

The Future of CPU Bus Architectures

A Cyrix Perspective

**Forrest Norrod, Senior Director
System Development & Strategic Planning**

Copies of slides @ www.cyrix.com



© 1998 Cyrix Corporation



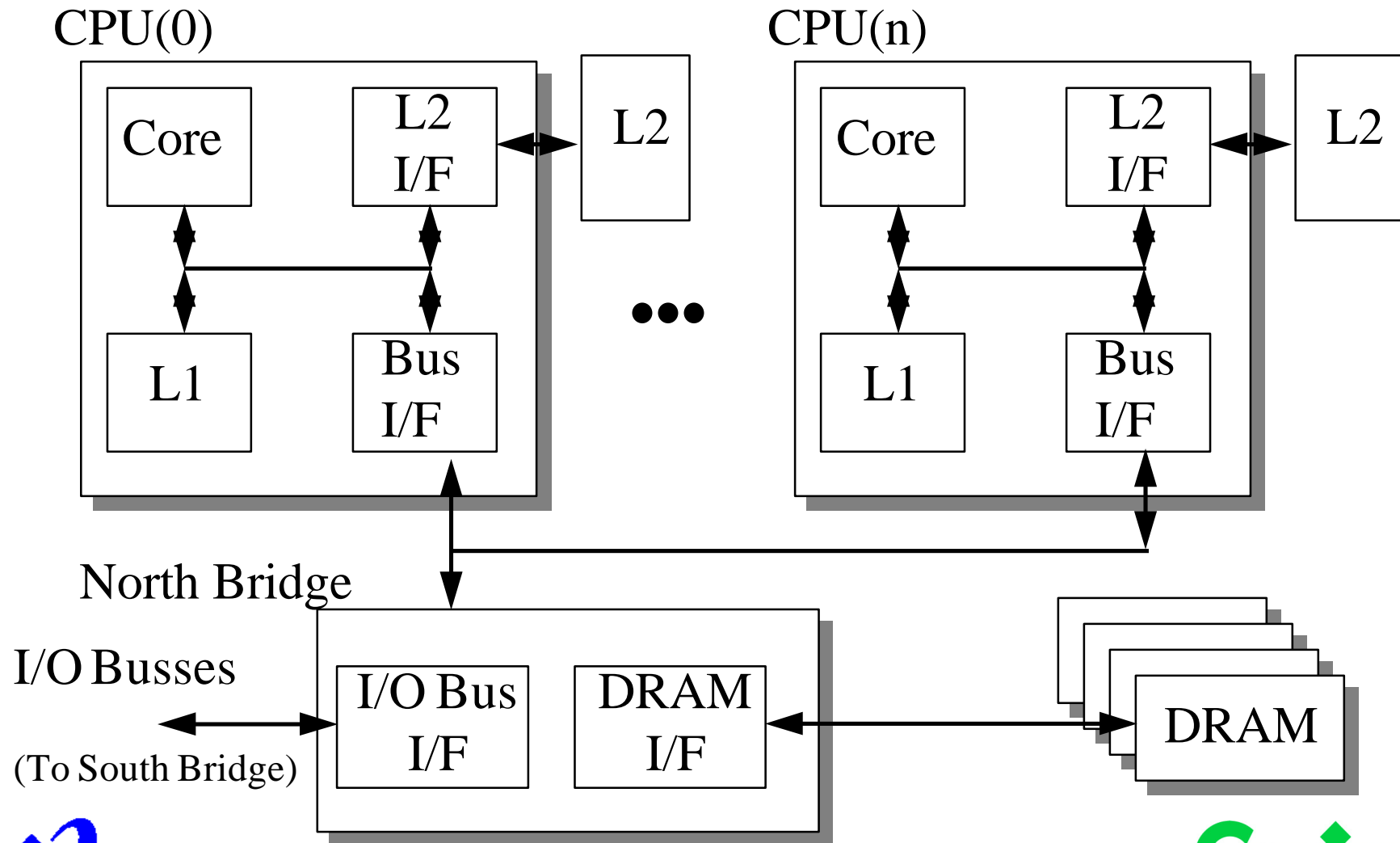
What are CPU Busses For?

- **CPU Busses have three main purposes:**
 - To communicate with external elements of the system's memory hierarchy
 - To link the CPU core with graphics, network and other peripheral devices
 - In multiprocessor systems, to link multiple CPUs
- **Bus Design Objectives:**
 - Minimize latencies to memory & I/O
 - Maximize bandwidth to all devices
 - Minimize cost & implementation complexity

Multiprocessor vs. Uniprocessor

- **MP optimized for multiple outstanding transactions**
 - **Deeply pipelined request queues**
 - » Striving for highest possible bandwidth utilization
 - **Pipelined arbitration**
 - **Extensible snooping capabilities**
 - **Deferred transaction support**
 - **Complex protocol for devices to implement & monitor**
- **Uniprocessor optimized for lowest latency**
 - **Typically, only one master on the bus**
 - **Limited pipelining**
 - » Trying to hide DRAM access times
 - **Simple, “light weight” protocol easy to implement**

Typical Multiprocessor System

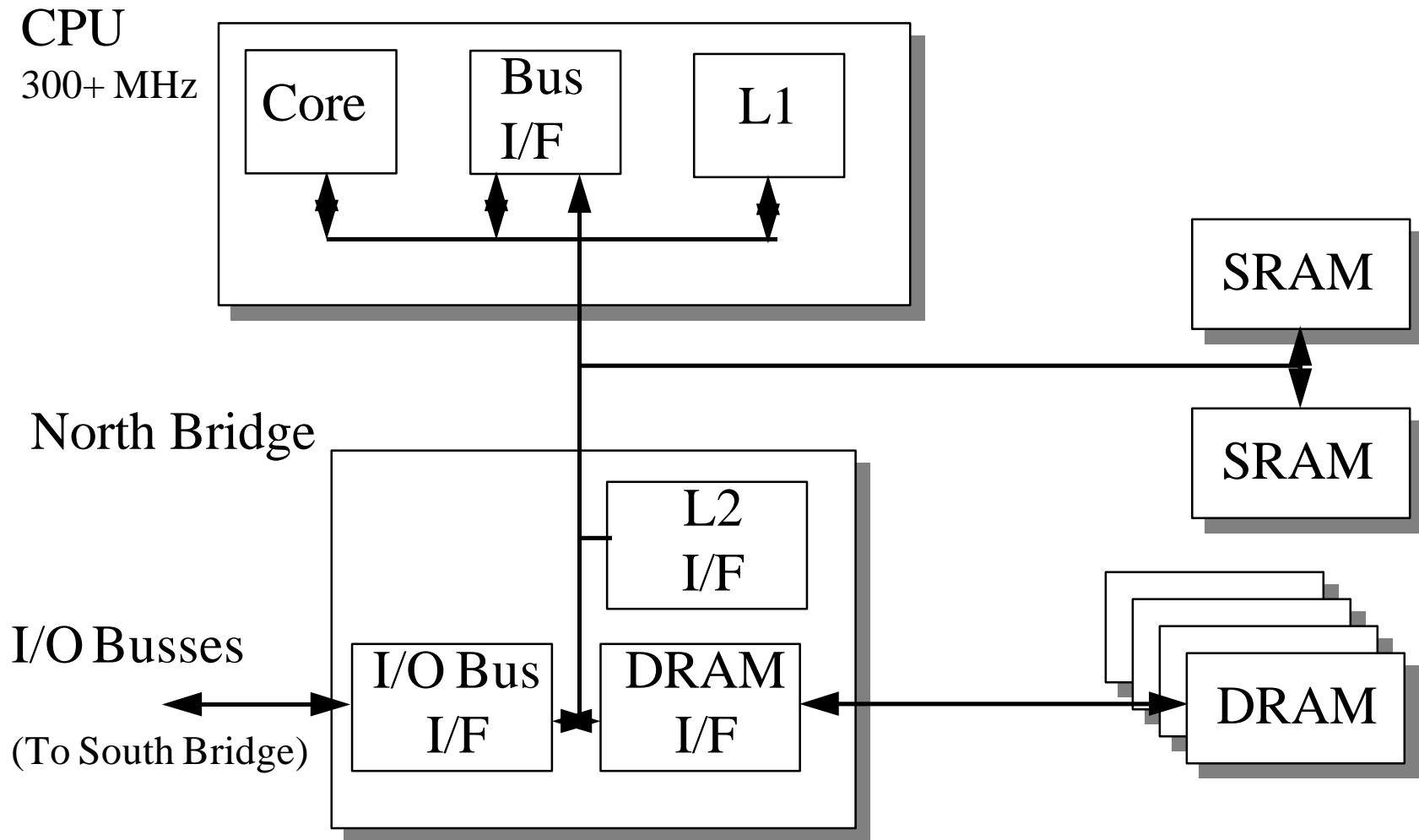


Multiprocessor Bus Protocols

- **Each read transaction must:**
 - Arbitrate for the address bus
 - Send transaction request
 - Check for errors
 - Check coherency on the transaction (snoop)
 - Request the data bus to return data
 - Send the data
- **Each bus agent must track all transactions**
 - Typically has to track 4-16 pending transactions
- **Complex structure is fairly expensive**
 - Significant pins & gates to implement bus
 - Increased latency to main memory

–

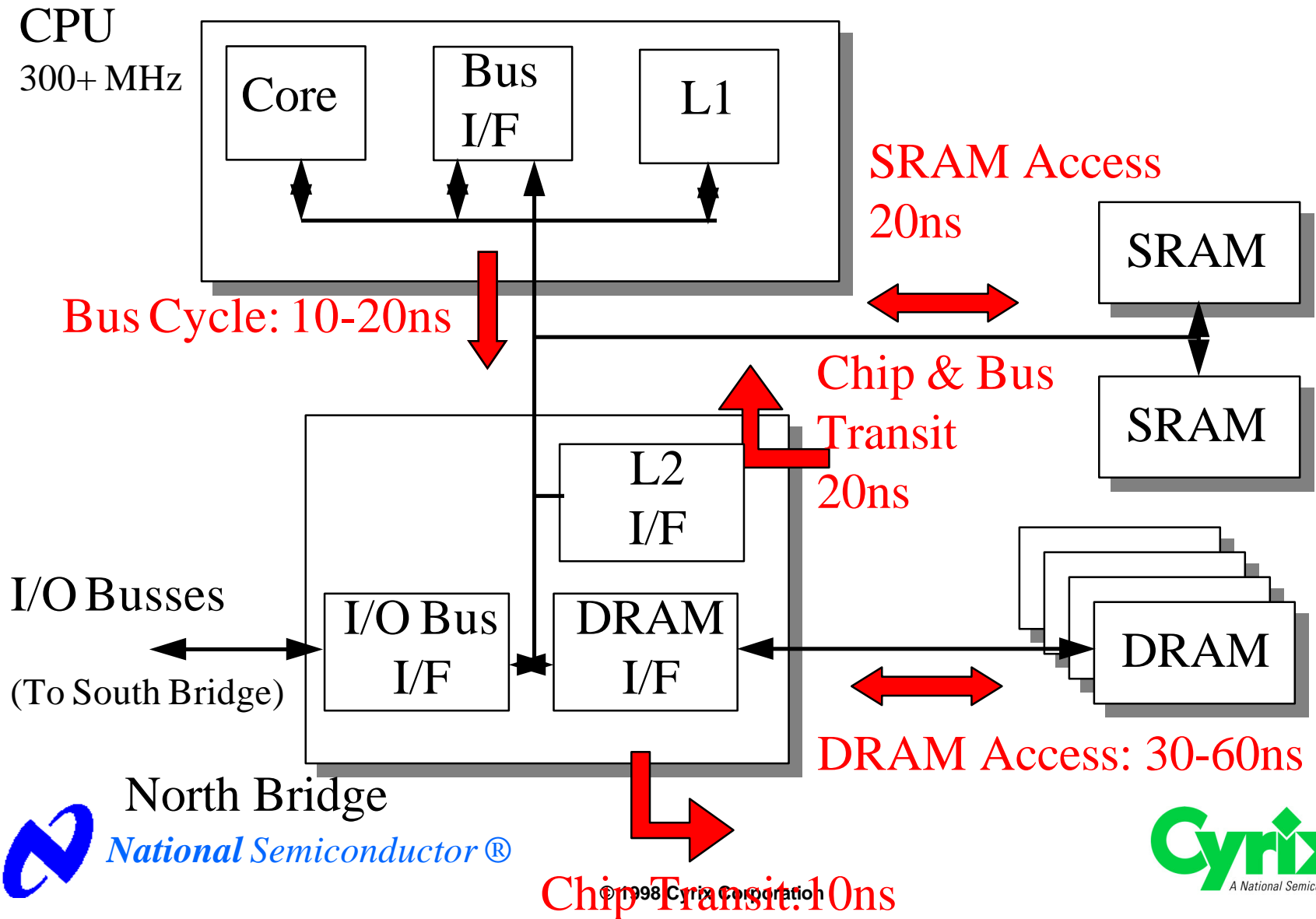
Typical Uniprocessor Architecture



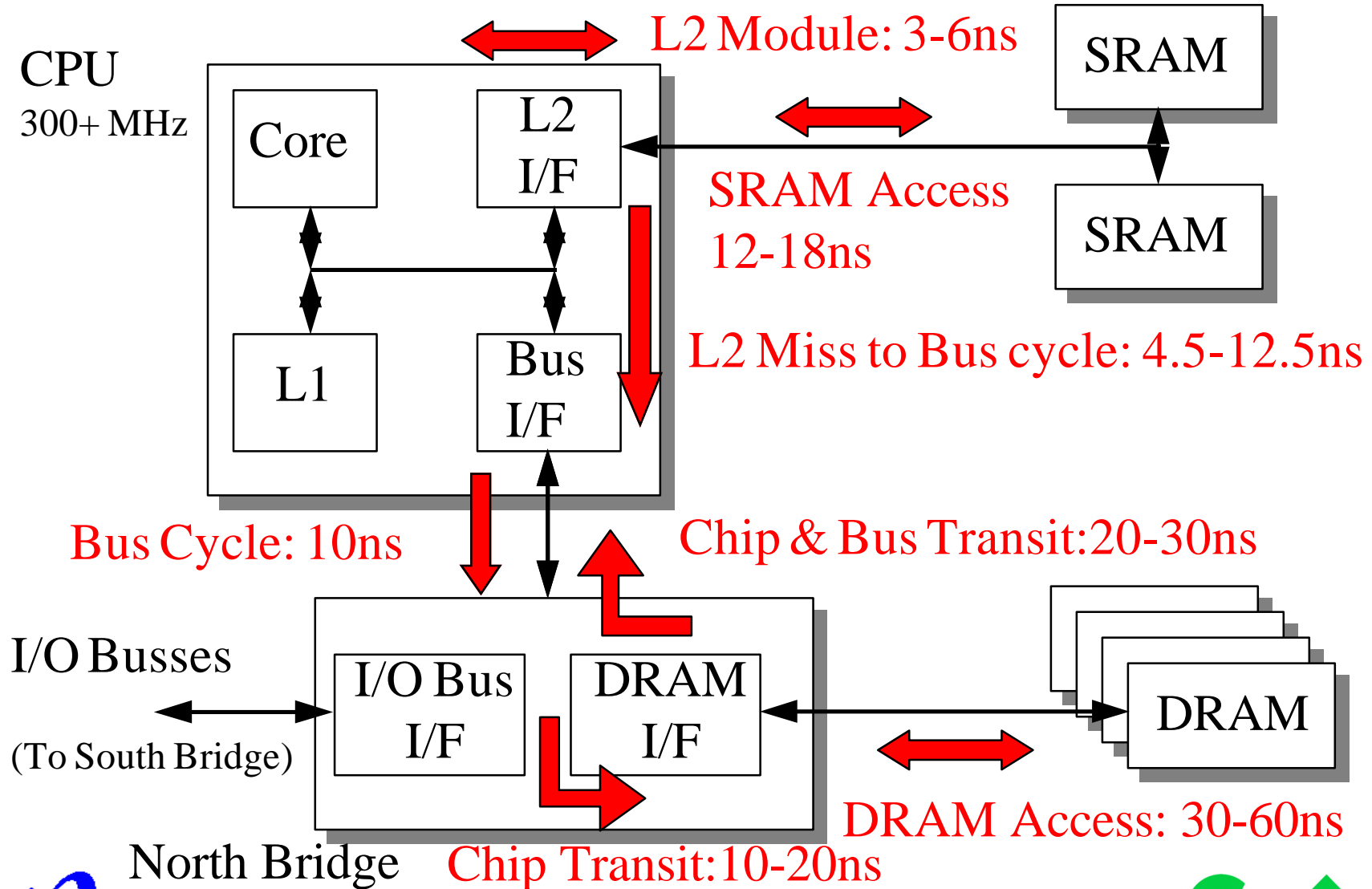
Uniprocessor Bus Protocols

- **Simple bus connection to north bridge & L2**
 - **CPU always (almost) the bus master**
 - » Eliminates arbitration time
 - » No need to snoop other bus agents (except L2)
 - » No need to support more than a few transactions
 - **Lowers latency to main memory**
- **Bus structure is easier to implement**
 - **Minimal fanout (point-to-point if no frontside L2)**
- **L2 can be frontside or backside**
 - **Backside: offers possibly lower access times to L2**
 - **Frontside: easy to start DRAM access in parallel**
-

Typical Frontside Architecture



“Traditional” Backside Architecture



North Bridge
National Semiconductor®

© 1998 Cyrix Corporation



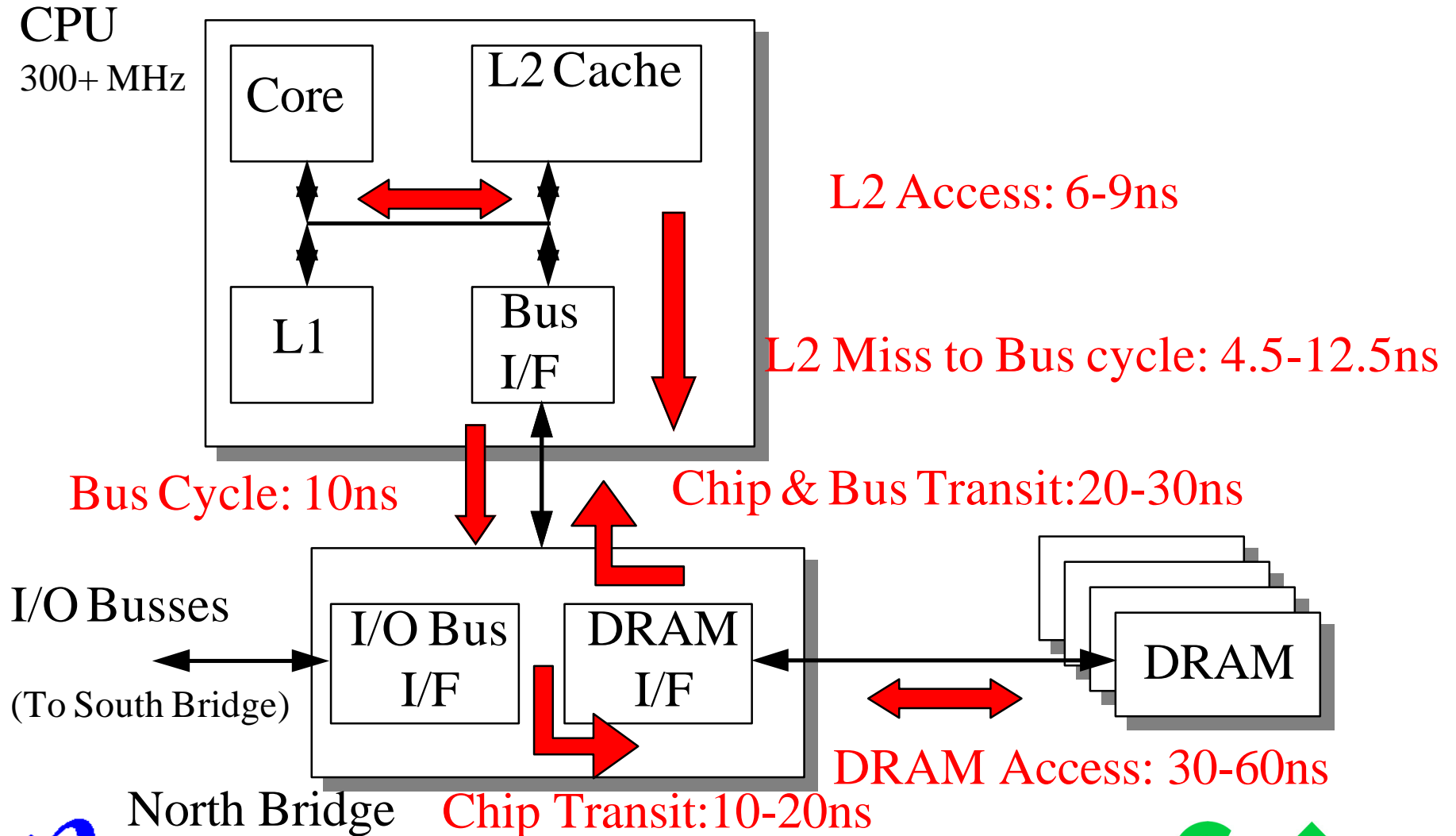
The Problems with Busses

- **System performance is determined by memory access time**
 - CPUs are tremendously faster than DRAM
 - CPU is, more and more, waiting for data
 - Avg Latency = $L1_hit_rate * L1_latency$
+ $L2_hit_rate * L2_latency$
+ $DRAM_hit_rate * DRAM_latency$
- **Backside L2 caches still have long core clock latencies**
 - L2 controller traversal: 1-2 core clocks
 - L2 lookup: 4-8 core clocks (at 300MHz)
 - Total: 5-10 core clocks (and increasing)

What Choice Do We Have?

- **Increase the size of the L1**
 - Great bang-for-the-buck - fastest possible access
 - BUT, increasing the size decreases the speed
 - » Looks like 16-32Kbyte L1's are here to stay
- **Increase the size of the L2**
 - Increases the hit rate, reducing DRAM accesses
 - Again, at some point it lowers the speed of the L2
 - Not very cost effective
- **Increase the speed of the L2**
 - Can be done with external parts, but very difficult
 - Integration of the L2 significantly increases the speed

Integrated L2 Architecture



North Bridge
National Semiconductor®

© 1998 Cyrix Corporation



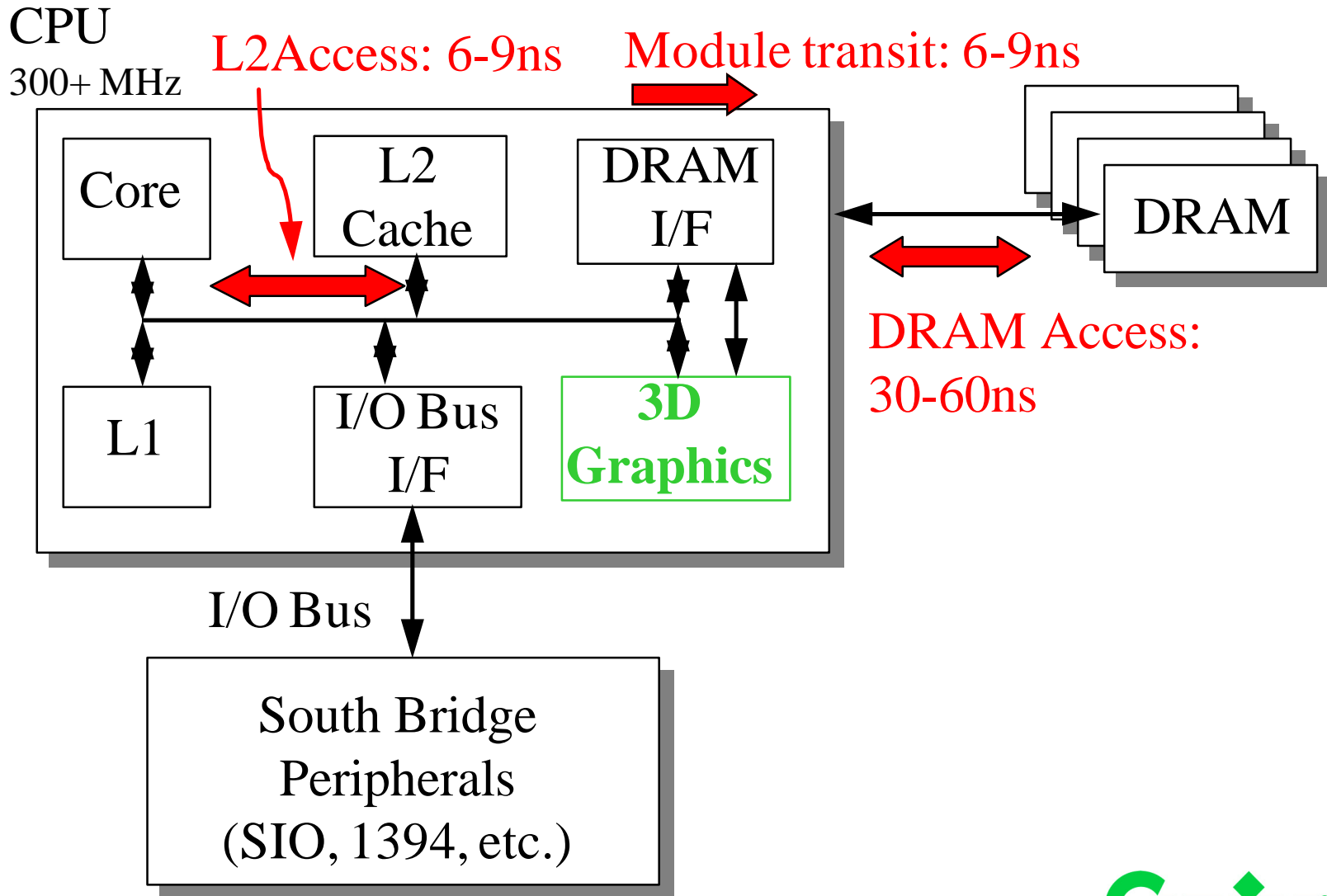
Integrated L2 vs. External

- **Compare an external 256K L2 to an internal 128K**
 - **256K external L2 (Winstone '98, 32K L1, 400 MHz)**
 - L1 hit rate of ~97%
 - L2 hit rate of ~84% latency of 8 core clocks
 - DRAM page hit rate of ~55%
 - Page hit latency of 37 core clocks
 - Page miss latency of 49 core clocks
 - Average latency = $.226 \text{ (L2)} + .082 \text{ (Page hit)} + .175 \text{ (miss)}$
= .475 core clocks
 - **For an integrated 128K L2:**
 - L2 hit rate decreases to 78%
 - L2 latency of 3 core clocks
 - Average latency = $.079 \text{ (L2)} + .107 \text{ (Page hit)} + .23 \text{ (miss)}$
= .418 core clocks **(12% reduction)**

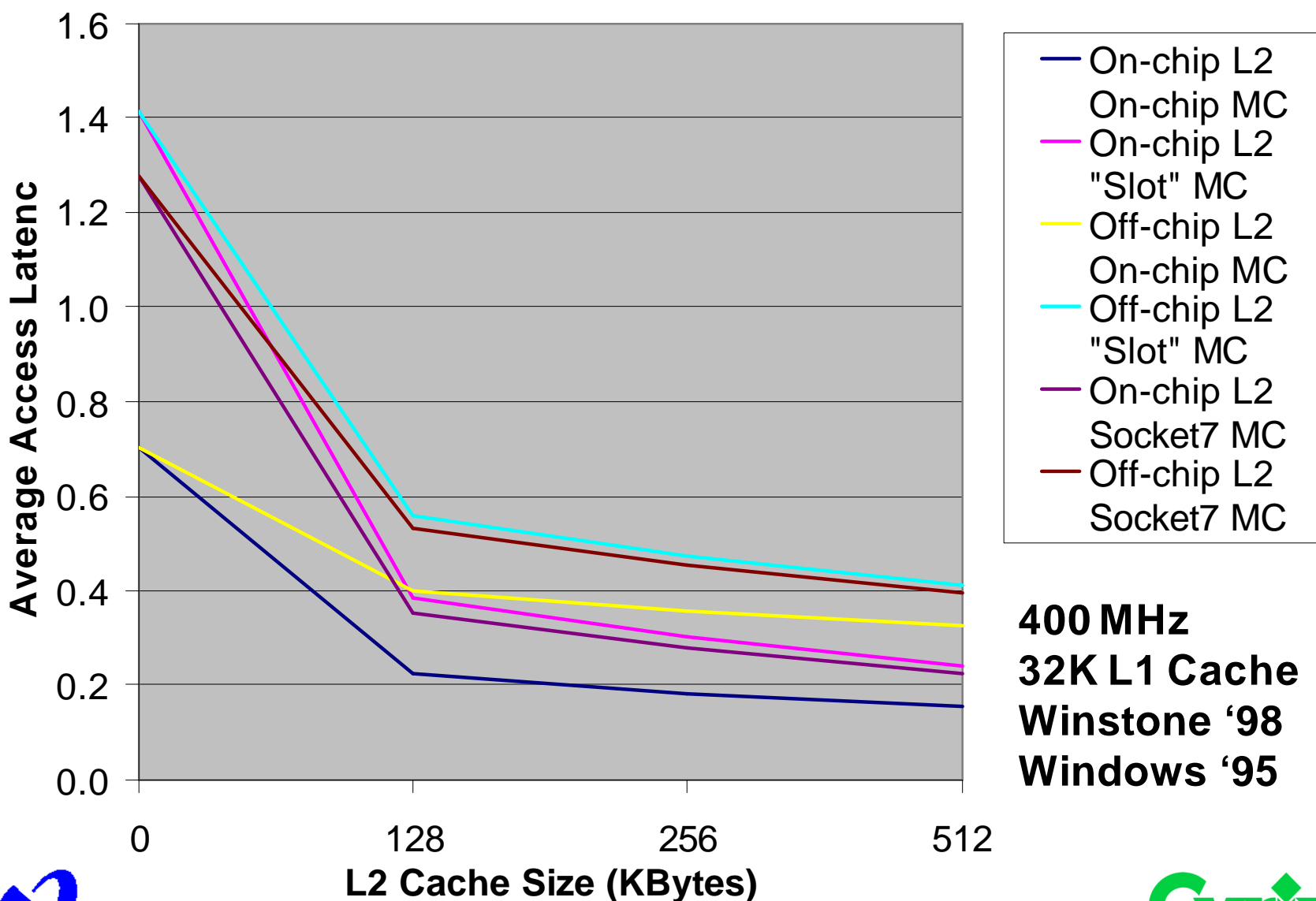
What about the DRAM?

- **The DRAM is still a long way away - & that matters**
 - **Example: (Winstone '98, 32K Unified L1, 400 MHz)**
 - L1 hit rate of ~97%
 - L2 hit rate of ~84% (256K), latency of 3 core clocks
 - DRAM page hit rate of ~55%
 - Page hit latency of 37 core clocks
 - Page miss latency of 49 core clocks
 - Average latency = $.111$ (L2) + $.082$ (Page hit) + $.175$ (miss)
= $.368$ core clocks
 - **If the DRAM were directly accessed by the CPU**
 - Page hit latency of 11 core clocks
 - Page miss latency of 23 core clocks
 - Average latency = $.111$ (L2) + $.024$ (Page hit) + $.082$ (miss)
= $.218$ core clocks (**41% reduction**)

Integrated DRAM Controller



Access Latencies



400 MHz
32K L1 Cache
Winstone '98
Windows '95



Summary

- **The overhead of multiprocessor busses makes them inappropriate for uniprocessor systems**
 - Any increase in latency to memory is unacceptable
- **Integration of L2 and memory controller onto the CPU yields the best uniprocessor performance**
 - Absolute lowest average latency to the memory hierarchy
- **This level of integration is now possible with .25 and .18 micron processes**
 - 128 Kbyte L2 cache is < 25 sq. mm in .18 micron
 - Cost of additional die area is smaller than the cost of the additional pins / packages of a traditional design
 -