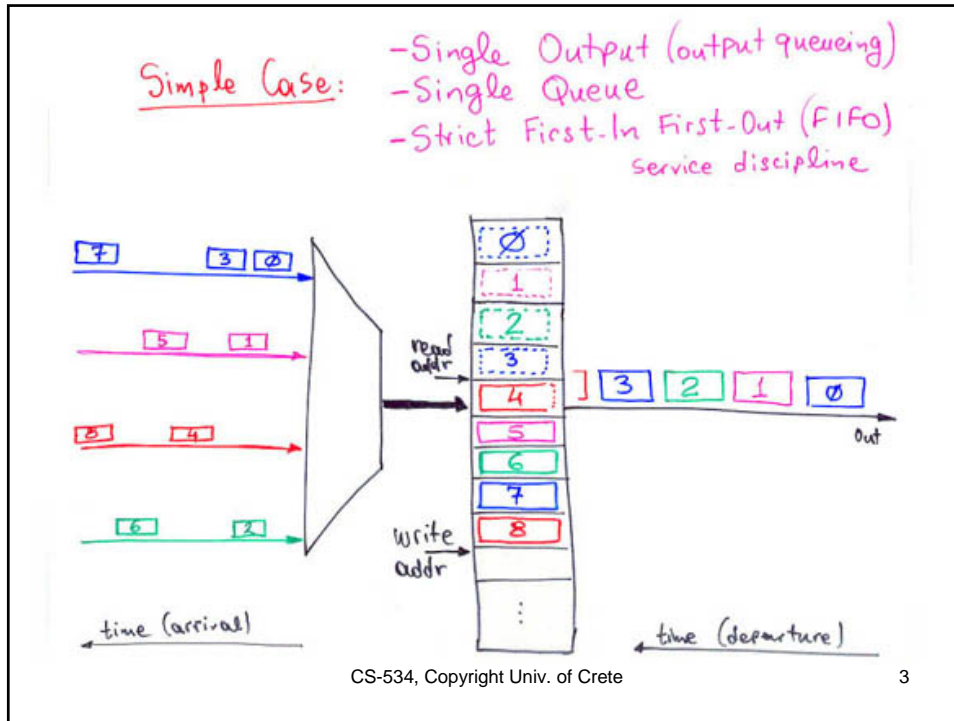
CS-534, Copyright Univ. of Crete

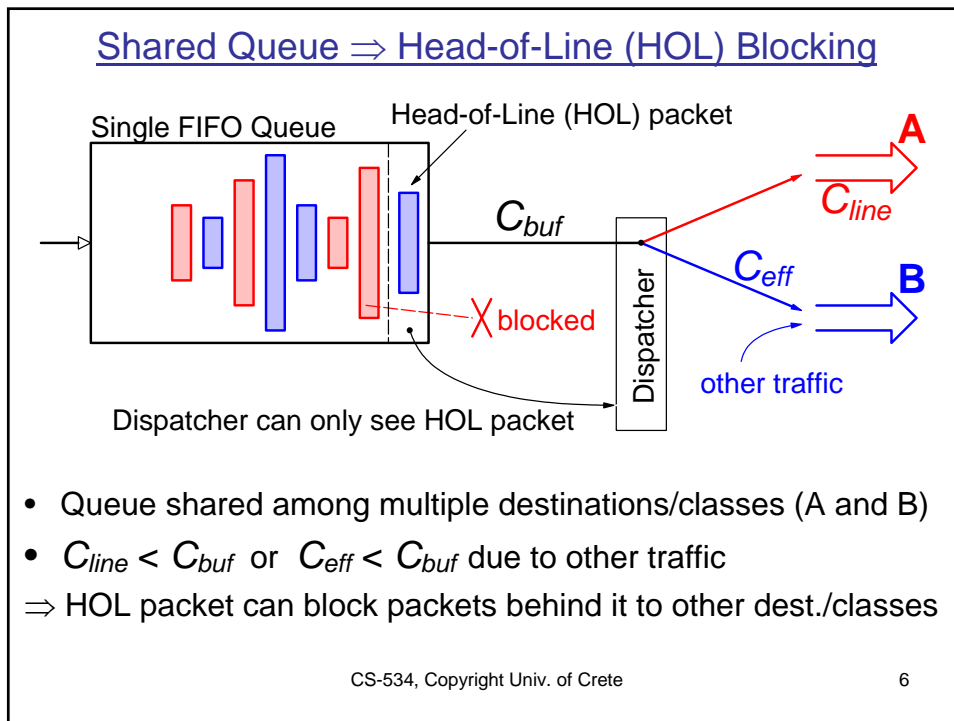
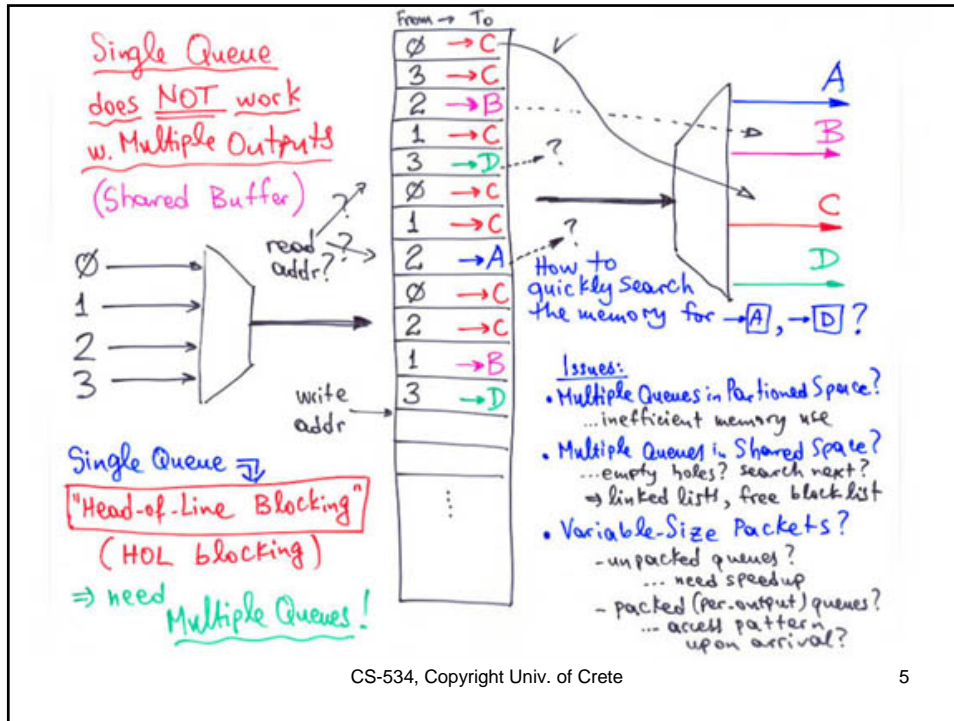
1

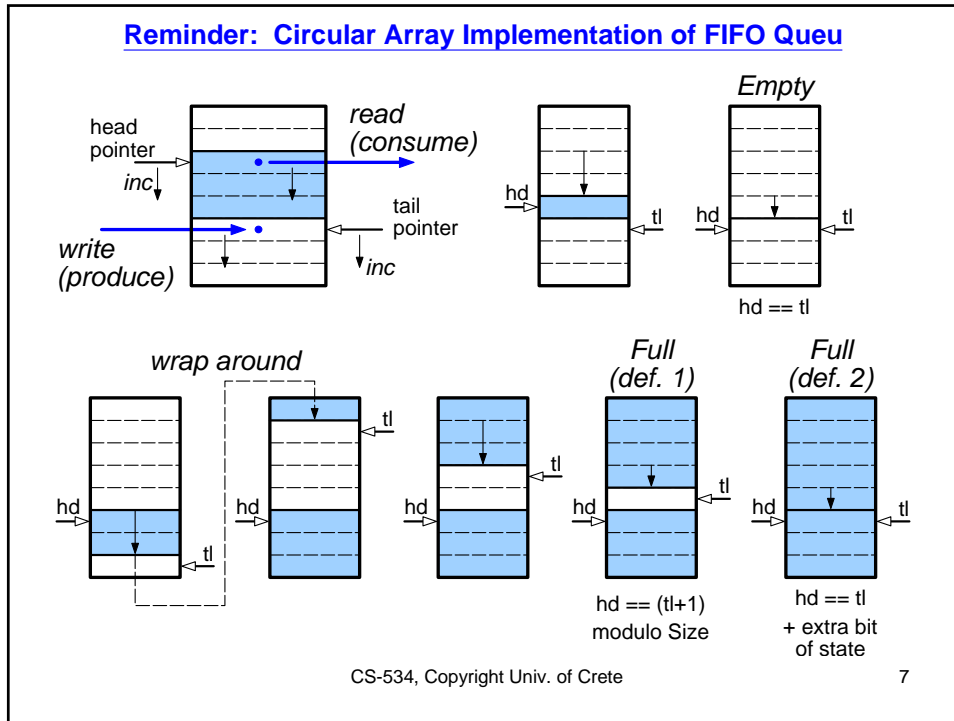


CS-534, Copyright Univ. of Crete

2







### B.2 Multiple FIFO Queues with Statically Partitioned Space

for each of them:

- One RAM for all queues;  
Two counters/pointers per queue.
- Limits between queue areas (partitions) can be "hardwired", or "configurable" off-line (when queues are empty); limit registers needed in the latter case
- Underutilization of memory space (one queue may overflow, while space exists on others)
- Variable-size packets OK.

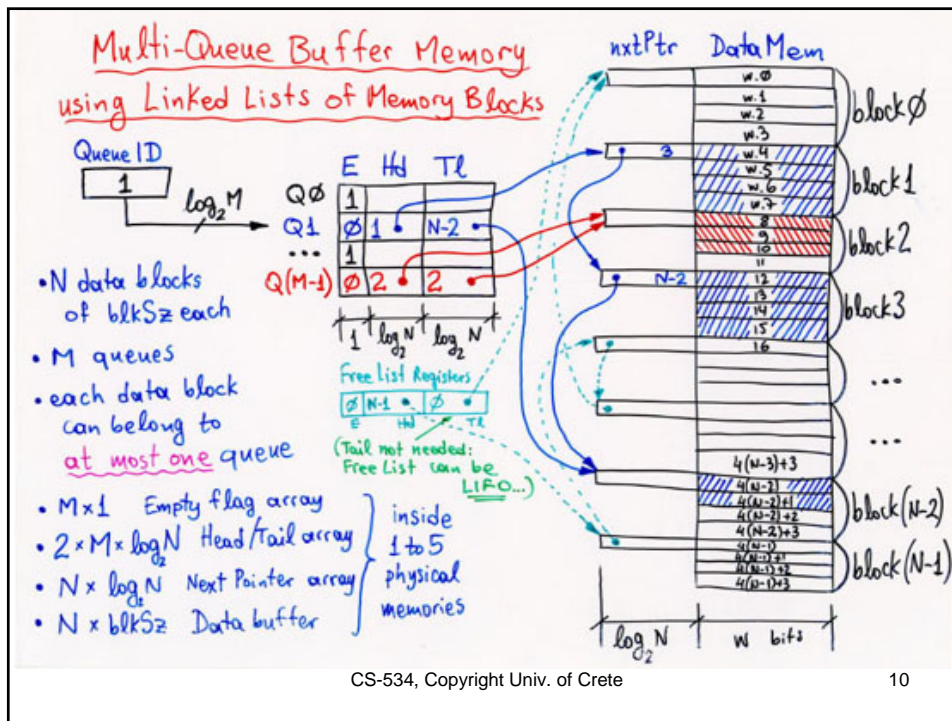
The diagram shows three statically partitioned FIFO queues (Q1, Q2, Q3) in a shared RAM. Each queue has its own head (h1, h2, h3) and tail (t1, t2, t3) pointers. The queues are separated by limits q1 and q2. The queues are shown as shaded regions within a larger container.

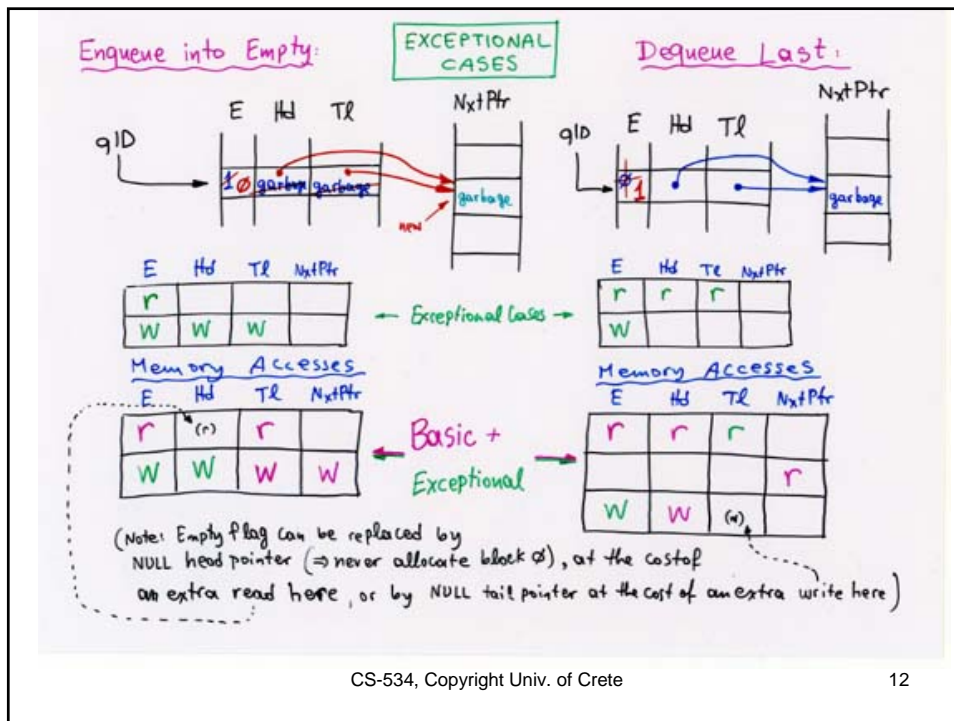
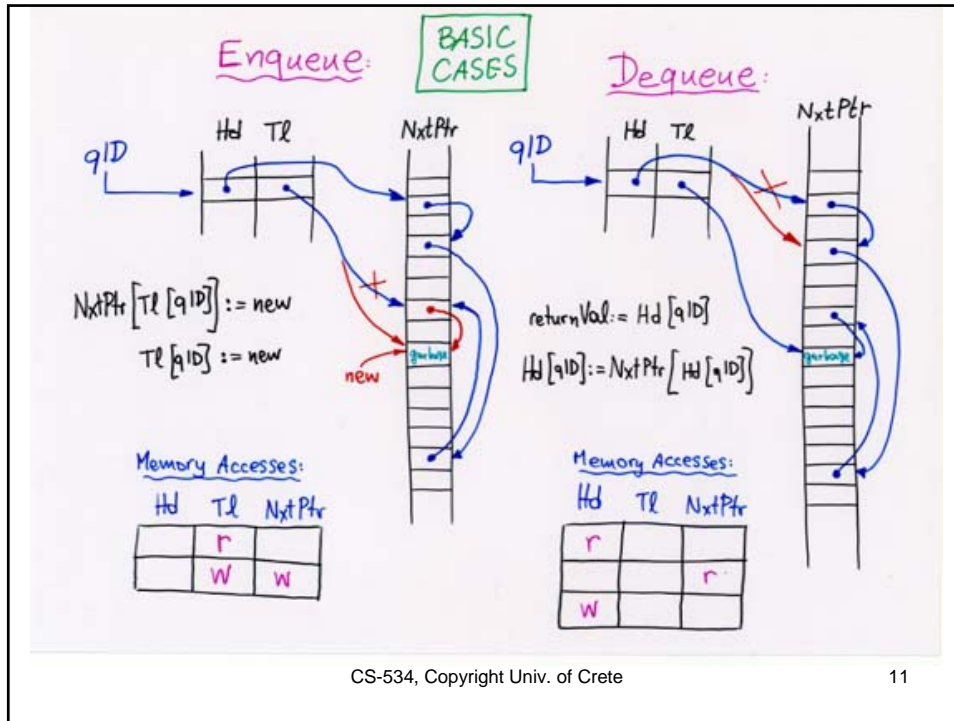
534  
12.6

CS-534, Copyright Univ. of Crete 8

### Multiple Queues Sharing a Buffer Memory Space

- Often need many queues where most have low occupancy (or empty) but few have high occupancy
  - Departures create holes  $\Rightarrow$  empty space is fragmented  $\Rightarrow$  neighboring packets physically non-contiguous  $\Rightarrow$  linked-list data structure
- Reference: Y. Tamir, G. Frazier: "High Performance Multi-Queue Buffers for VLSI Communication Switches", ISCA 1988*
- "malloc( )" works on fixed-size blocks
    - block size is a tradeoff between fragmentation cost (for very large blocks) and memory address rate (for very small blocks): see exercise set 4
    - ... unless multiple packet fragments in a same block - ex. 6.1





Cost-Performance Tradeoffs:

Assumptions:

- "off-chip" means enq/deq controller FSM & memories are located on separate chips
- "on-chip" means all together
- 1 memacc./cycle/port
- off-chip dependent accesses (read, then use data as next address) cost one extra cycle of latency between them

Note:

- \* enqueue latency may increase when mixing enqueue and dequeue operations in the same pipeline.

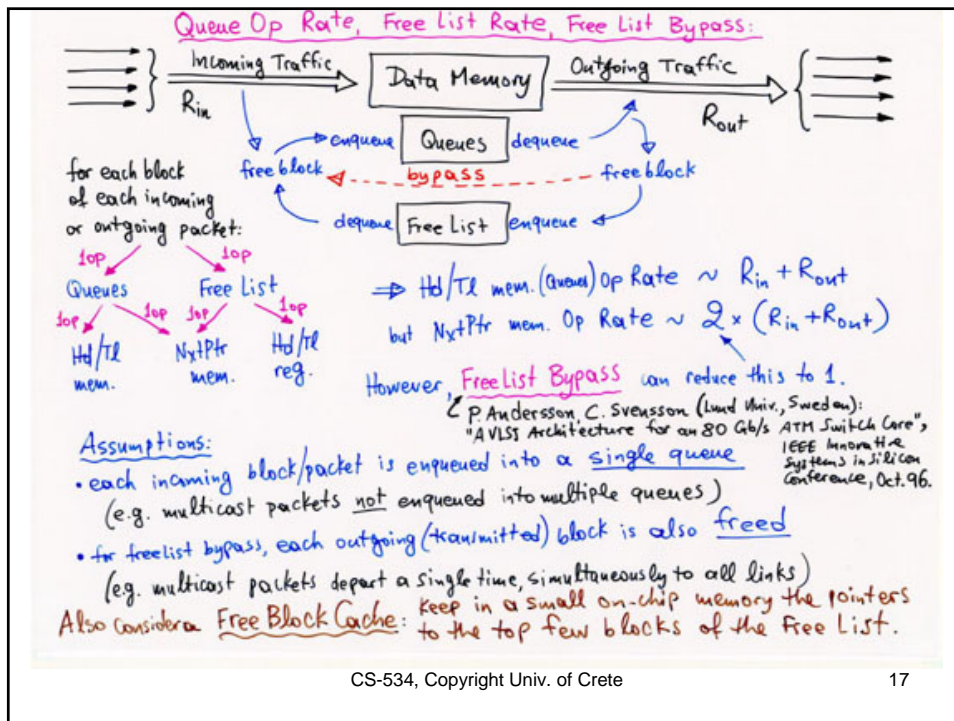
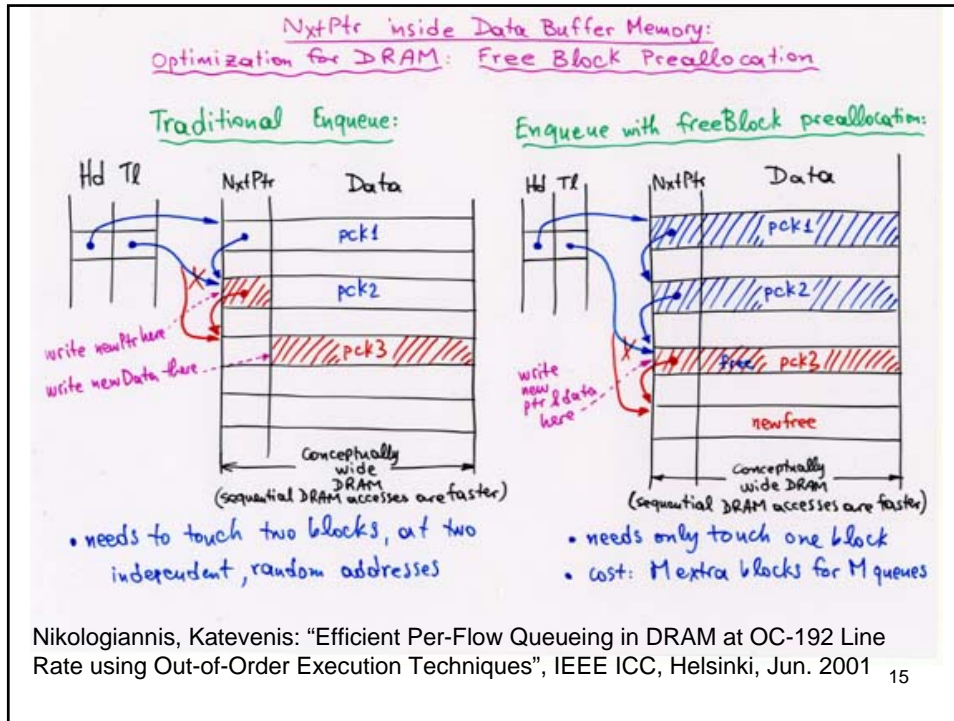
E	Hd	TL	NxtPtr	Enqueue		Dequeue	
				Latency (clock cycles)	Thruput (op's/cycle)	Latency (clock cycles)	Thruput (op's/cycle)
				4 or 5	$\frac{1}{4 \text{ or } 5}$	5	$\frac{1}{5 \text{ or } 4}$
				in a single, 1-port off-chip memory			
1-port on-chip				4	$\frac{1}{4}$	5	$\frac{1}{4}$
				in a single, 1-port off-chip memory			
1-port on-chip				3*	$\frac{1}{3}$	5	$\frac{1}{3}$
				1-port, off-chip log <sub>2</sub> N wide			
1-port on-chip				3*	$\frac{1}{2}$	5	$\frac{1}{2}$
				1-port off-chip			
1-port on-chip				2*	$\frac{1}{2}$	3	$\frac{1}{2}$
				1-port off-chip			
2-port on-chip				2*	1	3	1
				2-port on-chip			
				2-port on-chip			
				1-port off-chip			

CS-534, Copyright Univ. of Crete 13

### Notes on Enqueue/Dequeue Performance Table

- Table assumes we always need to perform both normal and exceptional accesses – see ex. 5.1
- Table assumes fall-through timing for off-chip SRAM: one-cycle latency per access, plus one extra cycle between dependent off-chip accesses – see ex. 5.2
- For peak throughput: overlap successive operations (latency of individual operations increases)

*For a highly pipelined implementation, refer to: Kornaros, Kozyrakis, Vatsolaki, Kateveni: "Pipelined Multi-Queue Management in a VLSI ATM Switch Chip with Credit-Based Flow Control", 17<sup>th</sup> Conf. on Advanced Research in VLSI, 1997*





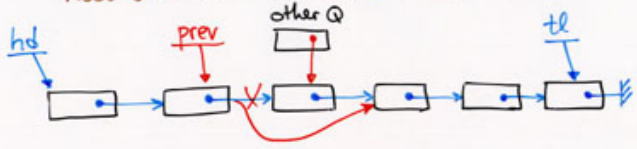
Dropping a Segment (= move to free list)  
or Moving a Segment to a Different Queue

① From the Head of the Source Queue:

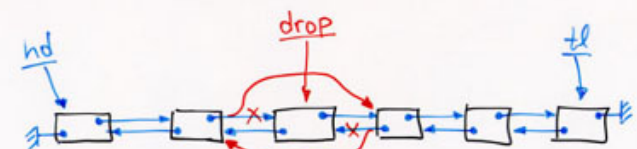
(a) dequeue from source queue,  
(b) enqueue into destination Q.

② From the "Middle" of the Source Queue:

- Need a Pointer to the Previous Segment:



- else • Need Doubly Linked List:



CS-534, Copyright Univ. of Crete 18

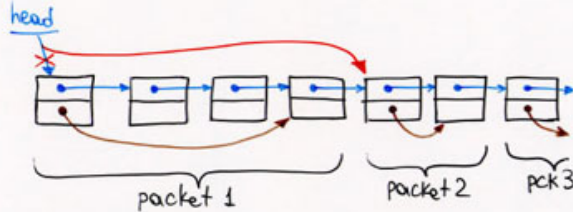
Dropping Entire Packet (multi-segment) or moving to other Queue:

① In  $O(\text{pktSize})$  Time:

- Drop or move one segment at a time...
- ⇒ takes time proportional to the size of the packet
- Usually, this consumes the same resources (time, memory ports) as if the packet were transmitted to the output port

② In  $O(1)$  Time:

- Need two next-pointers per segment:



(the utilization of the second set of pointers is only:  $\frac{1}{\text{average number of segments per packet}}$ )

→ alternative: separate "packet descriptor" data structures, like for multicast segment descriptors - see below).

CS-534, Copyright Univ. of Crete 19

## Queueing for Multicast Traffic

- Multicast traffic is expected to become very important in the future –but so has it been for many years in the past...
- Supporting multicast traffic usually increases complexity and cost
- Queueing for Multicast Traffic:
  - Each segment (block) allowed in only one queue  $\Rightarrow$  HOL blocking
  - Each segment allowed in multiple queues  $\Rightarrow$  need many nxtPtr's
  - Enqueue throughput and nxtPtr space: static vs. dynamic sharing
- References:
  - F. Chiussi, Y. Xia, V. Kumar: "Performance of Shared-Memory Switches under Multicast Bursty Traffic", IEEE Jour. Sel. Areas in Communications (JSAC), vol. 15, no. 3, April 1997, pp. 473-487.
  - D. Stiliadis: "Efficient Multicast Algorithms for High-Speed Routers", Proc. IEEE Workshop on High Performance Switching and Routing (HPSR 2003), Torino, Italy, June 2003, pp. 117-122.

Multicast Traffic: Same or Different Queues with Unicast Traffic?

Case 1: Each segment is only allowed to belong to a single queue:

per-output unicast queues

Single Multicast Queue

• impractical to have per-output (=combination) multicast queues: would need  $O(2^n)$  queues!

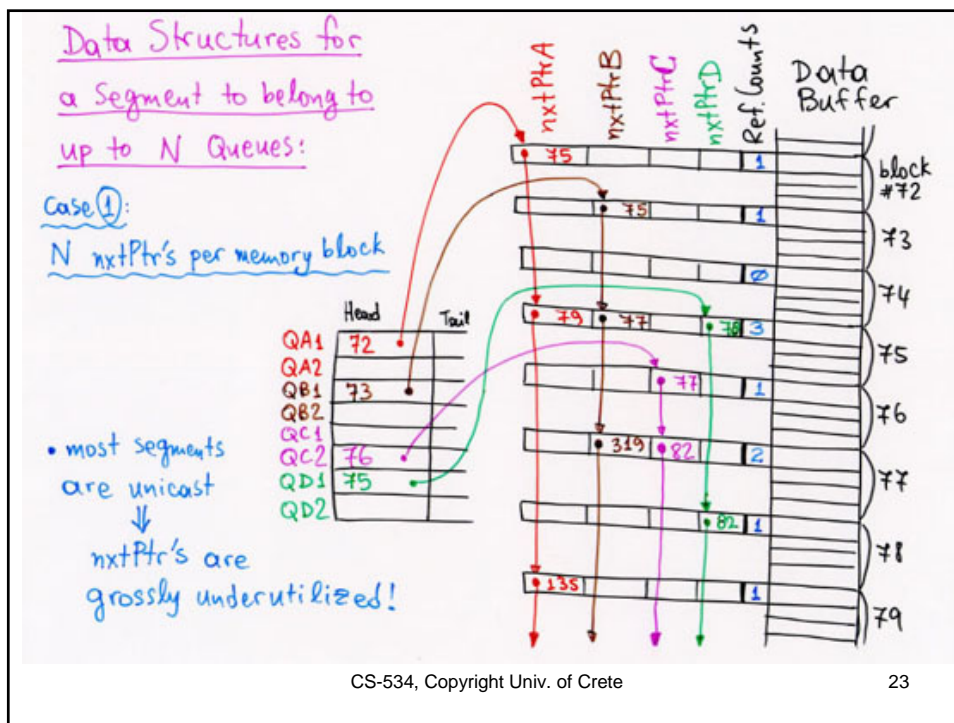
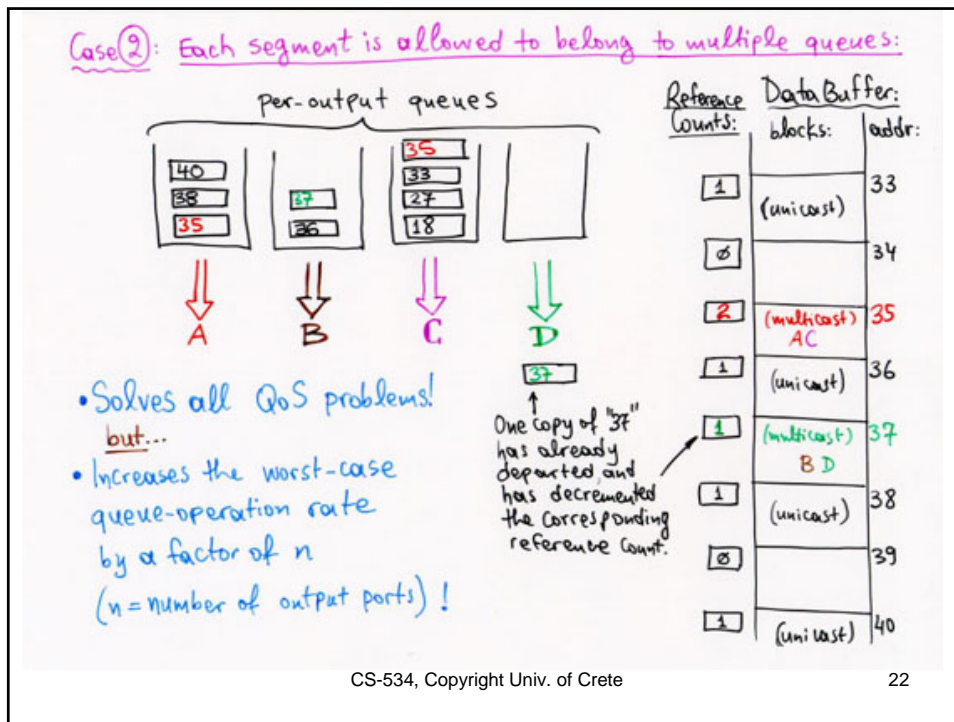
Head-of-Line Blocking!

Other Priority Level Queues

- we may have a different set of these queues (including multicast) per priority level, but it may still happen that traffic destined to outputs A and C currently exists at priority levels higher than "our" cell A C while all queues destined to B and D at priority levels above "our" cell B D are empty...

CS-534, Copyright Univ. of Crete

21



Case(2): Decouple Linked List Nodes from Data Buffer Addresses

- twice the cost per nextPtr (need a segPtr as well now)
- but...
- much fewer than N\*S descriptors (based on average ratio of unicast-to-multicast segments, and average fan-out of multicast segments - e.g. f=2)

Optimization:  
Partition the address space of queue member descriptors into two parts:

- 0 to S-1: unicast-only segments, no segPtr needed (segPtr[i] ≤ i);
- S to fS-1: full queue member descriptors, with nextPtr and segPtr, intended for use by multicast segments.

CS-534, Copyright Univ. of Crete 24

Enqueue Operation Rate for multicast segments onto multiple per-output queues

Dequeue Rate:  $\frac{1+S-U}{f} = \frac{(1-U)+S}{f} \geq \frac{m+S}{f} \geq m \Rightarrow S \geq m \cdot (f-1)$

for the replication buffer not to overflow

References:

- F. Chiussi, Y. Xia, V. Kumar: IEEE JSAC, April 1997, pp. 473-487
- D. Stiliadis: HPSR 2003, pp. 117-122

CS-534, Copyright Univ. of Crete 25