

**STAM Center**  
SECURE, TRUSTED, AND ASSURED MICROELECTRONICS

**ASU Engineering**  
Arizona State University

# CSE 520

## Computer Architecture II

### Graphic Processing Units

Prof. Michel A. Kinsy

1

---

---

---

---

---

---

---

**STAM Center**  
SECURE, TRUSTED, AND ASSURED MICROELECTRONICS

**ASU Engineering**  
Arizona State University

## Graphics Processing Units

```
graph TD; A[Scene Transformations] --> B[Lighting & Shading]; B --> C[Viewing Transformations]; C --> D[Rasterization]
```

2

---

---

---

---

---

---

---

**STAM Center**  
SECURE, TRUSTED, AND ASSURED MICROELECTRONICS

**ASU Engineering**  
Arizona State University

## Graphics Processing Units

```
graph TD; A[Geometry data] --- B[Transform & lighting]; B --- C["Culling, perspective divide, viewport mapping"]; C --- D[Rasterization]; D --- E[Simple texturing]; E --- F[Depth test]; F --- G[Frame buffer blending];
```

3

---

---

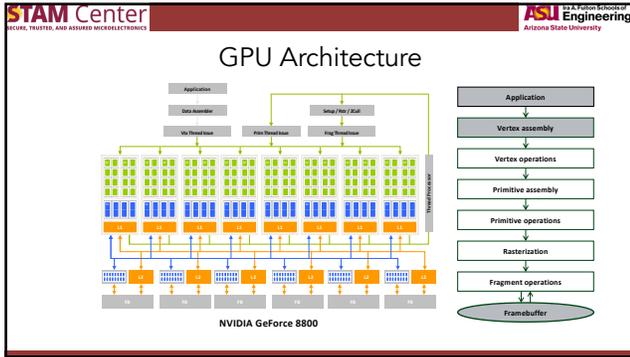
---

---

---

---

---



4

---

---

---

---

---

---

---

---

- ### GPU Computing
- Most successful commodity accelerator
  - GPUs combine two useful strategies to increase efficiency
    - Massive parallelism
    - Specialization
  - Illustrates tension between performance and programmability in accelerators

5

---

---

---

---

---

---

---

---

- ### Graphics Processors Timeline
- Till mid-90s
    - VGA controllers used to accelerate some display functions
  - Mid-90s to mid-2000s
    - Fixed-function accelerators for the OpenGL and DirectX APIs
    - 3D graphics: triangle setup & rasterization, texture mapping & shading

6

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Graphics Processors Timeline

- Modern GPUs
  - Programmable multiprocessors optimized for data-parallelism
    - OpenGL/DirectX and general purpose languages (CUDA, OpenCL, ...)
  - Some fixed-function hardware (texture, raster ops, ...)
  - GPUs in Modern Systems
    - Discrete GPUs
      - PCIe-based accelerator
      - Separate GPU memory
    - Integrated GPUs
      - CPU and GPU on same die
      - Shared main memory and last-level cache

7

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### GPU Architecture

- Warp
  - A set of threads that execute the same instruction on different data elements
  - All threads run the same kernel

8

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### GPU Architecture

- Fine-grained multithreading
  - One instruction per thread in pipeline at a time (No branch prediction)
  - Interleave warp execution to hide latencies
- Register values of all threads stay in register file
- No OS context switching
- Memory latency hiding
  - Graphics has millions of pixels

9

---

---

---

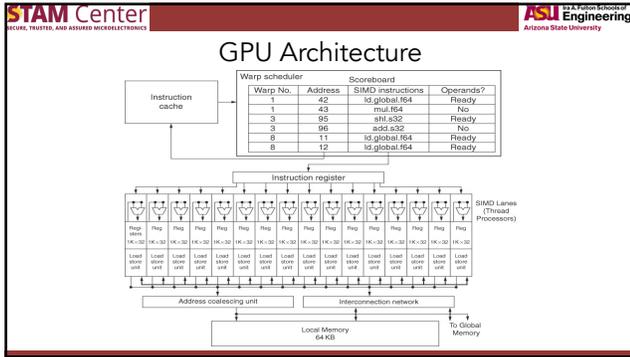
---

---

---

---

---



10

---

---

---

---

---

---

---

---

**GPU Architecture**

- Similarities to vector machines
  - Works well with data-level parallel problems
  - Scatter-gather transfers
  - Mask registers
  - Large register files
- Differences:
  - No scalar processor
  - Uses multithreading to hide memory latency
  - Has many functional units, as opposed to a few deeply pipelined units like a vector processor

11

---

---

---

---

---

---

---

---

**GPU Architecture**

- Each thread has local memory
- Parallel threads packed in warps
  - Access to per-warp shared memory
  - Can synchronize with barrier
- Grids include independent warps
  - May execute concurrently

The diagram shows a thread connected to per-thread local memory. A warp is a group of threads that can access per-block shared memory. Multiple warps are organized into a grid, which can execute multiple kernels (Kernel 0, Kernel 1) sequentially. The grid is connected to per-device global memory.

12

---

---

---

---

---

---

---

---



**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### GPU ISA and Compilation

- GPU microarchitecture and instruction set change very frequently
- To achieve compatibility:
  - Compiler produces intermediate pseudo-assembler language (e.g., Nvidia PTX)
  - GPU driver JITs kernel, tailoring it to specific microarchitecture
- In practice, little performance portability
  - Code is often tuned to specific GPU architecture

16

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Thread Scheduling

- The scheduler selects a ready thread of SIMD instructions and issues an instruction synchronously to all the SIMD Lanes executing the SIMD thread
- Because threads of SIMD instructions are independent, the scheduler can select a different SIMD thread each time

The diagram illustrates the thread scheduling process. A 'SIMD thread scheduler' box at the top has an arrow pointing to a vertical timeline labeled 'Time'. The timeline shows a sequence of instruction blocks: 'SIMD thread 8 instruction 11', 'SIMD thread 1 instruction 42', 'SIMD thread 3 instruction 95', an ellipsis, 'SIMD thread 8 instruction 12', 'SIMD thread 3 instruction 96', and 'SIMD thread 1 instruction 43'. Each instruction block is represented by a horizontal bar with vertical tick marks underneath, indicating the execution of multiple SIMD lanes.

17

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Instruction & Thread Scheduling

- In theory, all threads can be independent
- For efficiency, 32 threads packed in warps
  - Warp: set of parallel threads that execute the same instruction
    - Warp = a thread of vector instructions
    - Warps introduce data parallelism
  - 1 warp instruction keeps cores busy for multiple cycles

18

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Instruction & Thread Scheduling

- Individual threads may be inactive
  - Because they branched differently
  - This is the equivalent of conditional execution (but implicit)
  - Loss of efficiency if not data parallel
- Software thread blocks mapped to warps
  - When hardware resources are available

19

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Handling Branch Divergence

- Similar to vector processors, but masks are handled internally
  - Per-warp stack stores PCs and masks of non-taken paths
- On a conditional branch
  - Push the current mask onto the stack
  - Push the mask and PC for the non-taken path
  - Set the mask for the taken path
- At the end of the taken path
  - Pop mask and PC for the non-taken path and execute

20

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Handling Branch Divergence

- At the end of the non-taken path
  - Pop the original mask before the branch instruction
- If a mask is all zeros, skip the block
- Dynamic Warp Formation
  - Dynamically merge threads executing the same instruction (after branch divergence)
  - Form new warp at divergence
  - Enough threads branching to each path to create full new warps

21

---

---

---

---

---

---

---

---



**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### GPU Execution Model

- Data transfers can dominate execution time
- Integrated GPUs with unified address space
  - No copies

- Transfer input data from CPU to GPU memory
- Launch kernel (grid)
- Wait for kernel to finish (if synchronous)
- Transfer results to CPU memory

The diagram shows a CPU and a GPU connected to a shared memory space. Arrows indicate data flow: from CPU to Memory, from Memory to GPU, from GPU to Memory, and from Memory to CPU. A bidirectional arrow also connects the CPU and GPU.

25

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Program Execution

- Same instruction in different threads uses thread id to index and access different data elements

```
GPU code
for (int i = 0; i < 100; i++) {
    C[i] = A[i] + B[i];
}
```

26

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Program Execution

- Same instruction in different threads uses thread id to index and access different data elements

```
GPU code
for (int i = 0; i < 100; i++) {
    C[i] = A[i] + B[i];
}
```

↓

```
CUDA code
// there are 100 threads
__global__ void KernelFunction() {
    int tid = blockDim.x * blockIdx.x + threadIdx.x;
    int varA = aa[tid];
    int varB = bb[tid];
    C[tid] = varA + varB;
}
```

27

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Loop-Level Parallelism

- Focuses on determining whether data accesses in later iterations are dependent on data values produced in earlier iterations
  - Loop-carried dependence
- Example 1
 

```
for (i=999; i>=0; i=i-1)
  x[i] = x[i] + s;
```
- No loop-carried dependence

28

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Loop-Level Parallelism

- Example 2
 

```
for (i=0; i<100; i=i+1) {
  A[i+1] = A[i] + C[i]; /* S1 */
  B[i+1] = B[i] + A[i+1]; /* S2 */
}
```
- S1 and S2 use values computed by S1 in previous iteration
- S2 uses value computed by S1 in same iteration

29

---

---

---

---

---

---

---

---

**STAM Center** SECURE, TRUSTED, AND ASSURED MICROELECTRONICS **ASU Engineering** The Future School of Arizona State University

### Loop-Level Parallelism

- Example 3
 

```
for (i=0; i<100; i=i+1) {
  A[i] = A[i] + B[i]; /* S1 */
  B[i+1] = C[i] + D[i]; /* S2 */
}
```
- S1 uses value computed by S2 in previous iteration but dependence is not circular so loop is parallel
- Transform to:
 

```
A[0] = A[0] + B[0];
for (i=0; i<99; i=i+1) {
  B[i+1] = C[i] + D[i];
  A[i+1] = A[i+1] + B[i+1];
}
B[100] = C[99] + D[99];
```

30

---

---

---

---

---

---

---

---

**STAM Center**  
SECURE, TRUSTED, AND ASSURED MICROELECTRONICS

**ASU Engineering**  
Arizona State University

### Loop-Level Parallelism

- Example 4

```
for (i=0; i<100; i=i+1) {  
  A[i] = B[i] + C[i];  
  D[i] = A[i] * E[i];  
}
```

  - No loop-carried dependence
- Example 5

```
for (i=1; i<100; i=i+1) {  
  Y[i] = Y[i-1] + Y[i];  
}
```

  - Loop-carried dependence in the form of recurrence

---

---

---

---

---

---

---

---

31

**STAM Center**  
SECURE, TRUSTED, AND ASSURED MICROELECTRONICS

**ASU Engineering**  
Arizona State University

### Next Class

- Memory Organization

---

---

---

---

---

---

---

---

32