

CS4803DGC Design Game Consoles

Spring 2010

Prof. Hyesoon Kim



**Georgia
Tech**



College of
Computing

Xbox 360 System Architecture, 'Anderews, Baker

Xbox 360 System Block Diagram

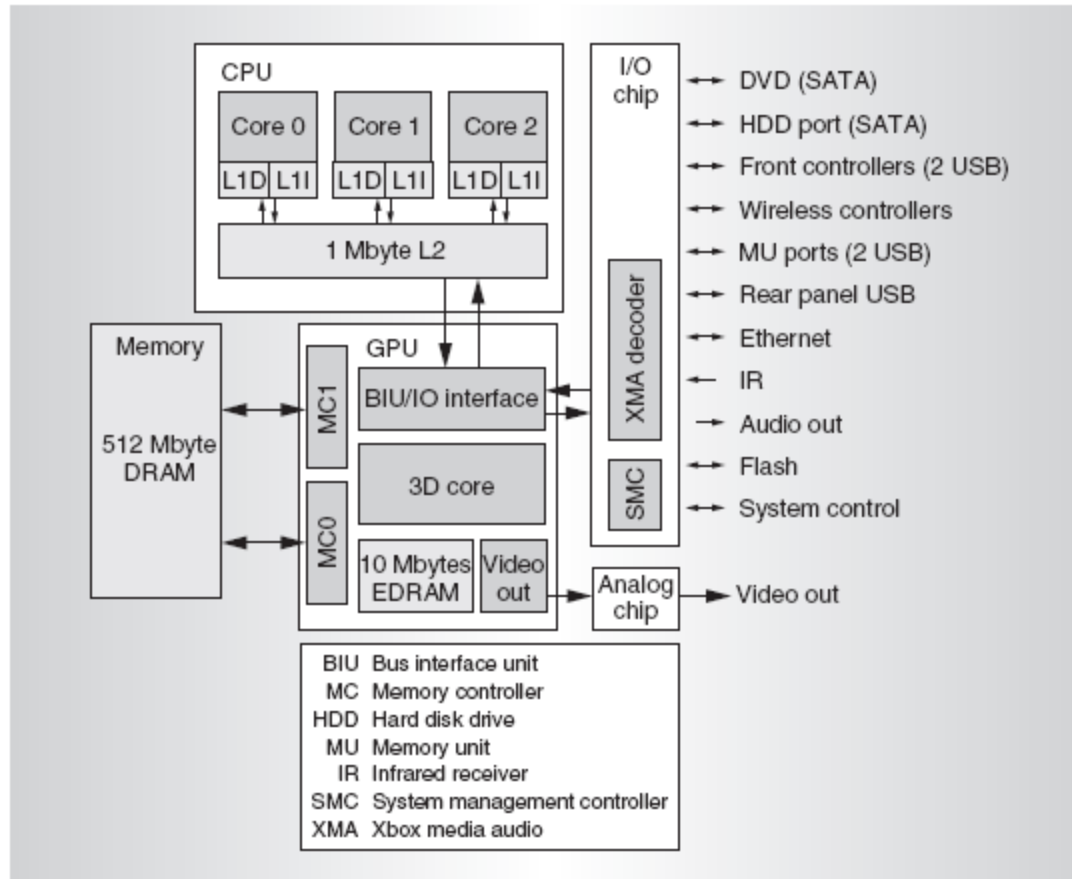


Figure 2. Xbox 360 system block diagram.



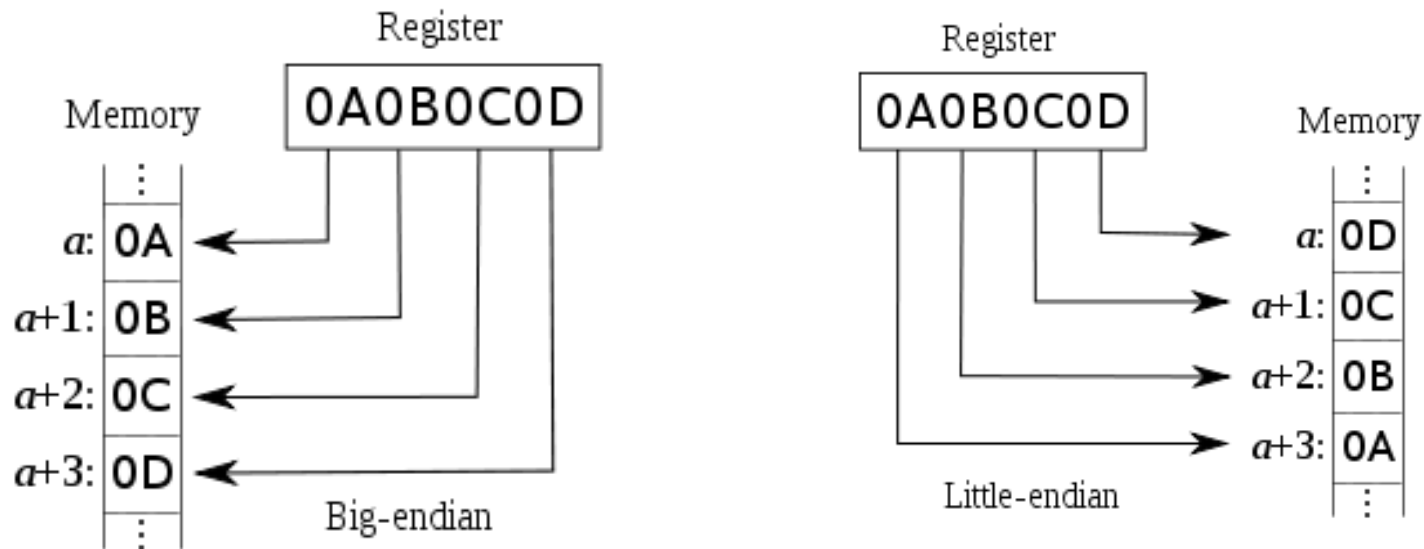
Xbox 360 Architecture

- 3 CPU cores
 - 4-way SIMD vector units
 - 8-way 1MB L2 cache (3.2 GHz)
 - 2 way SMT
- 48 unified shaders
- 3D graphics units
- 512-Mbyte DRAM main memory
- FSB (Front-side bus): 5.4 Gbps/pin/s (16 pins)
- 10.8 Gbyte/s read and write

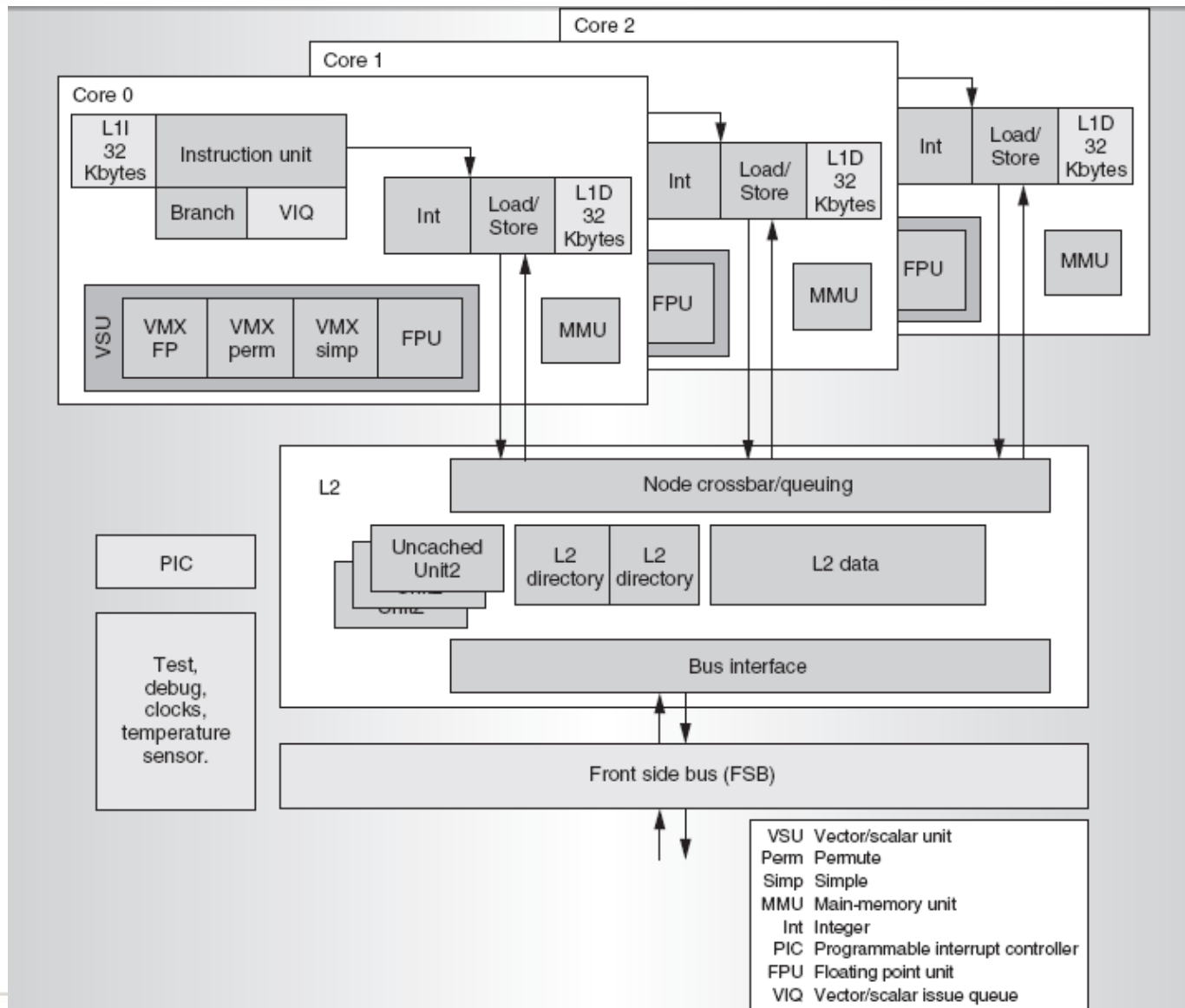


Xbox 360 vs. Windows

- Xbox 360: Big endian
- Windows: Little endian



Xbox 360 CPU Block Diagram





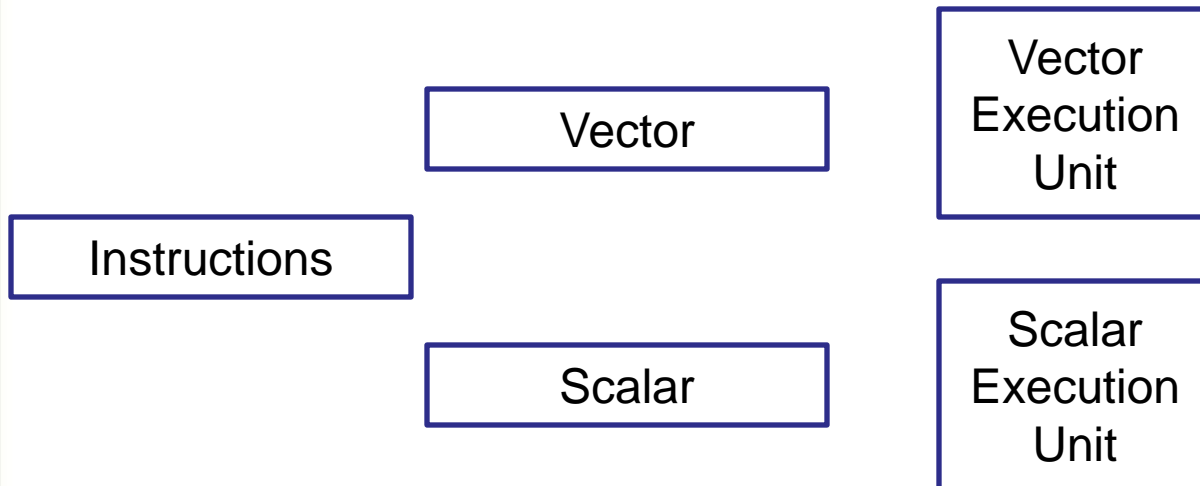
On-chip caches

- L2 cache :
 - Greedy allocation algorithm
 - Different workloads have different working set sizes
- 2-way 32 Kbyte L1 I-cache
- 4-way 32 Kbyte L1 data cache
- Write through, no write allocation
- Cache block size :128B (high spatial locality)



Core

- 2-way SMT,
- 2 insts/cycle,
- In-order issue
- Separate vector/scalar issue queue (VIQ)





VMX 128

- Four-way SIMD VMX 128 units:
 - FP, permute, and simple
- 128 registers of 128 bits each per hardware thread
- Added dot product instruction (simplifying the rounding of intermediate multiply results)
- 3D compressed data formats . Use compressed format to store at L2 or memory. 50% of space saving.

A Brief History



- First game console by Microsoft, released in 2001, \$299
 - Glorified PC
 - 733 Mhz x86 Intel CPU, 64MB DRAM, NVIDIA GPU (graphics)
 - Ran modified version of Windows OS
 - ~25 million sold
- XBox 360
 - Second generation, released in 2005, \$299-\$399
 - All-new custom hardware
 - 3.2 Ghz PowerPC IBM processor (custom design for XBox 360)
 - ATI graphics chip (custom design for XBox 360)
 - 34+ million sold (as of 2009)
- Design principles of XBox 360 [Andrews & Baker]
 - Value for 5-7 years
 - ! big performance increase over last generation
 - Support anti-aliased high-definition video (720*1280*4 @ 30+ fps)
 - ! extremely high pixel fill rate (goal: 100+ million pixels/s)
 - Flexible to suit dynamic range of games
 - ! balance hardware, homogenous resources
 - Programmability (easy to program)



Xenon

- Code name of Xbox 360's core
- Shared cell (playstation procesor) 's design philosophy.
- 2-way SMT
- Good: Procedural synthesis is highly multi-thread
- Bad: three types of game-oriented tasks are likely to suffer from the lack of high ILP support: game control, artificial intelligence (AI), and physics.



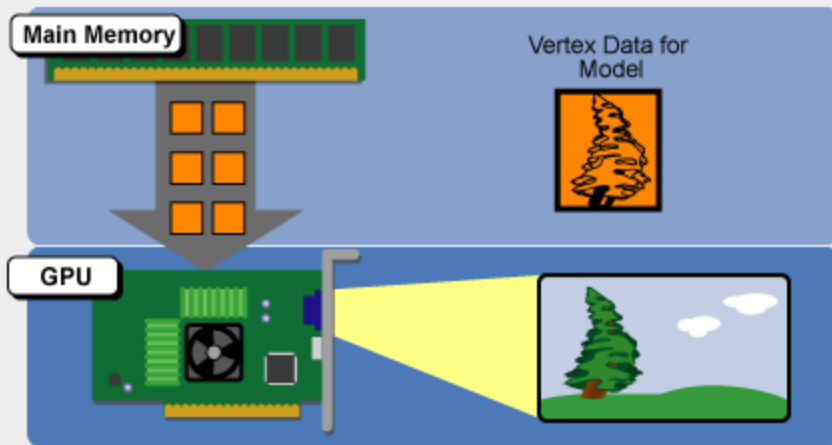
Xenon Processor

- ISA: 64-bit PowerPC chip
 - RISC ISA
 - Like MIPS, but with condition codes
 - Fixed-length 32-bit instructions
 - 32 64-bit general purpose registers (GPRs)
- ISA++: Extended with VMX-128 operations
 - **128 registers, 128-bits each**
 - Packed “vector” operations
 - Example: four 32-bit floating point numbers
 - One instruction: VR1 * VR2 ! VR3
 - Four single-precision operations
 - Also supports conversion to MS DirectX data formats
- Works great for 3D graphics kernels and compression
- 3.2 GHZ
- Peak performance Peak performance: ~75 gigaflops

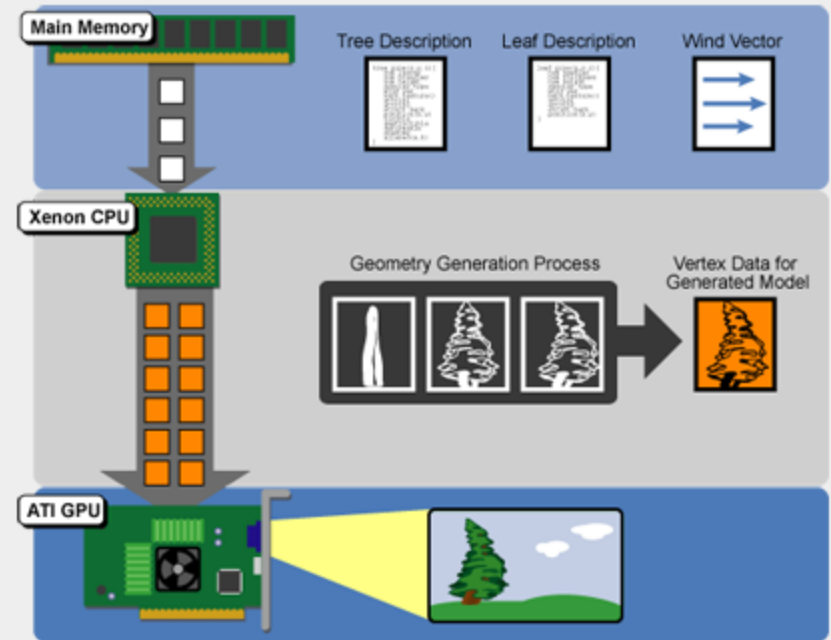


Procedural Synthesis

- Microsoft refers to this ratio of stored scene data to rendered vertex data as a **compression ratio**, the idea being that main memory stores a "compressed" version of the scene, while the GPU renders a "decompressed" version of the scene.



Rendering a wind-blown tree, the conventional way



Rendering a wind-blown tree on the Xbox 360



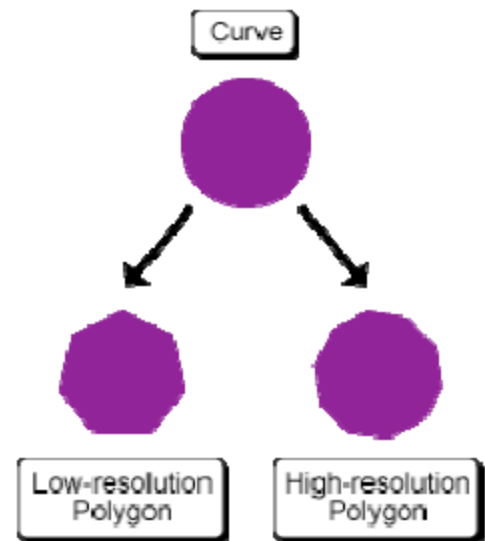
The Benefits of Procedure Synthesis

- Scalable “virtual” artists
- Reduction of bandwidth from main memory to GPUs



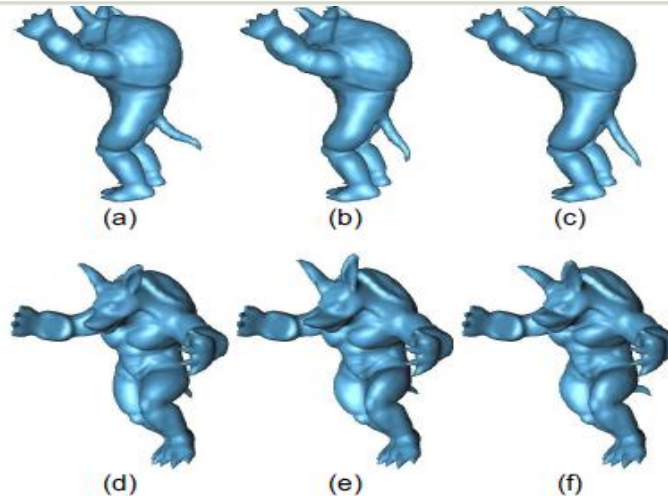
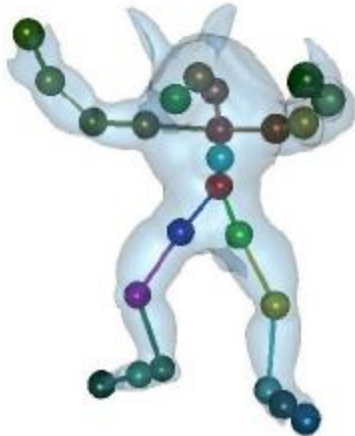
Real-time Tessellation

- Tessellation: The process of taking a higher order curve and approximating it with a network of small flat surfaces is called tessellation.
- Traditional GPU: Artist
- Xbox 360: using Xeon
- Real time tessellation
 - Another form of data compression
 - Instead of list of vertex, stores them as higher order of curves
 - Dynamic Level of Detail (LOD)
 - Keep the total number of polygons in a scene under control





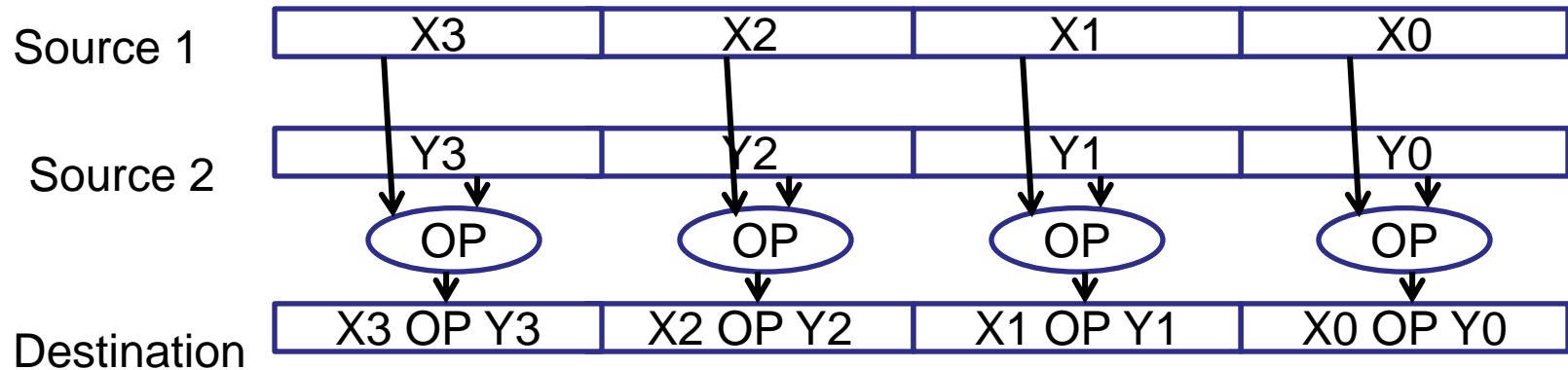
Real-time Skinning



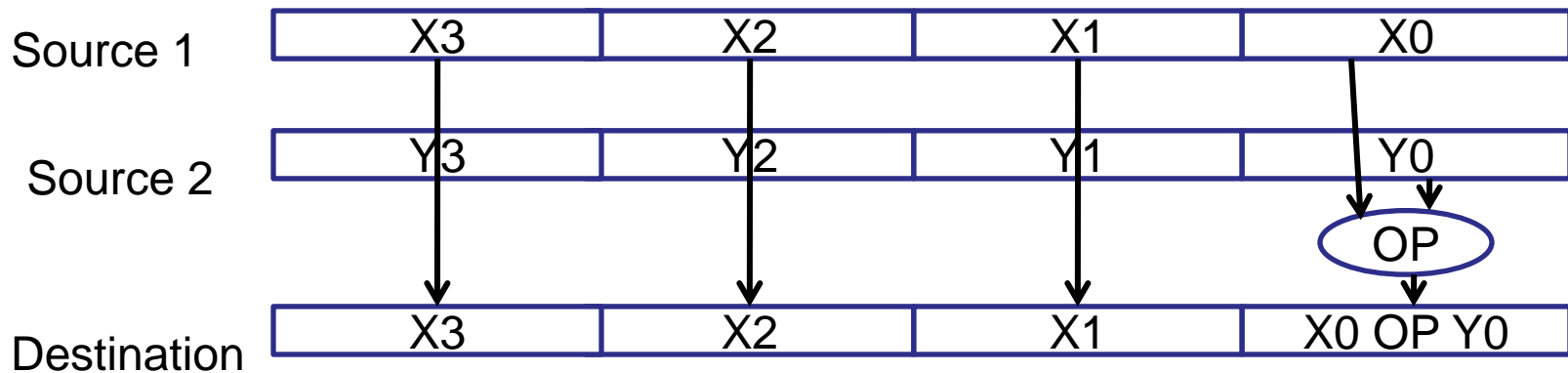
Images are from shi et al.'s "Example-based Dynamic Skinning in Real Time"

- Artists use standard tools to generate a character model a long with a series of key poses
- Model: a set of bones + deformable skins
- Xenon interpolate new poses as needed
- Skins are generated on the fly
- Xenon only sends the vertices that have changed to save bandwidth

Background: Packed and Scalar Floating-Point Instructions

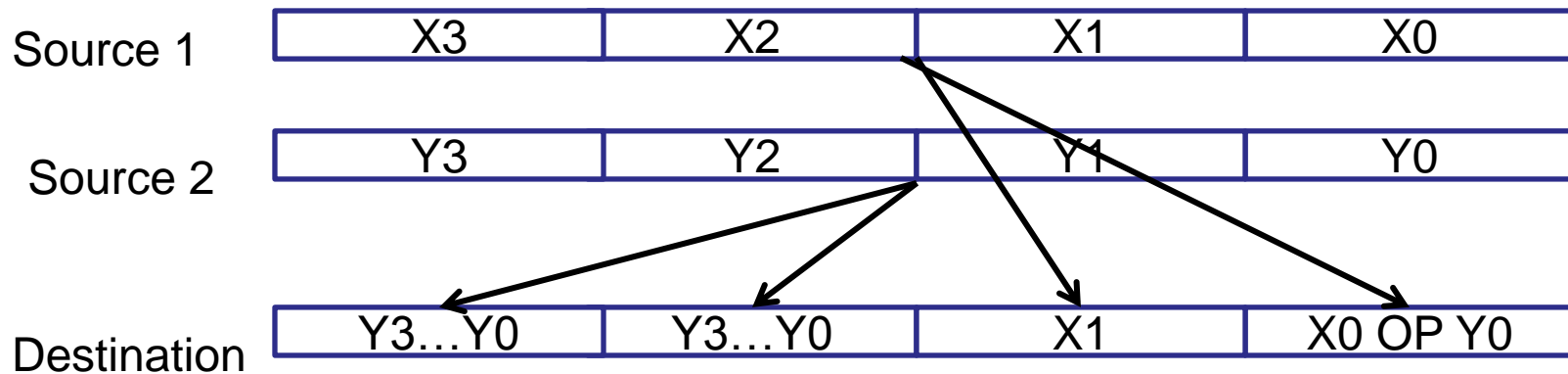


Packed single-precision floating-point operation



Scalar single-precision floating-point operation

Background: Shuffle and Unpack Instructions



Scalar single-precision floating-point operation



SIMD Background: Loop unrolling

```
for (i = 1; i < 12; i++) x[i] = j[i]+1;
```

```
for (i = 1; i < 12; i=i+4)
```

```
{
```

```
    x[i] = j[i]+1;
```

```
    x[i+1] = j[i+1]+1;
```

```
    x[i+2] = j[i+2]+1;
```

```
    x[i+3] = j[i+3]+1;
```

```
}
```

SSE ADD



SIMD Background: Swizzling

- Changing the order of vector elements by calling some operands
- Vector2 foo;
Vector4 bar = Vector4(1.0f, 3.0f, 1.0f, 1.0f);
foo.xy = bar.zw;

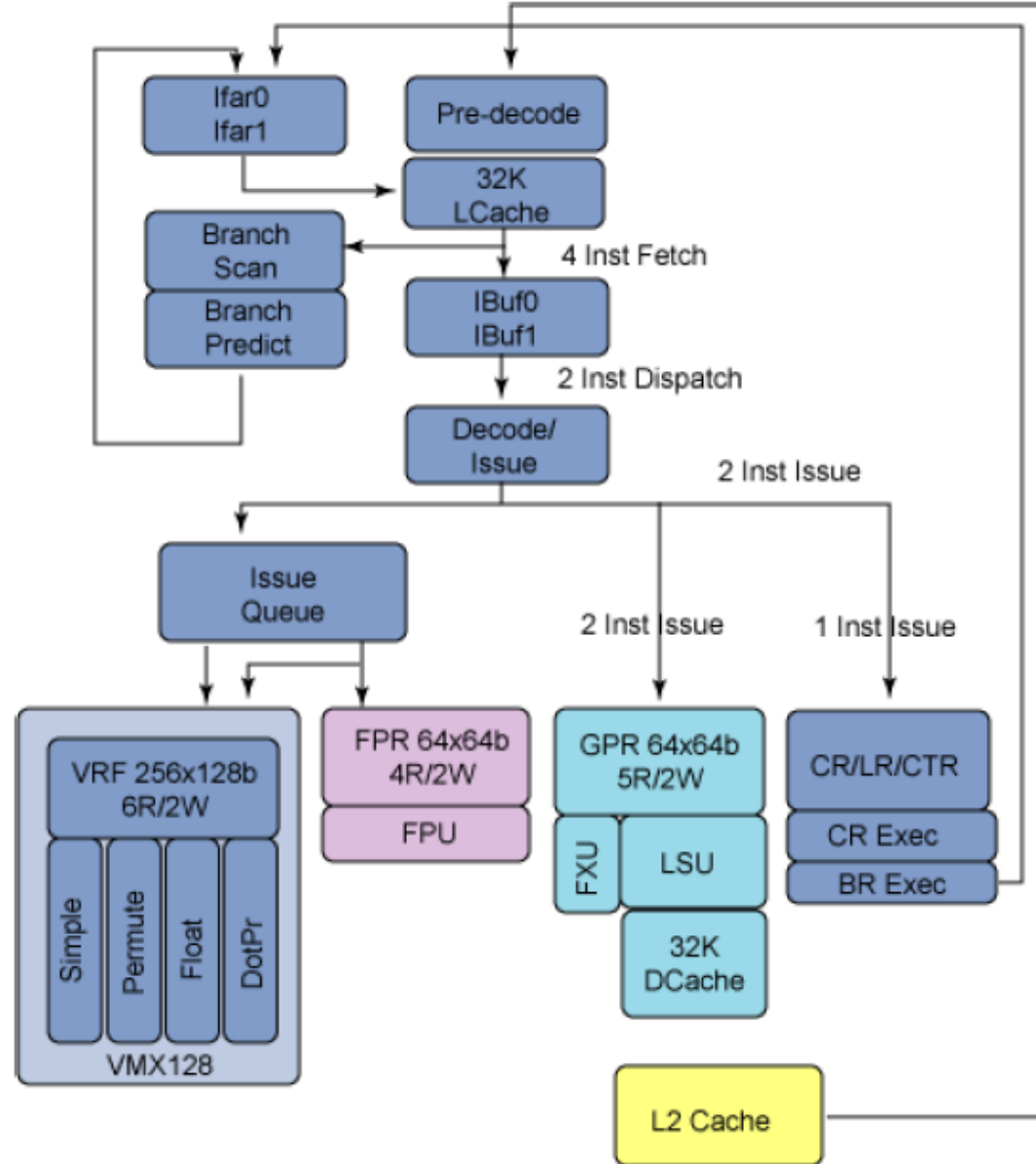


SOA & AOS

- Array of structures (AOS)
 - $\{x_1, y_1, z_1, w_1\}$, $\{x_2, y_2, z_2, w_2\}$, $\{x_3, y_3, z_3, w_3\}$
 , $\{x_4, y_4, z_4, w_4\}$
 - Intuitive but less efficient
 - What if we want to perform only x axis?
- Structure of array (SOA)
 - $\{x_1, x_2, x_3, x_4\}$, ..., $\{y_1, y_2, y_3, y_4\}$, ... $\{z_1, z_2, z_3, z_4\}$,
 ... $\{w_1, w_2, w_3, w_4\}$...
 - Better SIMD unit utilization, better cach
 - Also called “swizzled data”

Data path

- Four-instruction fetch
- Two-instruction “dispatch”
- Five functional units
- “VMX128” execution “decoupled” from other units
- 14-cycle VMX dot-product
- Branch predictor:
- “4K” G-share predictor
- Unclear if 4KB or 4K 2-bit counters
- Per thread





BACKGROUND: G-SHARE BRANCH PREDICTOR



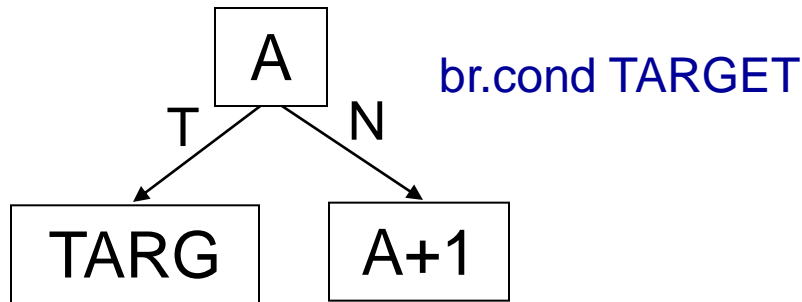
Branches...



- Movement of Kung Fu Panda is dependent on user inputs
- What happened to the previous scenes
- “Branches” in the code makes a decision
- Draw all the motions and new characters after an user input
 - Requires fast computing,
 - May be we can prepare speculatively



Branch Code







- Depending on the direction of branch in basic block A, we have to decide whether we fetch TARG or A+1

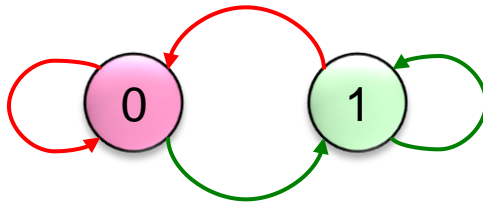


Branches: Prediction

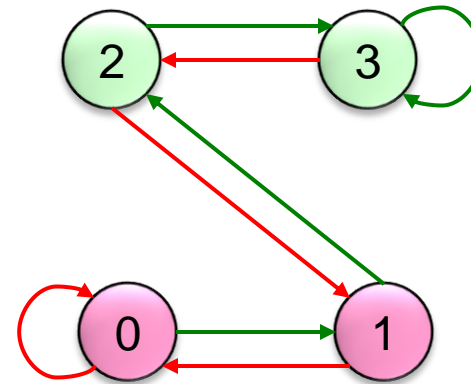
- Predict Branches
 - Predict the next fetch address
- Fetch, decode, etc. on the predicted path
 - Execute anyway (speculation)
- Recover from mispredictions
 - Restart fetch from correct path
- How?
 - Based on old history
- Simple example: last time predictor

Two Bits Counter Based Prediction

-  Predict NT
-  Predict T
-  Transition on T outcome
-  Transition on NT outcome



FSM for Last-time Prediction



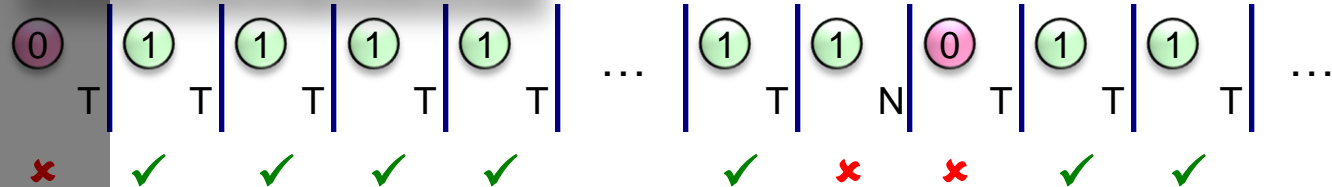
FSM for 2bC
(2-bit Counter)

Example

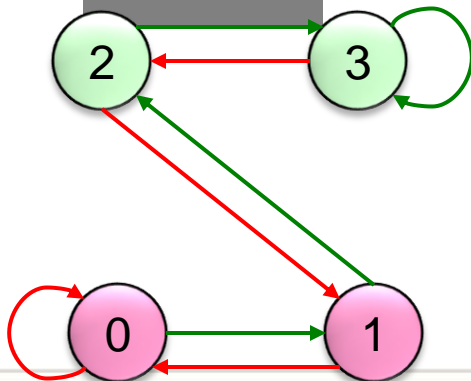
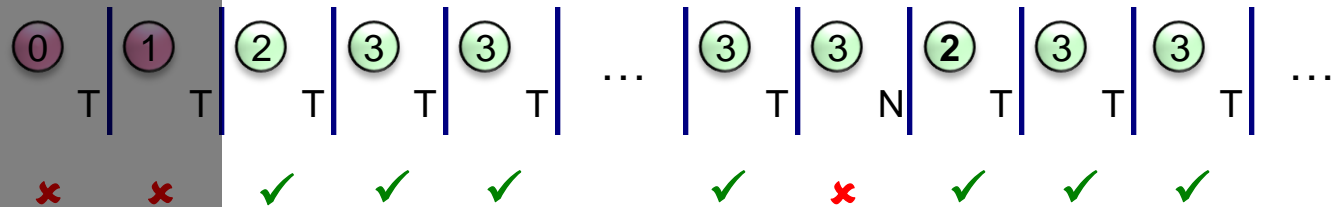


1bC:

Initial Training/Warm-up



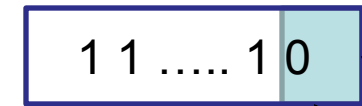
2bC:



Only 1 Mispredict per N branches now!
DC08: 99.999% DC44: 99.0%



Two-level Branch Predictor



previous one

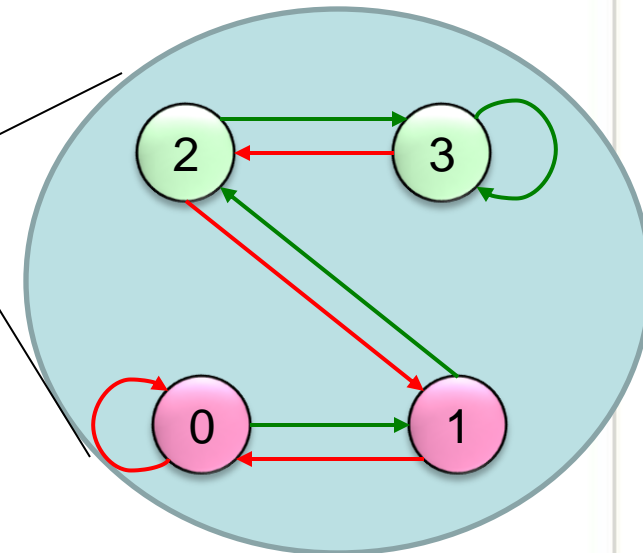
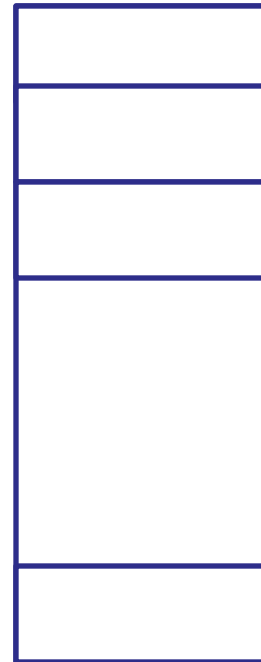
BHR
(branch
history
register)

index

Pattern History Table

00 00
00 01
00 10

11 11



Yeh&patt'92

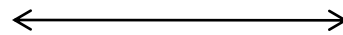
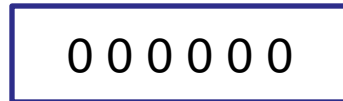


BHR (Branch History Register)

Initialization value (0 or 1)

Old history

New history



History length

1 : branch is taken

0: branch is not-taken

$$\text{New BHR} = \text{old BHR} \ll 1 \mid (\text{br_dir})$$

Example

BHR: 00000

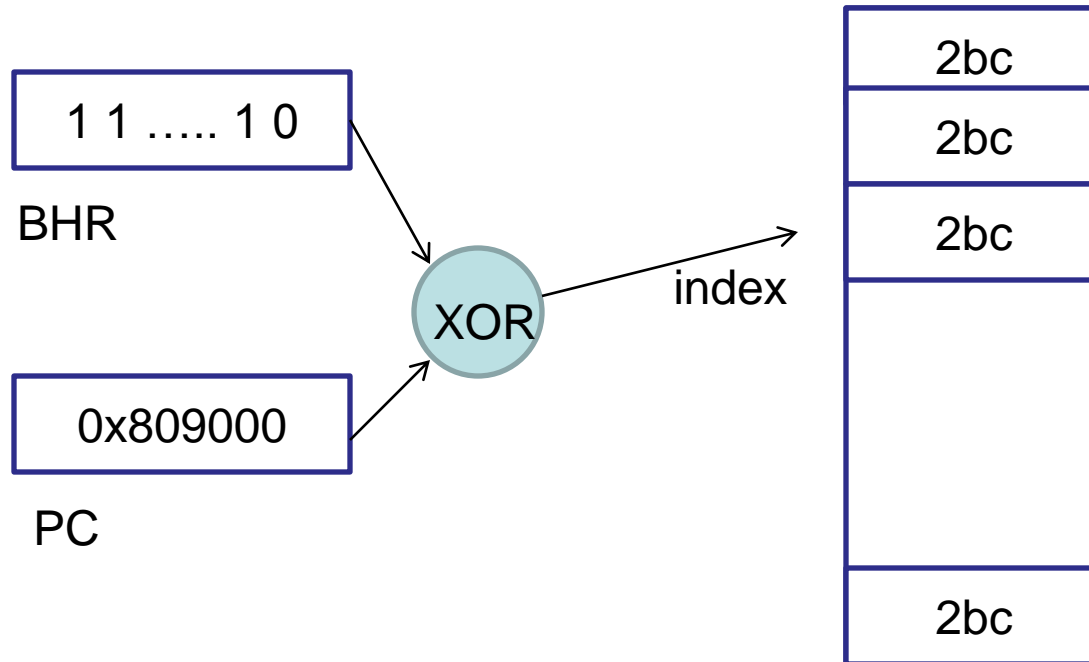
Br1 : taken → BHR 00001

Br 2: not-taken → BHR 00010

Br 3: taken → BHR 00101



Gshare Branch Predictor



McFarling'93

Predictor size: $2^{(\text{history length})} * 2\text{bit}$



Why Branch Predictor Works ?

- Repeated history
 - Could be user actions
 - Many generic regularity in many applications,
- Correlations
 - Panda acquired a new skill it will use it later
 - E.g.
 - If (skill > higher)
 - Panda gets a new fancy knife
 - If (panda has a new fancy knife)
 - draw it. etc..



MEMORY SYSTEM: STREAM OPTIMIZATIONS



Stream Optimizations

- 128B cache line size
- Write streaming:
 - L1s are write through, write misses do not allocate in L1
 - 4 uncacheable write gathering buffers per core
 - 8 cacheable, non-sequential write gathering buffers per core
- Read streaming:
 - 8 outstanding loads/prefetches.
 - xDCBT: Extended data cache block touch, bringing data directly to L1 , never store L2
 - Useful for non-shared data



CPU/GPU

- CPU can send 3D compressed data directly to the GPU w/o cache
- Geometry data
- XPS support:
 - (1): GPU and the FSB for a 128-byte GPU read from the CPU
 - (2) From GPU to the CPU by extending the GPU's tail pointer write-back feature.



Cache-set-locking

- Threads owns a cache sets until the instructions retires.
- Reduce cache contention.
- Common in Embedded systems
- Use L2 cache as a FIFO buffer: sending the data stream into the GPU



Tail Pointer write-back

- Tail pointer write-back: method of controlling communication from the GPU to the CPU by having the CPU poll on a cacheable location, which is updated when a GPU instruction writes an updated to the pointer.
- Free FIFO entry
- System coherency system supports this.
- Reduce latency compared to interrupts.
- Tail pointer backing-store target

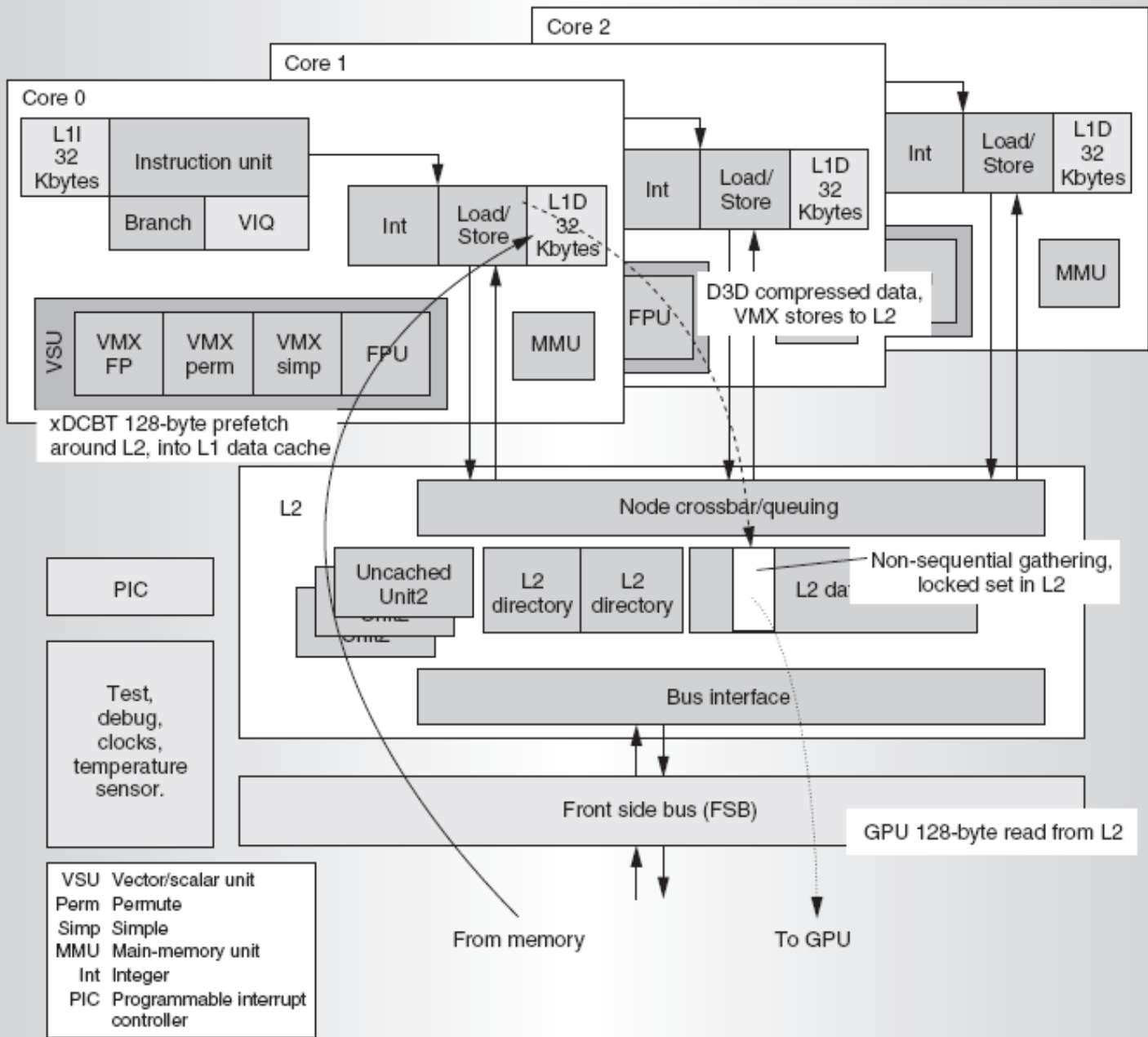
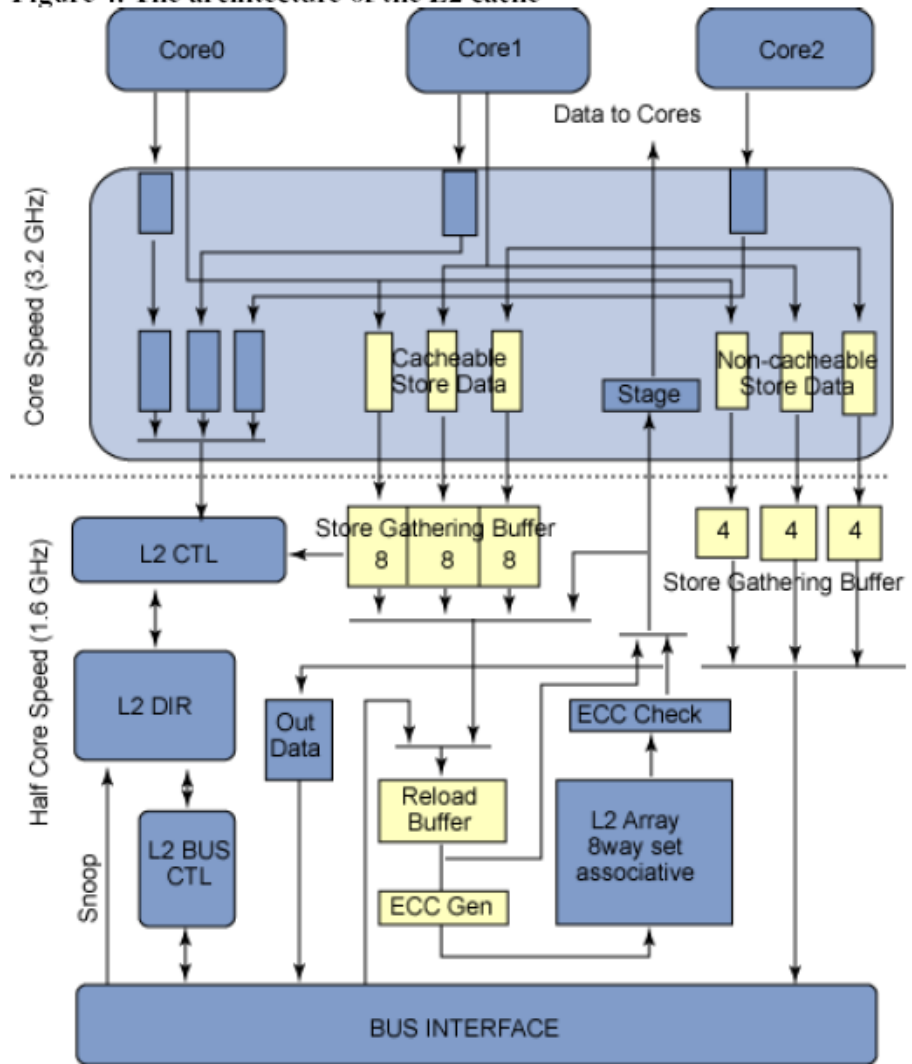


Figure 4. CPU cached data-streaming example.

Memory systems

Figure 4. The architecture of the L2 cache





Non-Blocking Caches

- Hit Under Miss
 - Allow cache hits while one miss in progress
 - But another miss has to wait
- Miss Under Miss, Hit Under Multiple Misses
 - Allow hits and misses when other misses in progress
 - Memory system must allow multiple pending requests
- MSHR (Miss Information/Status Holding Register):
Stores unresolved miss information for each miss that will be handled concurrently.