

P51: High Performance Networking

Lecture 5: Low Latency Devices

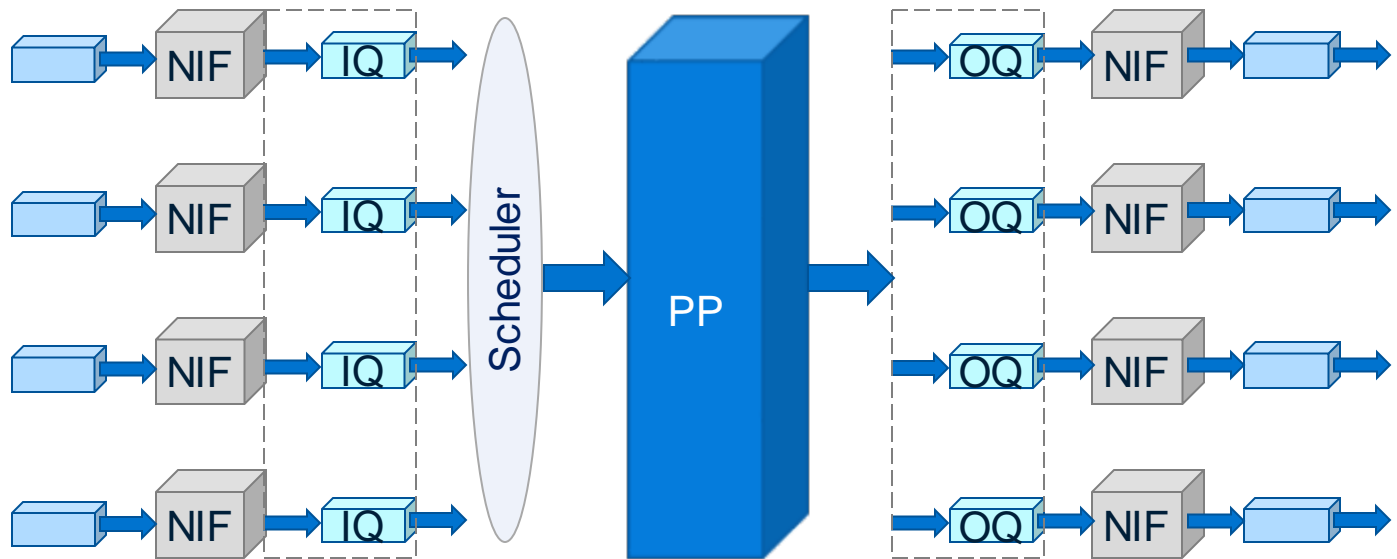
Project – Evaluation Plan

- For the next lab, you should prepare an evaluation plan for your project.
- The following evaluation tests are expected:
 - Functional testing (using the NetFPGA test infrastructure)
 - Performance testing, using synthetic traffic (using OSNT or equivalent to measure latency at low-rate and maximum throughput)
 - Performance comparison to other solutions (software, hardware)

Low Latency Switches

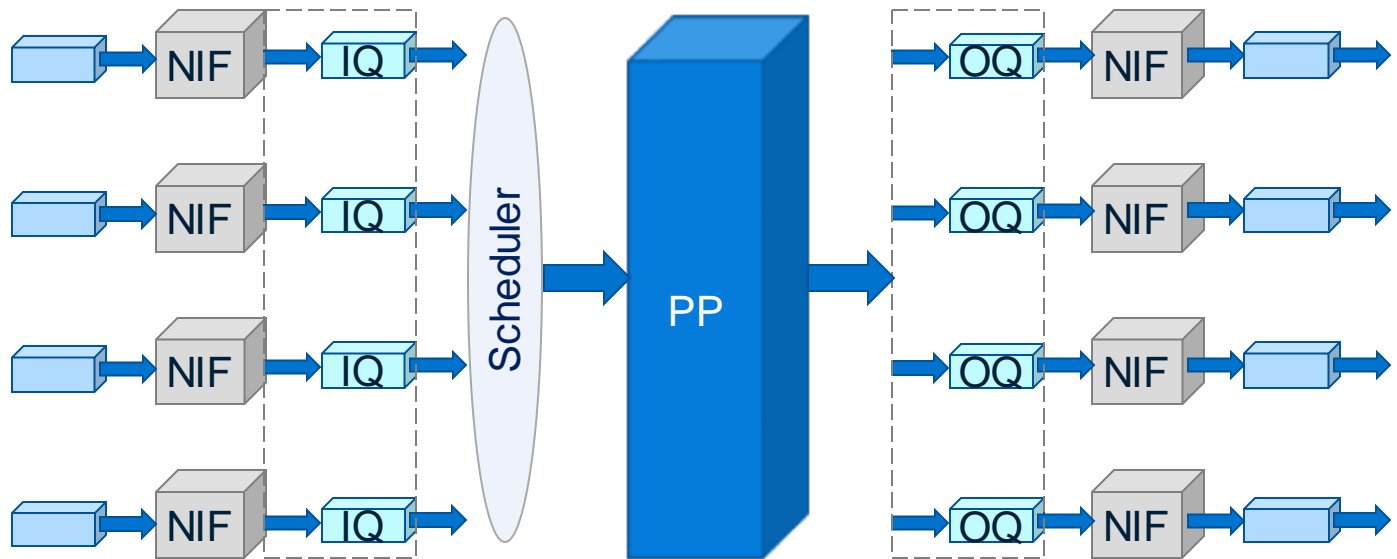
How to lower the latency of a switch?

- Obvious option 1: Increase clock frequency
 - E.g. change core clock frequency from 100MHz to 200MHz
 - Half the time through the pipeline



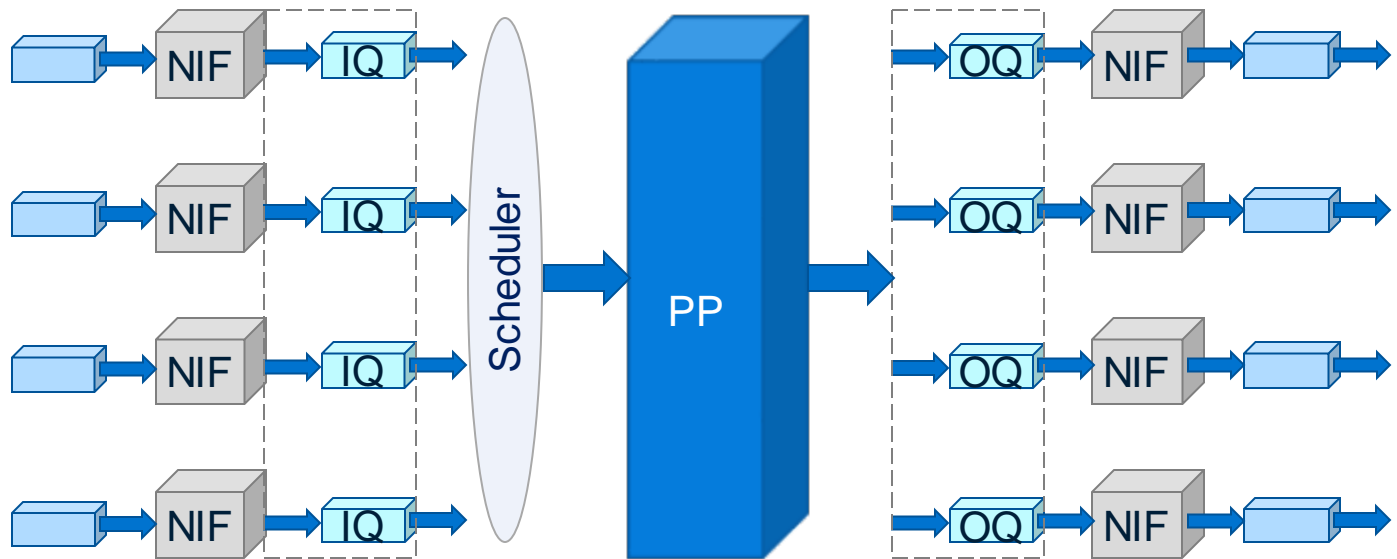
How to lower the latency of a switch?

- Obvious option 1: Increase clock frequency
- Limitations:
 - Frequency is often a property of manufacturing process
 - Some modules (e.g. PCS) must work at a specific frequency (multiplications)



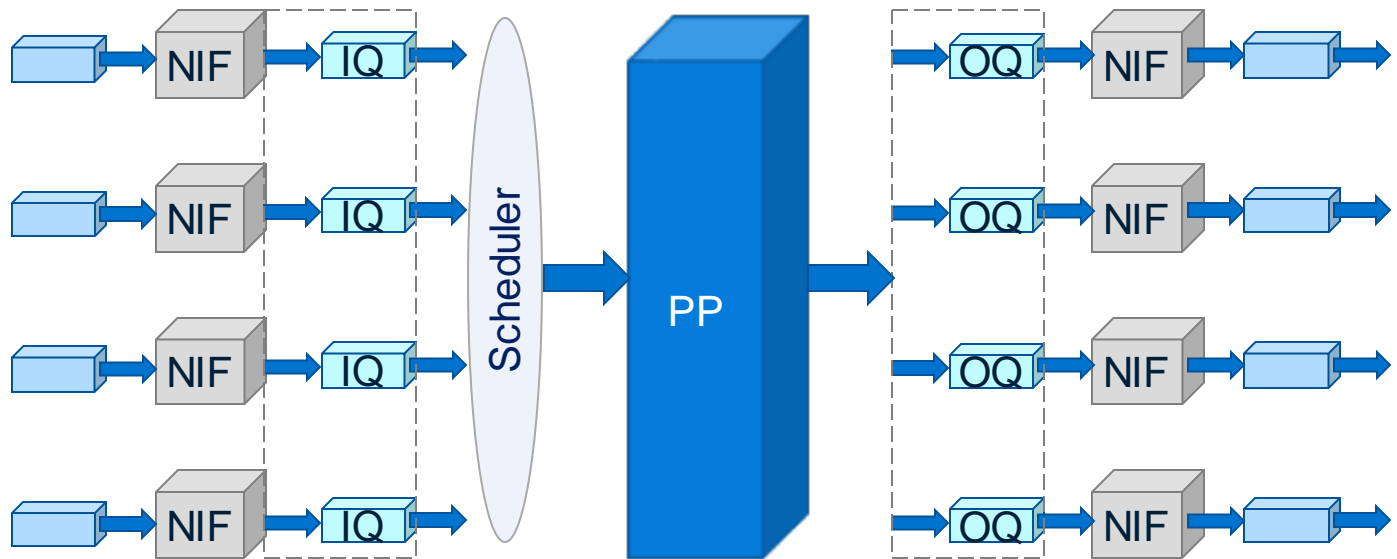
How to lower the latency of a switch?

- Obvious option 2: Reduce the number of pipeline stages
 - Can you do the same in 150 pipeline stages instead of 200?
 - Limitation: hard to achieve.



How to lower the latency of a switch?

- Can we achieve ~ 0 latency switch?
 - Is there a lower bound on switch latency?



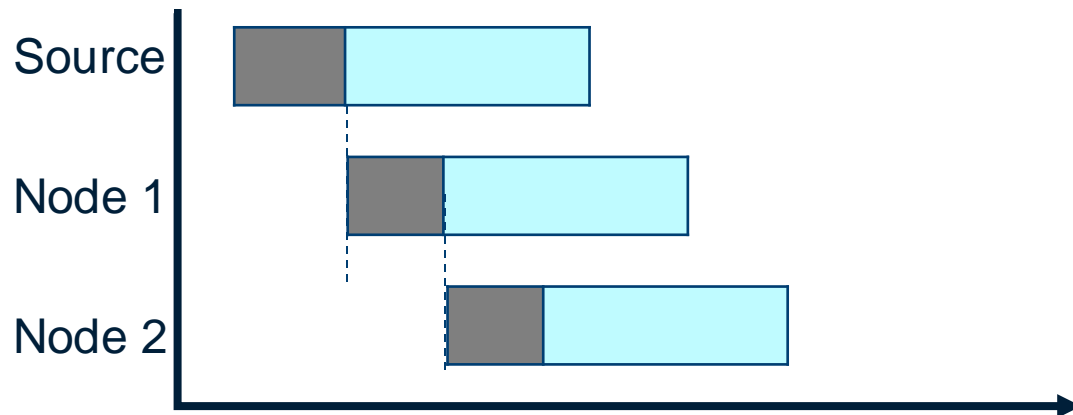
Cut Through Switching

Cut Through Switch

- Cut through switch \neq Low latency switch
 - A cut through switch can implement a very long pipeline...
- But:
 - For the smallest packet, the latency is ~same as longest packet
 - As packet size grows, latency saving grows

What is a cut-through switch?

- Kermani & Kleinrock, “Virtual cut-through: A new computer communication switching technique”, 1976
- “when a message arrives in an intermediate node **and its selected outgoing channel is free (just after the reception of the header)**, then, in contrast to message switching, the message is sent out to the adjacent node towards its destination **before it is received completely** at the node; only if the message is blocked due to a busy output channel is a message buffered in an intermediate node.”

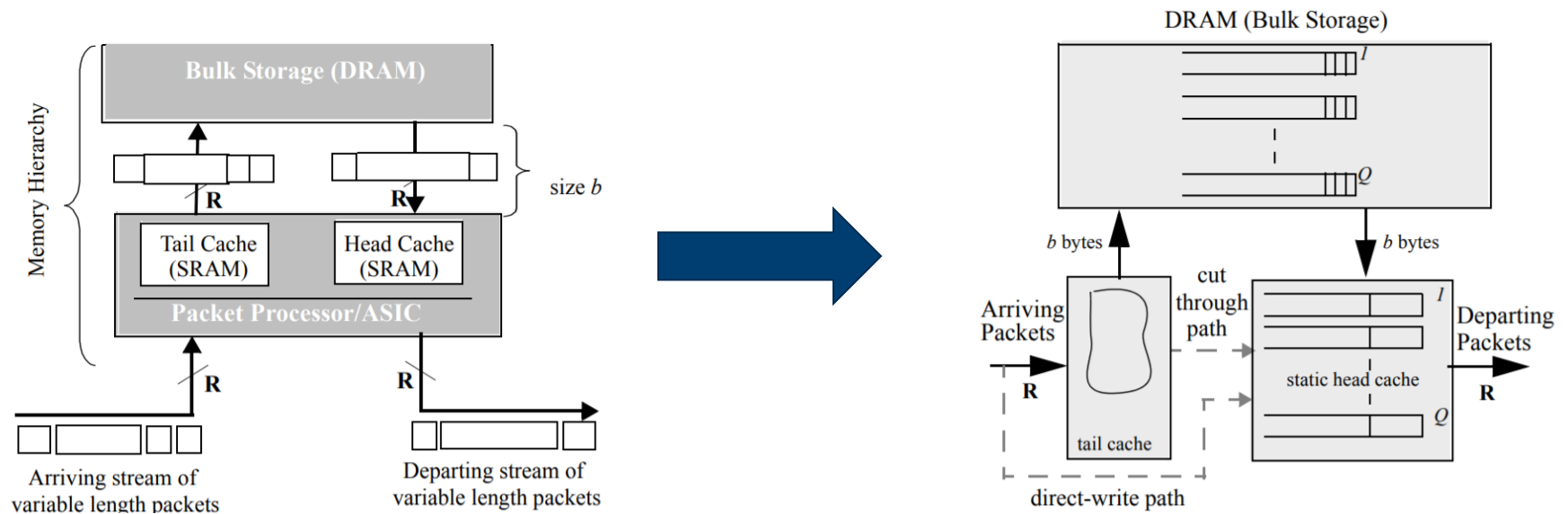


What is a cut-through switch?

- Past (far back):
 - Networks were slow
 - Memory was fast
 - Writing packets to the DRAM took “negligible” time
- With time:
 - Networks became faster
 - Memory access time is no longer “negligible”

What is a cut-through switch?

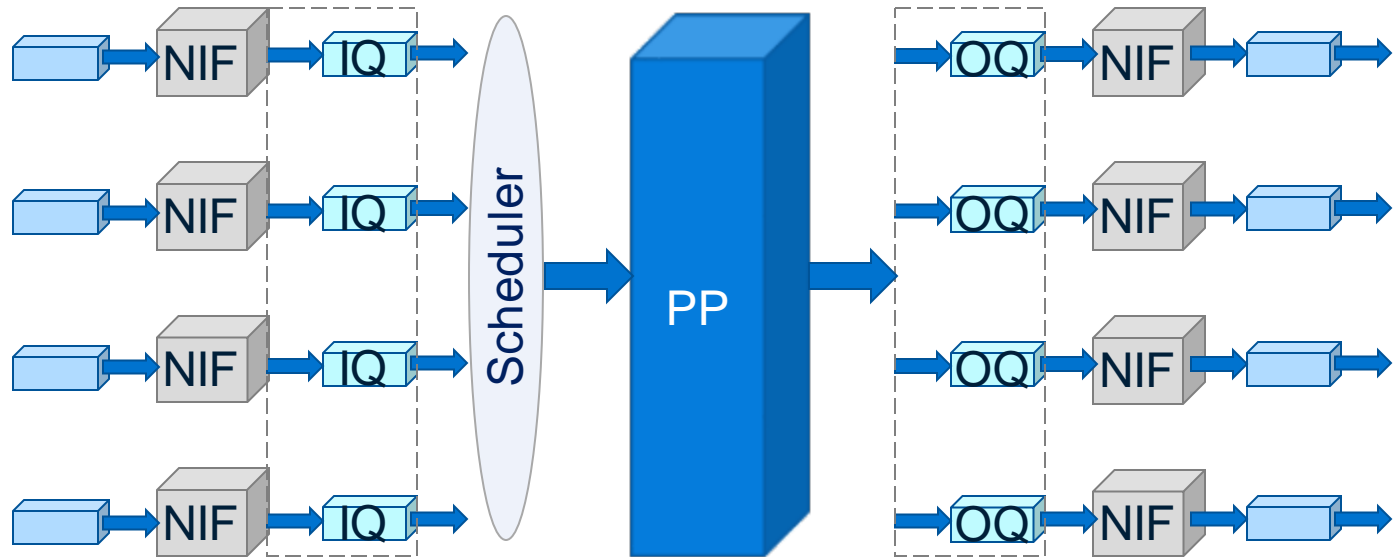
- Sundar, Kompella, and McKeown. "Designing packet buffers for router line cards." 2002.



Cut Through

- Start processing the packet as soon as the first chunk arrives
 - Does not wait for the FCS
- If FCS error is detected, the packet is dropped somewhere along the pipeline

PK13



Latency considerations within modules

Network Interfaces

- Data arrives at (up to) ~50Gbps per link.
- Let us ignore clock recovery, signal detection etc.
- Feasible clock rate is ~1GHz
- But if data rate is $\times 50$ times faster...

- Observation: data bus width will be determined by incoming data rate and feasible device clock rate

Network Interfaces

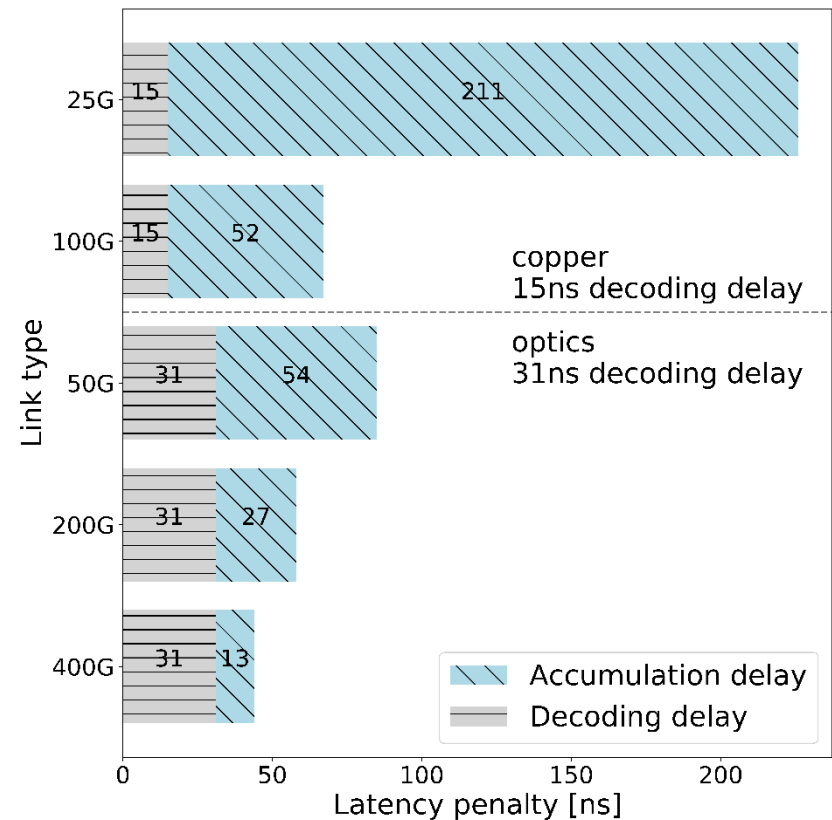
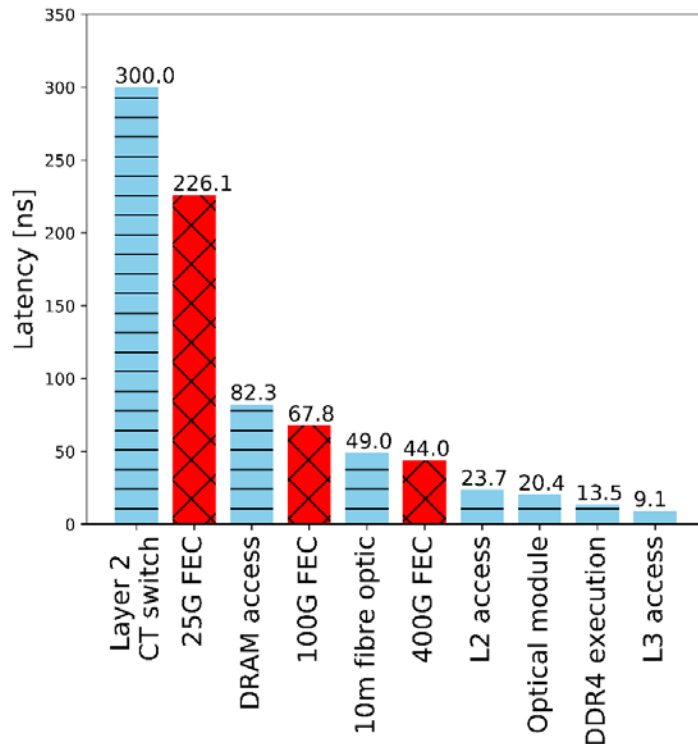
- Line coding often directs the bus widths:
 - E.g., 8b/10b coding led to bus widths of 16b (20b) or 64b (80b)
- A port is commonly an aggregation of multiple serial links
 - 10G XAUI = $4 \times 3.125\text{Gbps}$
 - 100G CAUI4 = $4 \times 25\text{Gbps}$
 - 400G PSM4 = $8 \times 50\text{Gbps}$
 - Need to take care of aligning the data arriving from multiple links on the same port.

Network Interfaces

- Role: check the validity of the packet (e.g., FCS)
- What to do if an error is detected?
 - Forward an error using a “fast path”
 - Mark the last cycle of the packet
 - E.g., to cause drop in the next hop
- Other roles need to be maintained too
 - Frame delimiting and recognition, flow control, enforcing IFG, ...

Network Interfaces - FEC

- FEC – Forward Error Correction
 - Current standards use Reed-Solomon block codes
- Can add significant latency penalty



Packet Processing

- A likely flow:



- Possible implementations:
 - The entire packet goes through the header processing unit
 - Just the header goes through the header processing unit
 - “Better” depends on your performance profile (what are the bottlenecks? Resource limitations?)

Packet Processing

- A likely flow:



- Challenges:

- A field may arrive over multiple clock cycles (e.g. 32b field, 16b on clock 2 and 16b on clock 3)
- Memory access taking more than 1 clock cycle
 - E.g. request on clock 1, reply on clock 3
 - Some memories allow multiple concurrent accesses, some don't
 - The bigger the memory, the more time it takes

Packet Processing

- A likely flow:



- Solutions:

- Pipelining!

Don't stall, add NOP stages in your pipe.

- Reorder operations (where possible)

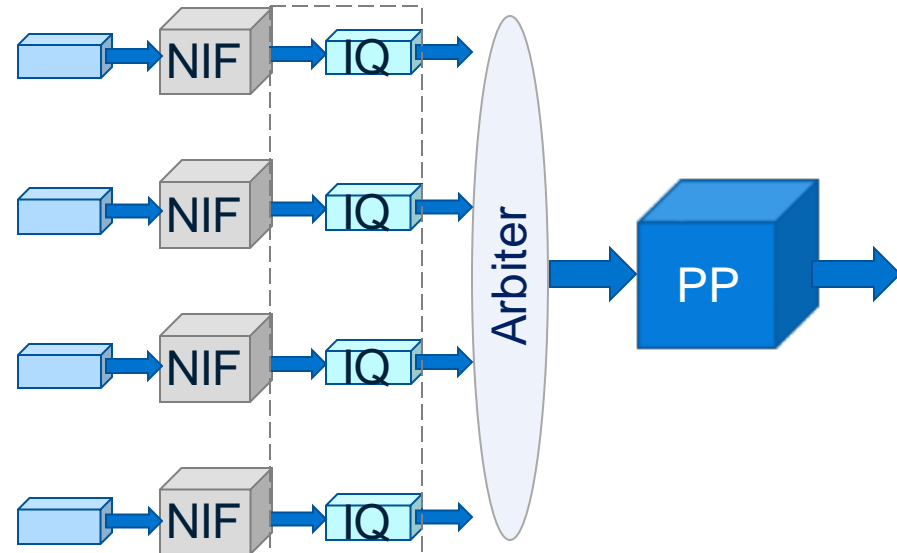
- E.g. Lookup 1 → Action 1 → Lookup 2 → Action 2 becomes:

Lookup 1 → Lookup 2 → Action 1 → Action 2

- Don't create hazards!

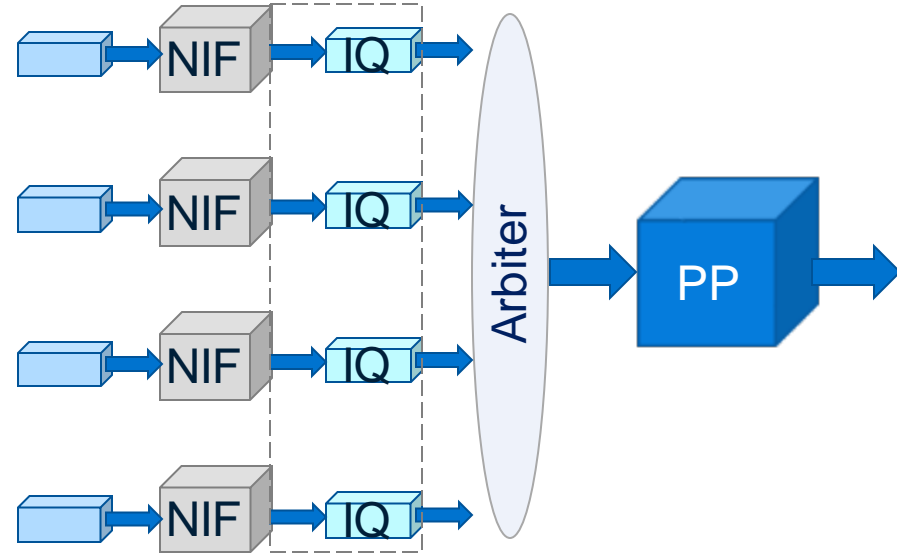
Arbitration

- Simple example:
 - Packets arriving from 4 ports
 - (approximately) same arrival time
 - Arbiter uses Round Robin
- Problem: arbitration on packet boundaries?
 - No: interleaved packets within the pipeline
Need to track which cycle belongs to which packet
May require multiple concurrent header lookups
Order is not guaranteed (e.g. P1-P2-P3-P1-P2-P2-...), due to NIF timing



Arbitration

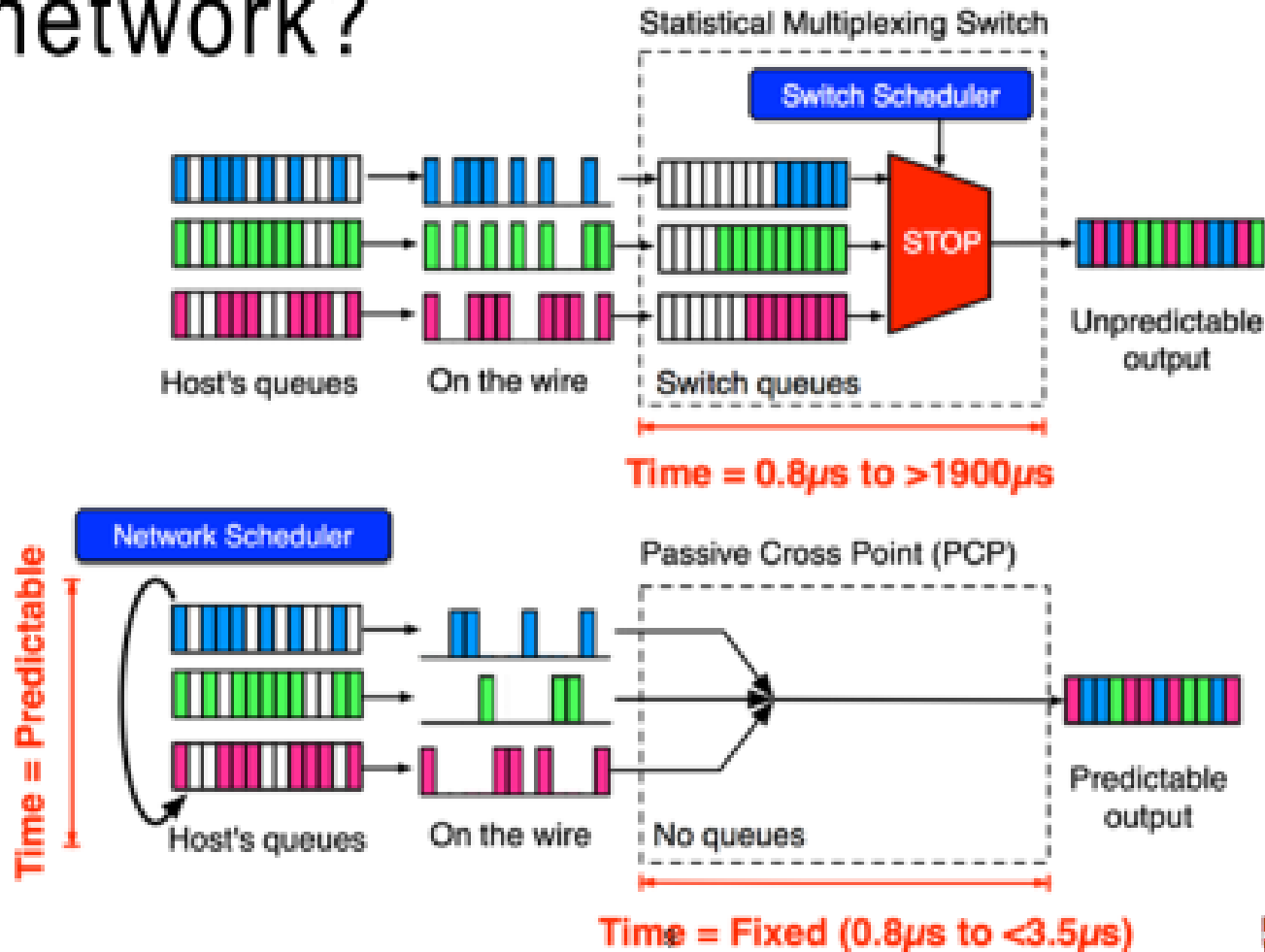
- Simple example:
 - Packets arriving from 4 ports
 - (approximately) same arrival time
 - Arbiter uses Round Robin
- Problem: arbitration on packet boundaries?
 - Yes: packets need to wait for previous packets to be handled before being admitted.
Worst case waiting with $\langle N \rangle$ inputs is $\langle N-1 \rangle \times \text{Packet time}$



Arbitration

- Solutions to the previous problem:
 - Scheduled (or slotted) traffic
 - Multiple pipelines
 - ...

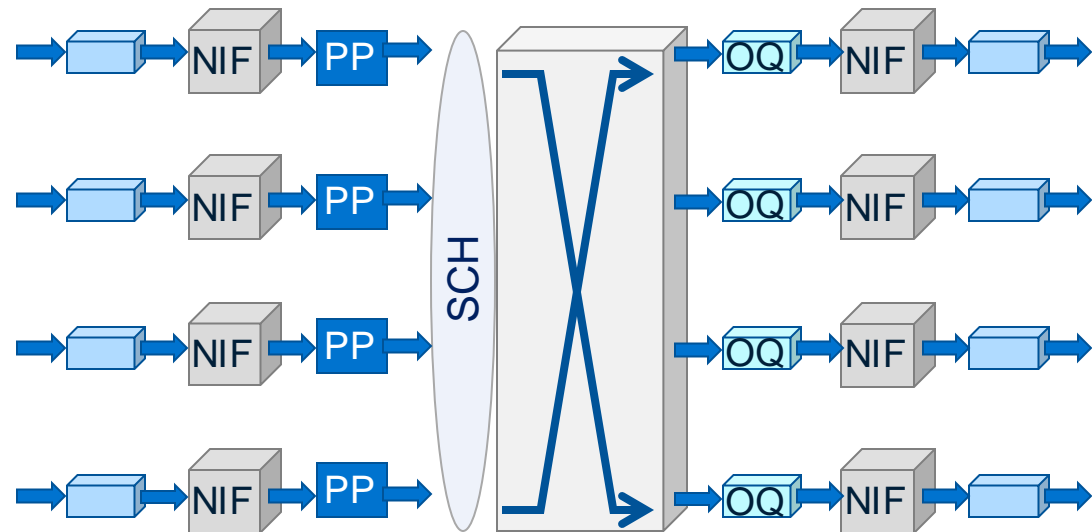
How do you build a bufferless network?



Thursday, 26 September 13

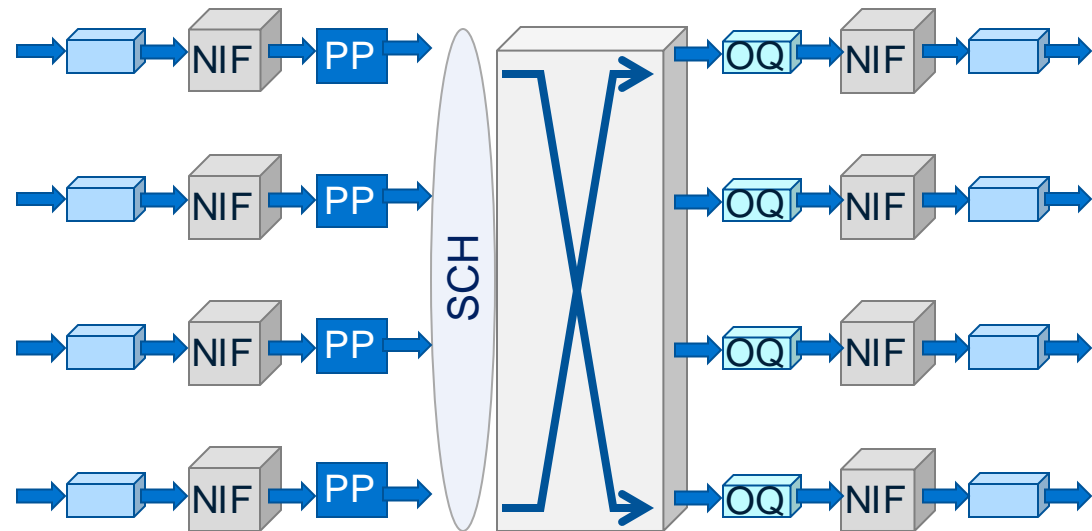
Arbitration

- This example solves the arbitration problem entering the device
- Resource inefficient:
 - Pipeline overdesign
 - Inefficient use of memories
- Concurrency issues
- One solution:
 - Shared memories / tables
 - Highly complex

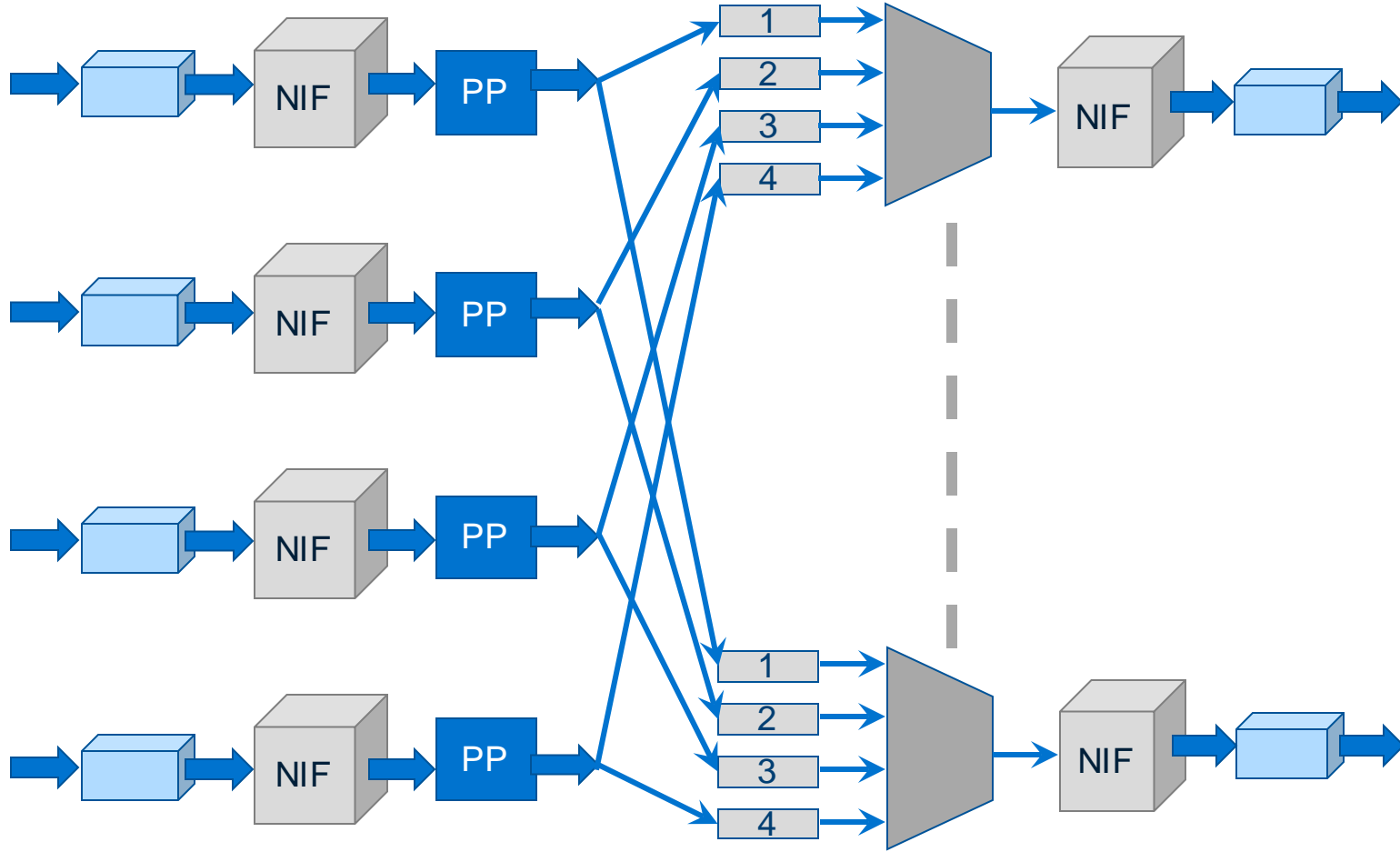


Switching

- The previous arbitration solution “pushed” the problem to the switching unit
- But now the problem is only when multiple packets compete over the same output – that’s fine!
- Assuming your switch can handle multiple packets per cycle
 - E.g. crossbar



Switching

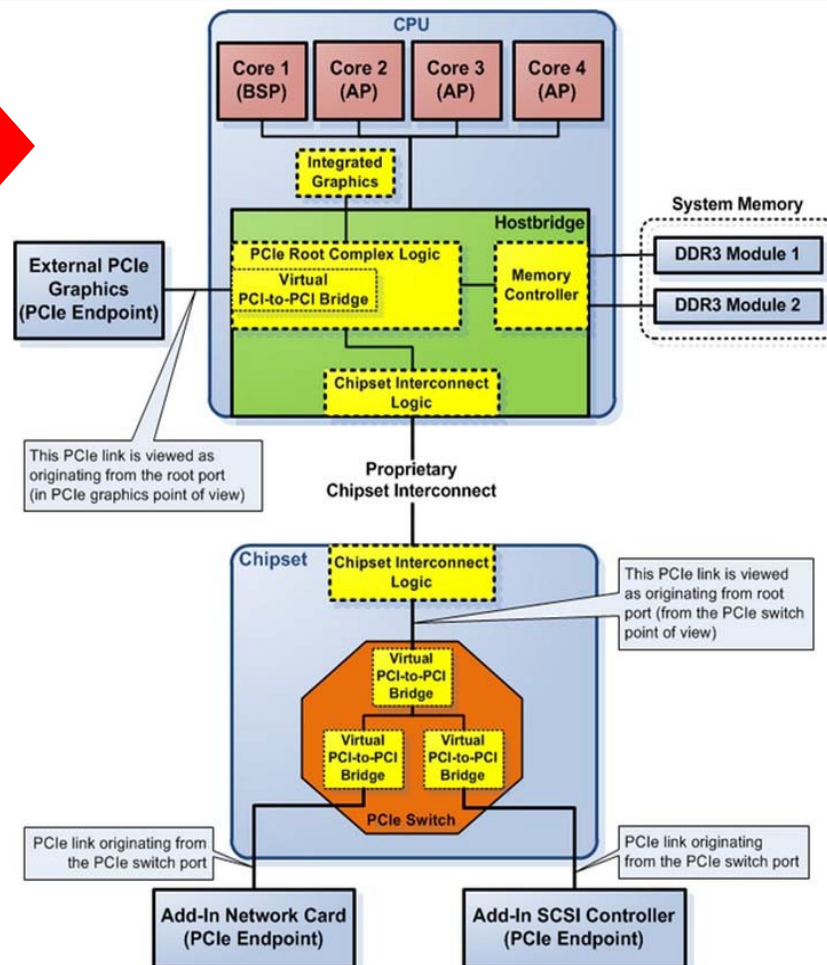
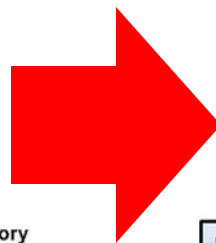
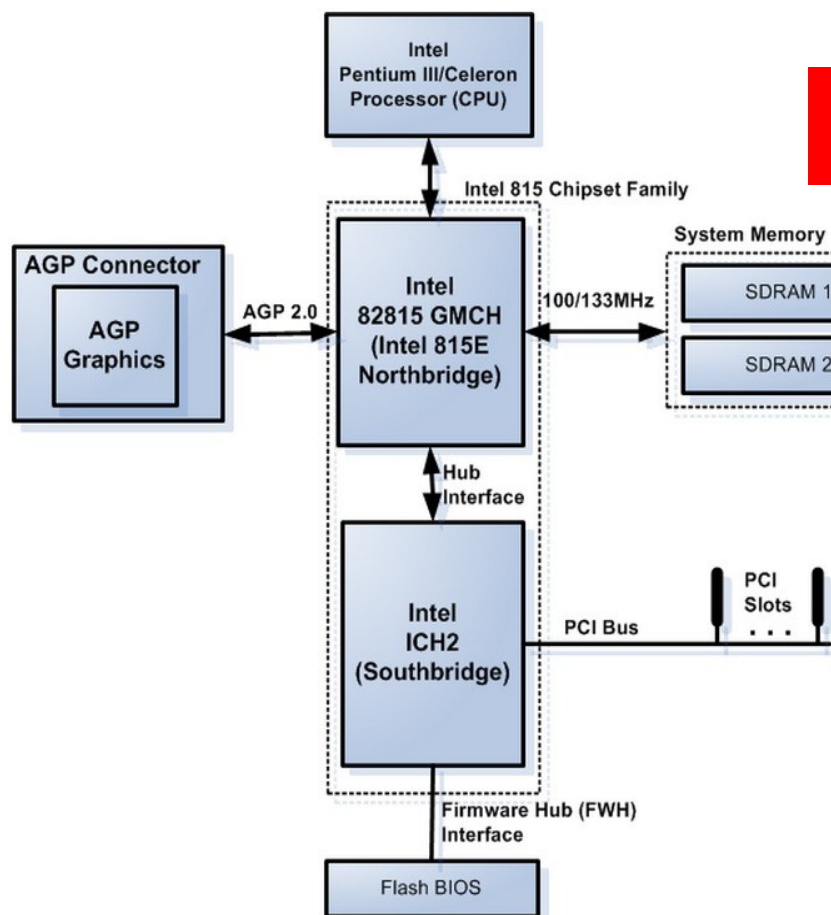


Switching

- ... This is also queueing
- Challenge: SCALE
 - So you can do it with 4 ports
 - Can you do it with 32? 128? 256?
- Not just resource / area
 - Computation time – being able to examine and choose between all available inputs
- Eventually:
Packets must be sent out on packet boundaries
Do not interleave packets!

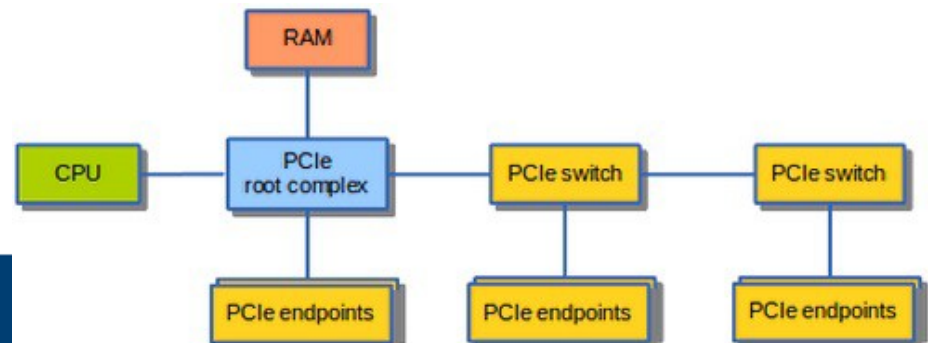
DMA

Host architecture



Interconnecting components

- **Need interconnections between**
 - CPU, memory, storage, network, I/O controllers
- **Shared Bus: shared communication channel**
 - A set of parallel wires for data and synchronization of data transfer
 - Can become a bottleneck
- **Performance limited by physical factors**
 - Wire length, number of connections
- **More recent alternative: high-speed serial connections with switches**
 - Like networks



I/O System Characteristics

- **Performance measures**
 - Latency (response time)
 - Throughput (bandwidth)
 - Desktops & embedded systems
 - Mainly interested in response time & diversity of devices
 - Servers
 - Mainly interested in throughput & expandability of devices
- **Reliability**
 - Particularly for storage devices (fault avoidance, fault tolerance, fault forecasting)

I/O Management and strategies

- **I/O is mediated by the OS**
 - Multiple programs share I/O resources
 - Need protection and scheduling
 - I/O causes asynchronous interrupts
 - Same mechanism as exceptions
 - I/O programming is fiddly
 - OS provides abstractions to programs

Strategies characterize **the amount of work** done by the CPU in the I/O operation:

- **Polling**
- **Interrupt Driven**
- **Direct Memory Access**

The I/O Access Problem

- Question: how to transfer data from I/O devices to memory (RAM)?
- Trivial solution:
 - Processor individually reads or writes every word
 - Transferred to/from I/O through an internal register to memory
- Problems:
 - Extremely inefficient – can occupy a processor for 1000's of cycles
 - Pollute cache