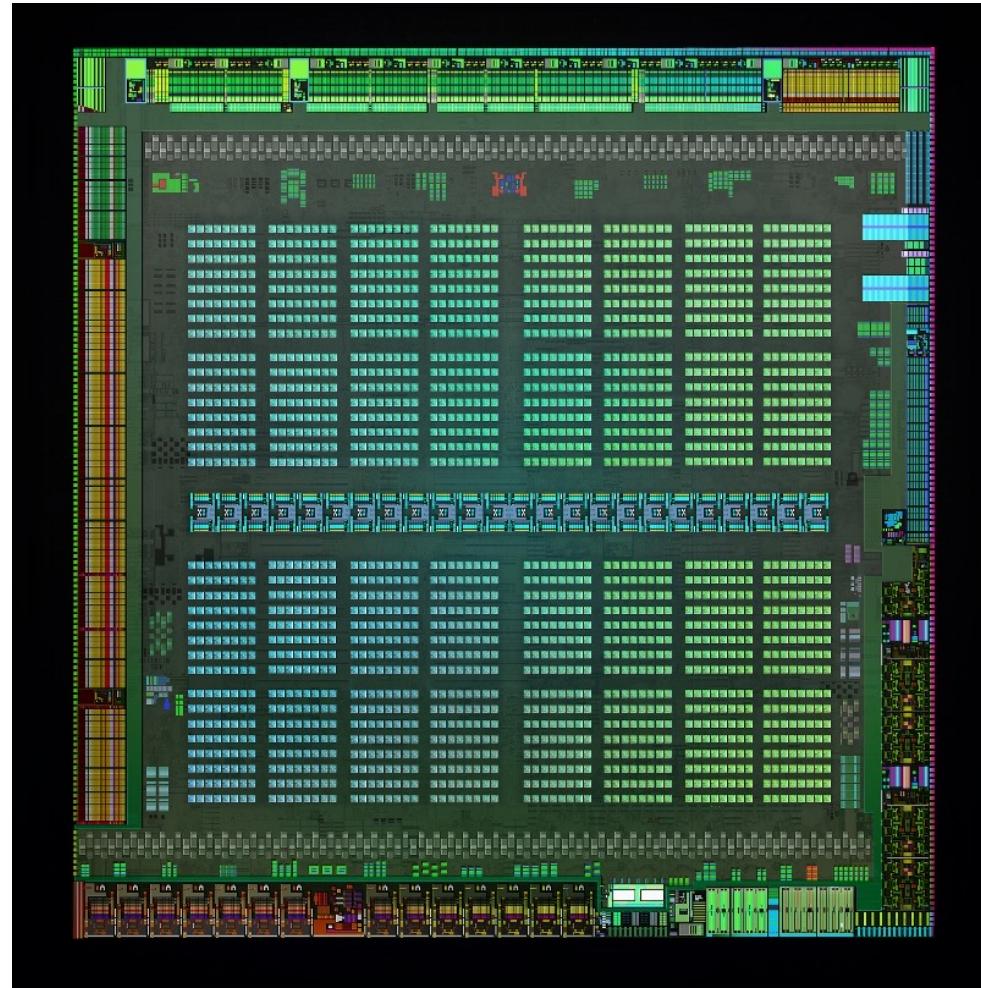


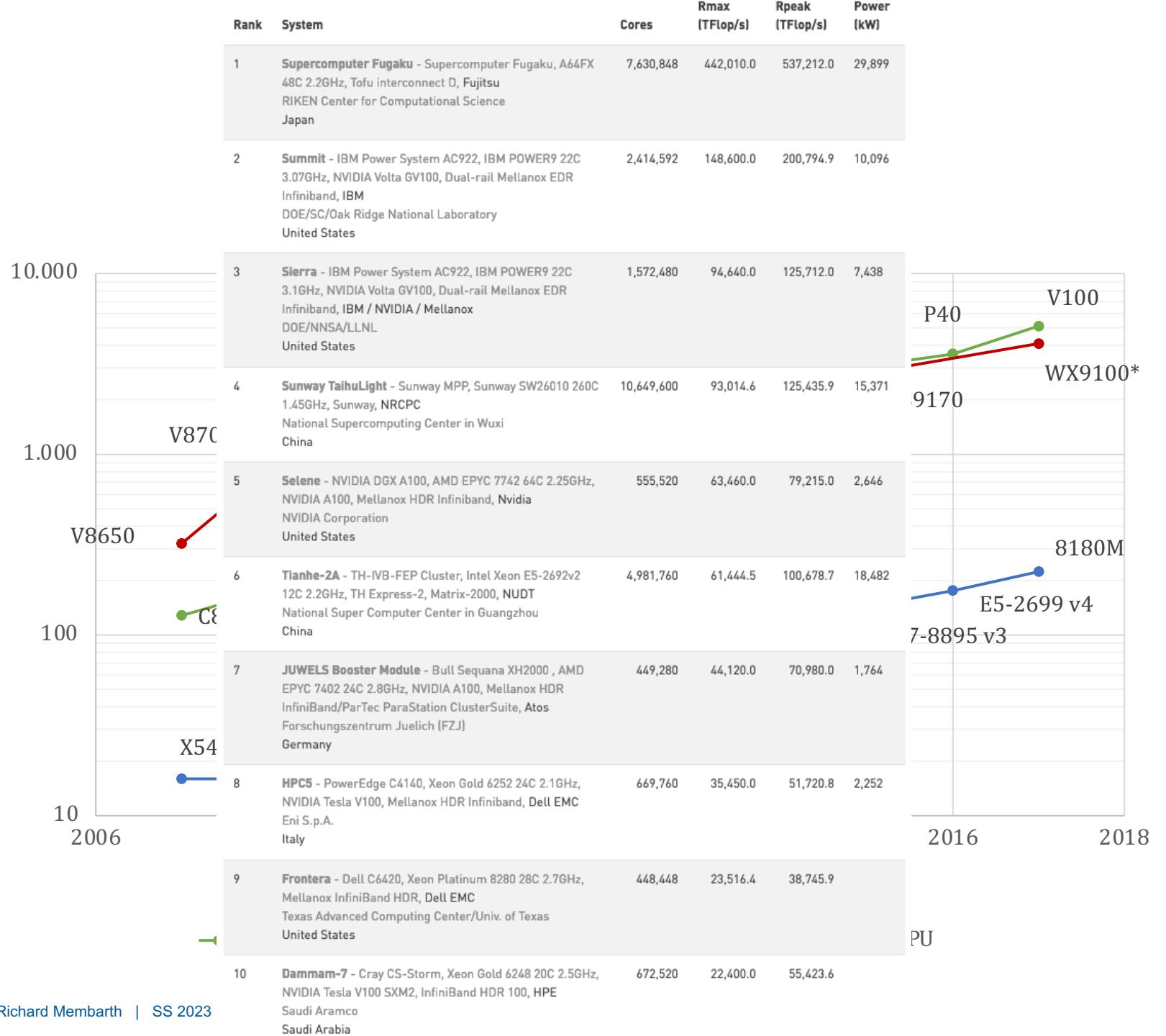
CUDA Programming Model



- Basic understanding of GPU hardware architecture and its difference to the CPU
- Ability to write simple parallel programs that run on the GPU
- Understanding of the different GPU memory types and the ability to use them in a simple program

- Compute Unified Device Architecture







- Given
 - a computation that takes time T_{seq} to complete sequentially
 - N processors
- How much faster can we go?

Speedup
$$S = \frac{T_{seq}}{T_{par}}$$



- Not every operation can be parallelized
- Speedup limited by
 - How much of our program can be parallelized
 - How many processors we can use to run code in parallel
- Assume:
 - p fraction of our program that is parallelizable
 - $s = 1 - p$ sequential fraction



- Purely sequential:

$$T_1 = (1 - p)T_{seq} + pT_{seq}$$

- Using N processors:

$$T_N = (1 - p)T_{seq} + \frac{p}{N}T_{seq}$$

$$S = \frac{T_1}{T_N} = \frac{1}{(1 - p) + \frac{p}{N}} = \frac{1}{s + \frac{p}{N}}$$

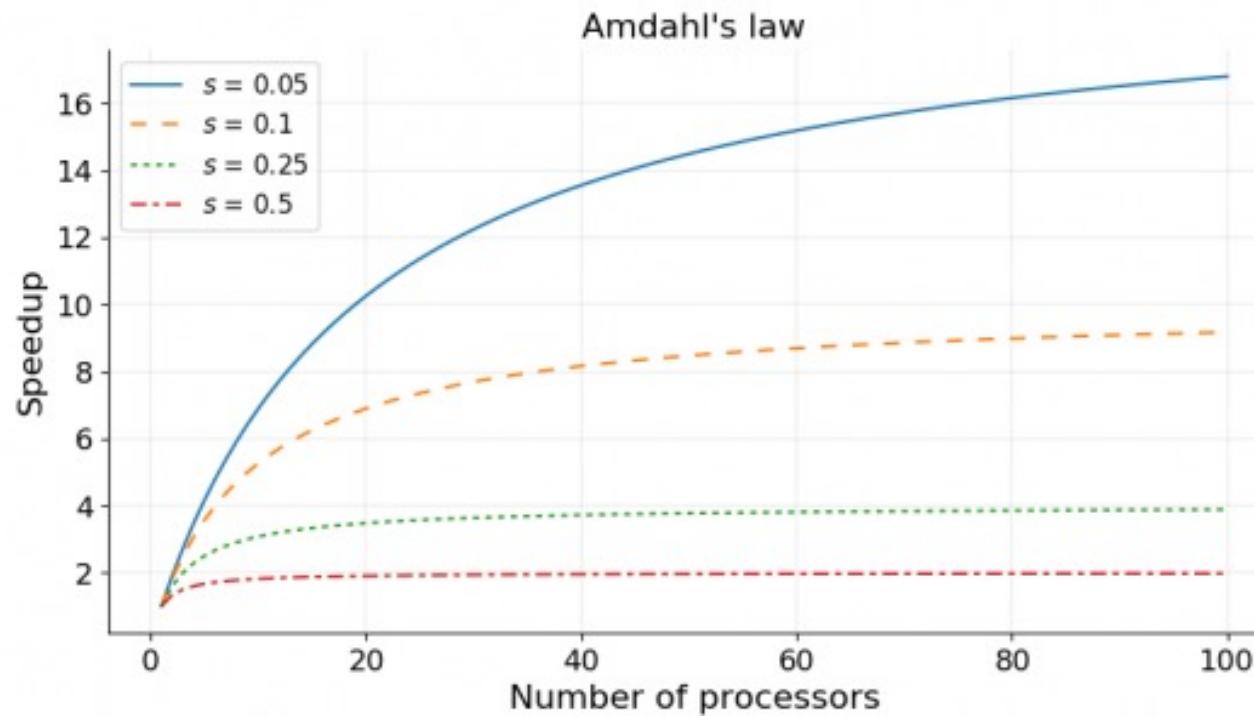
- Amdahl's law from 1967
- Speedup S over purely sequential implementation

$$S = \frac{1}{s + \frac{p}{N}}$$

s : fraction of sequential code

p : fraction of parallel code

N : # of processors



- Strong scaling
 - Limited by sequential fraction
 - Problem size fixed \rightarrow scale processors

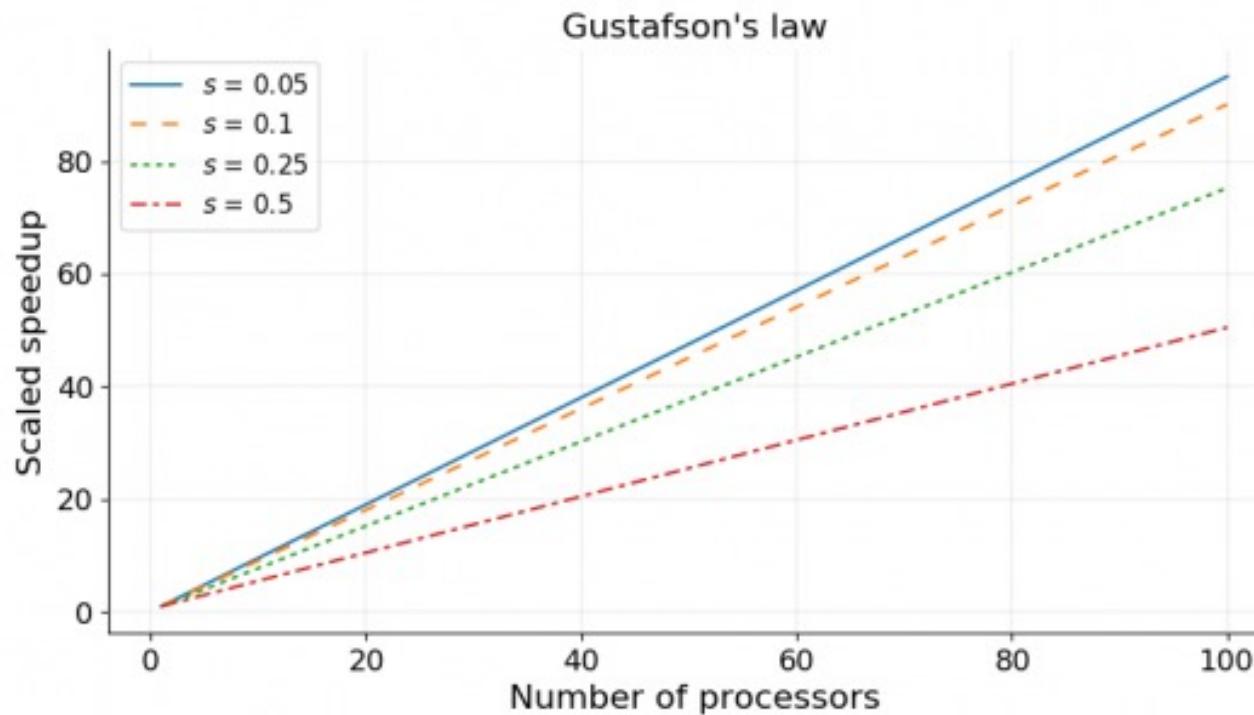
- Gustafson's law from 1988
- Speedup S over purely sequential implementation

$$S = s + p \cdot N$$

s : fraction of sequential code

p : fraction of parallel code

N : # of processors



- Weak scaling
- Sequential fraction less important
- Problem size scales with processors



Task Parallelism

- Parallelize different, independent computation
- Distribute tasks to processors
- Examples: multitasking, pipeline parallelism

Data Parallelism

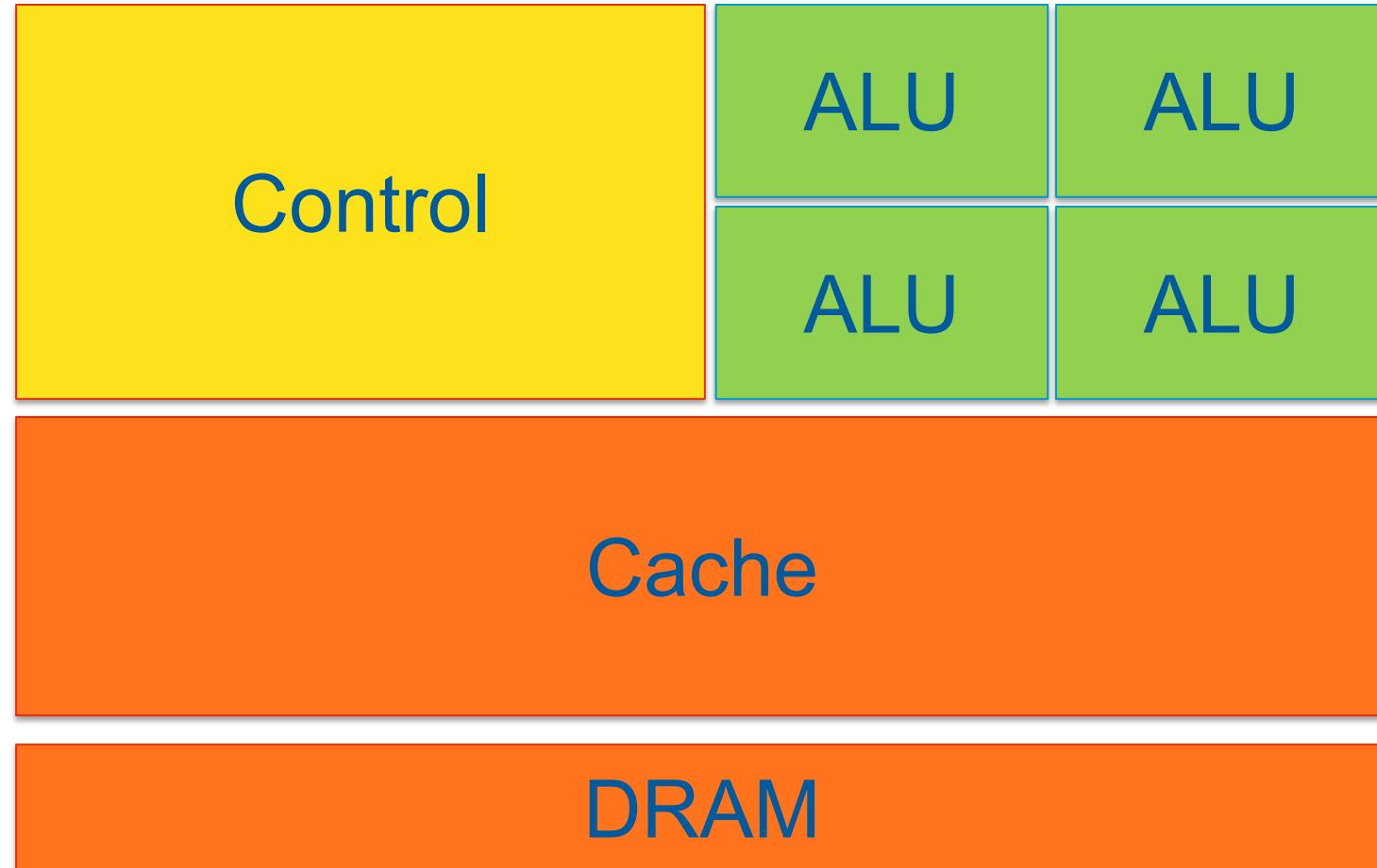
- Parallelize same computation on different, independent data
- Distribute data to processors
- Examples: image processing, loop-level parallelism, divide and conquer



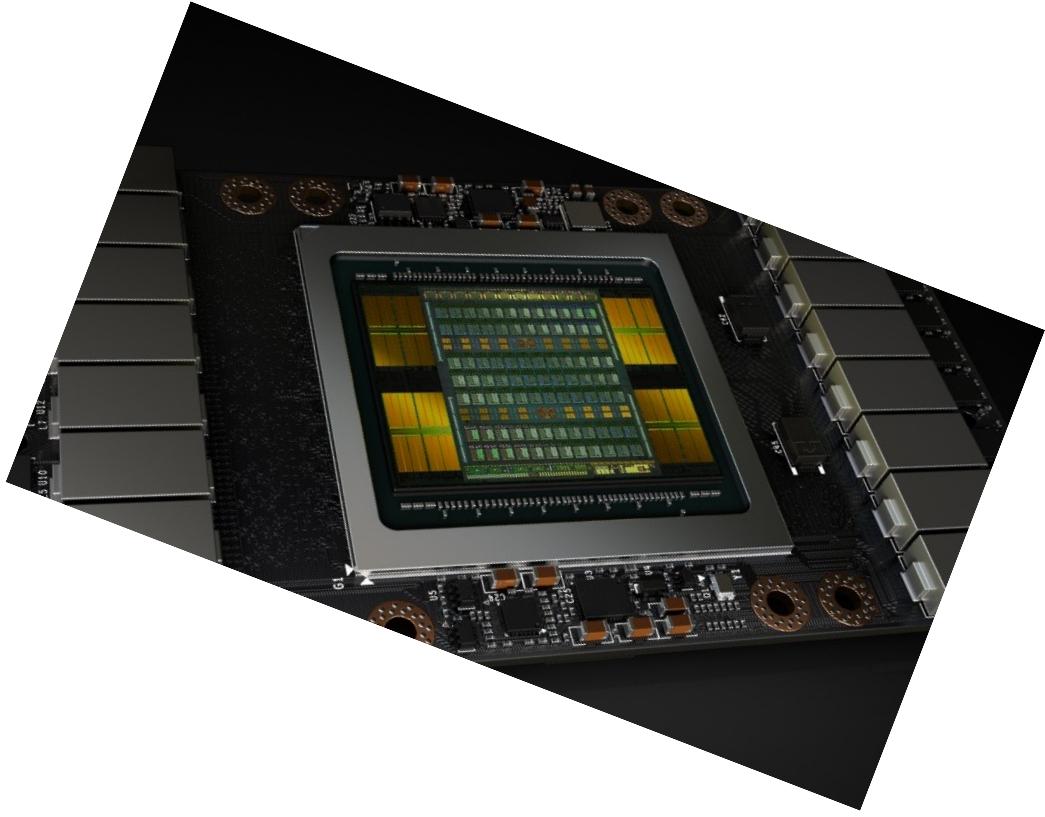
GPU Architecture

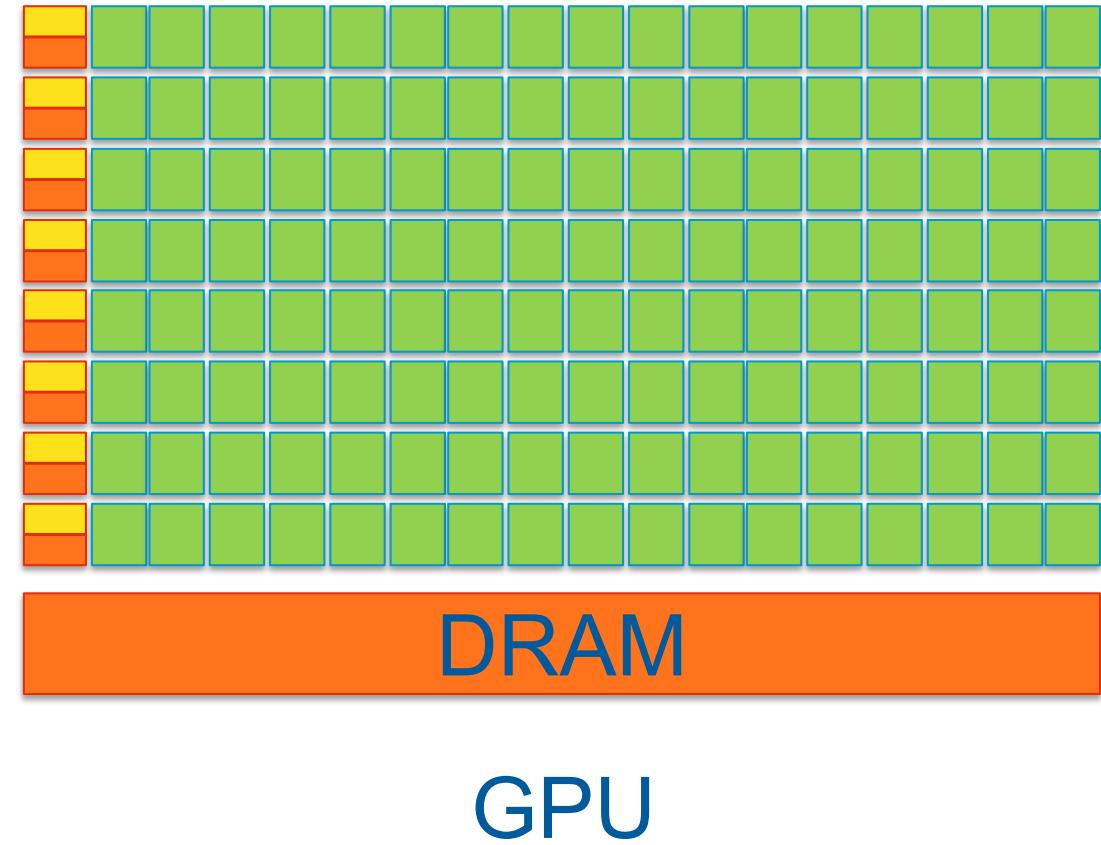
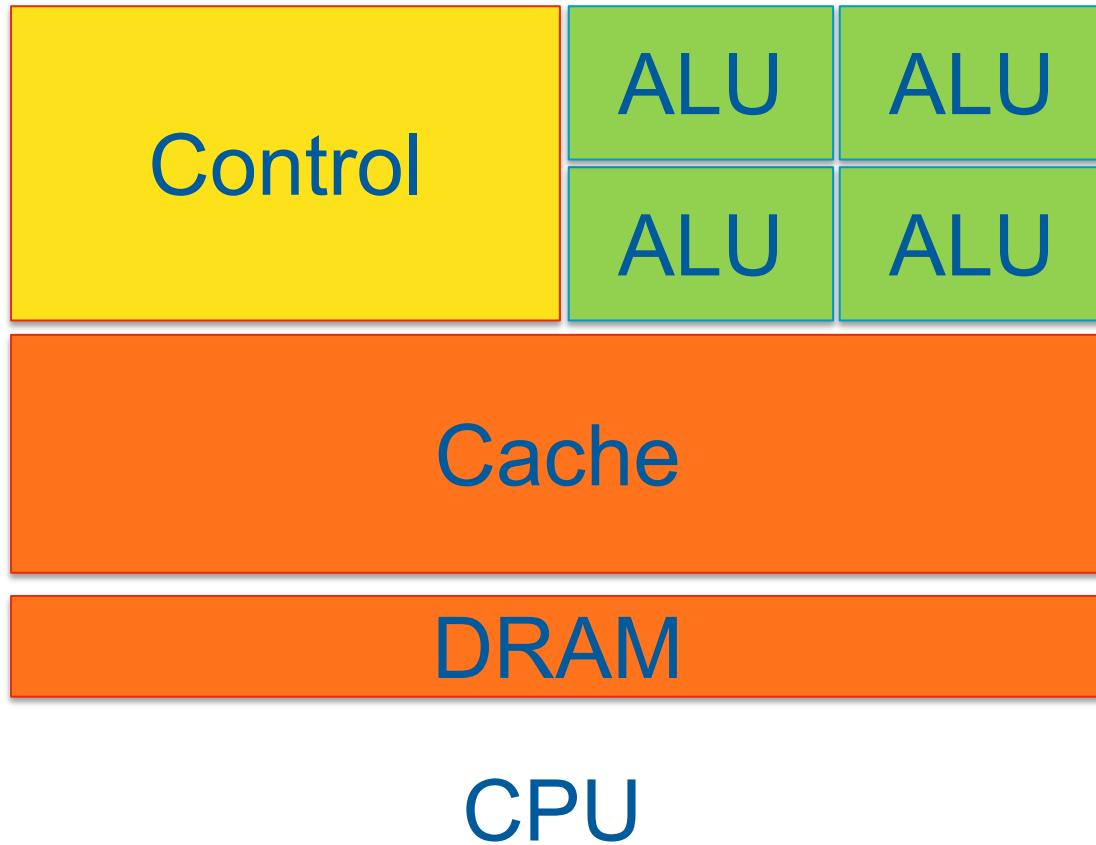
- Goal
 - Get the best performance for a single heavy-weight thread
 - Big data caches
 - Low-latency arithmetic units
 - Complex control logic
 - Branch prediction
 - Out-of-order execution
- latency-oriented design





- Goal
 - Get the best performance for thousands of simple threads
- Small caches
- Hide latency with computation
- In-order execution with no branch prediction
- Issue the same command to multiple cores
→ throughput-oriented design







- SISD single-instruction, single-data
(e.g., single core CPU)
- MIMD multiple-instruction, multiple-data
(e.g., multi core CPU)
- ~~SIMD~~
~~SIMT~~ single-instruction, multiple-data
(e.g., data-based parallelism)
- MISD multiple-instruction, single-data
(e.g., fault-tolerant computers)

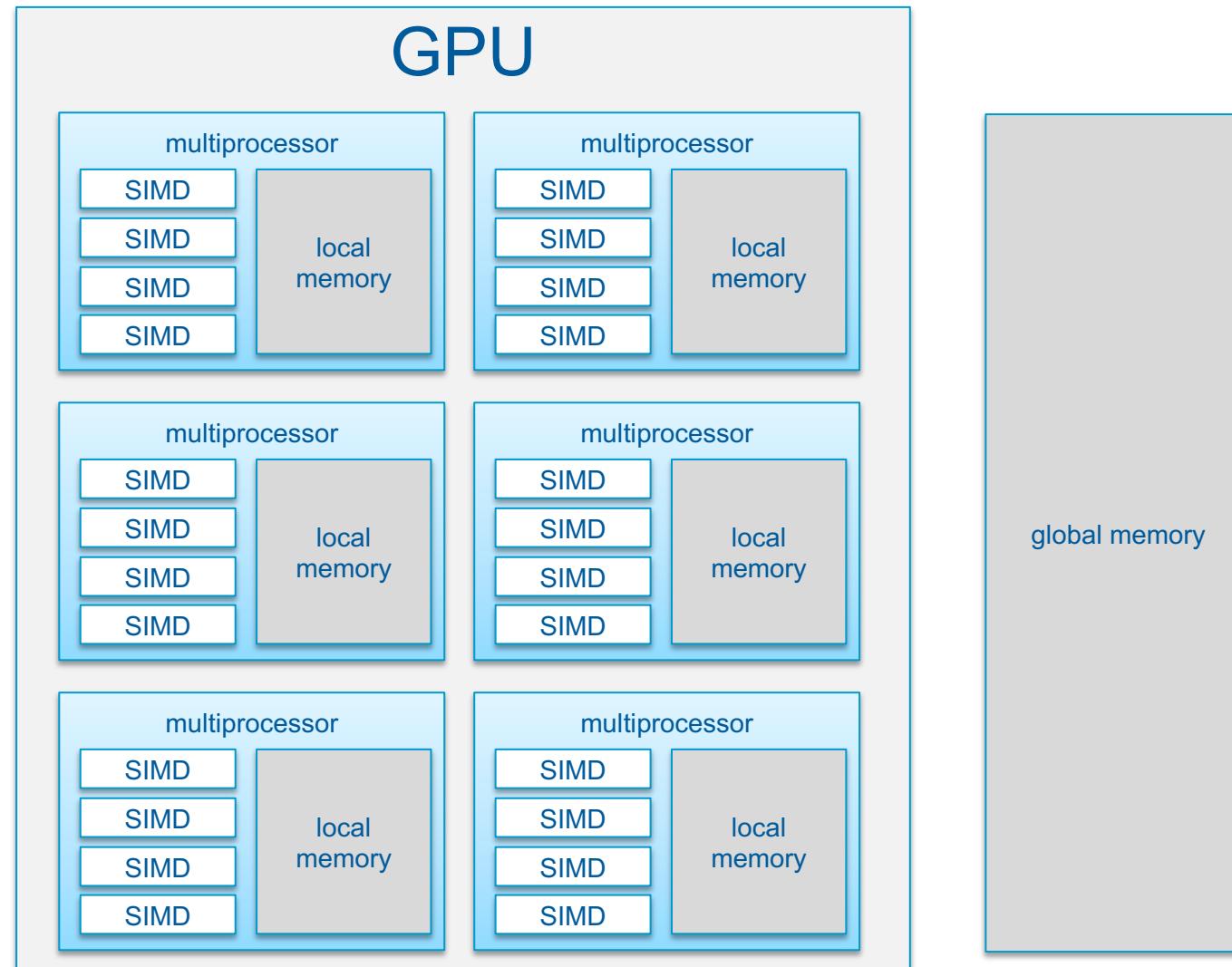


Both execute a single instruction on multiple processing elements

- SIMD exposes SIMD width to software
- SIMT exposes thread-level
 - Write code for independent, scalar threads
 - Threads can diverge within SIMD width

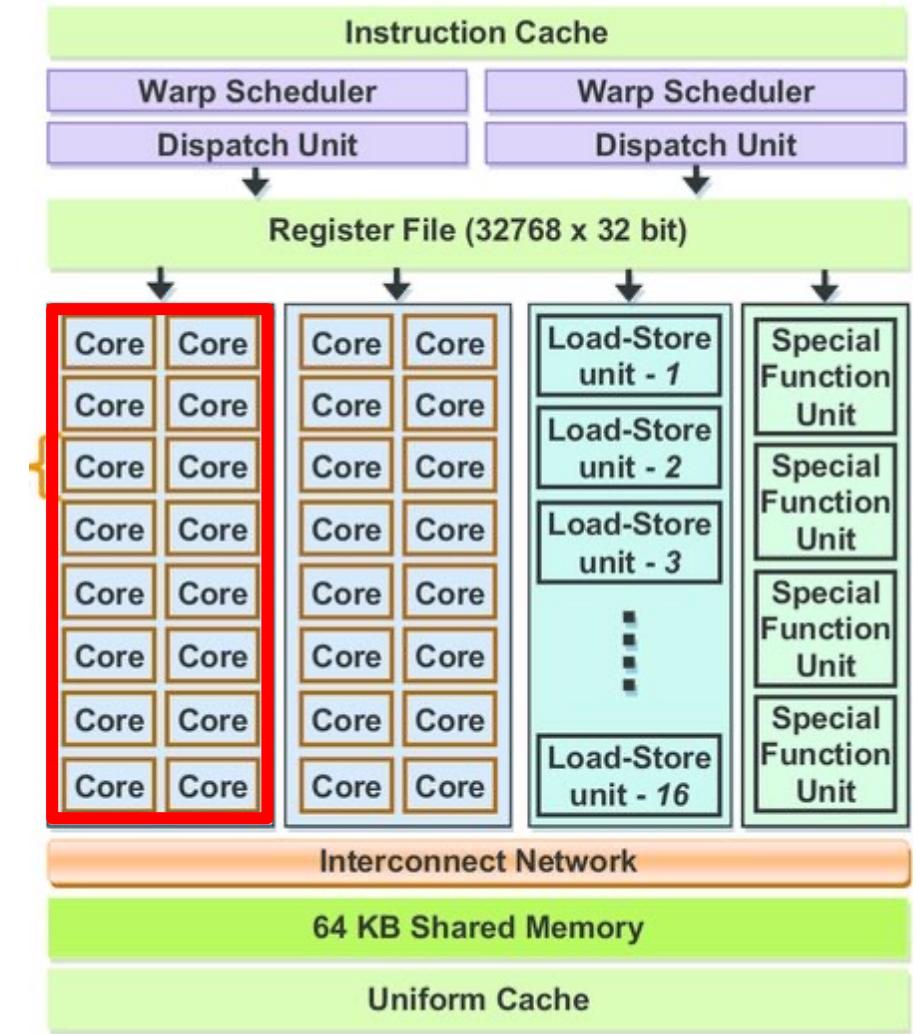
CUDA GPU (Volta)





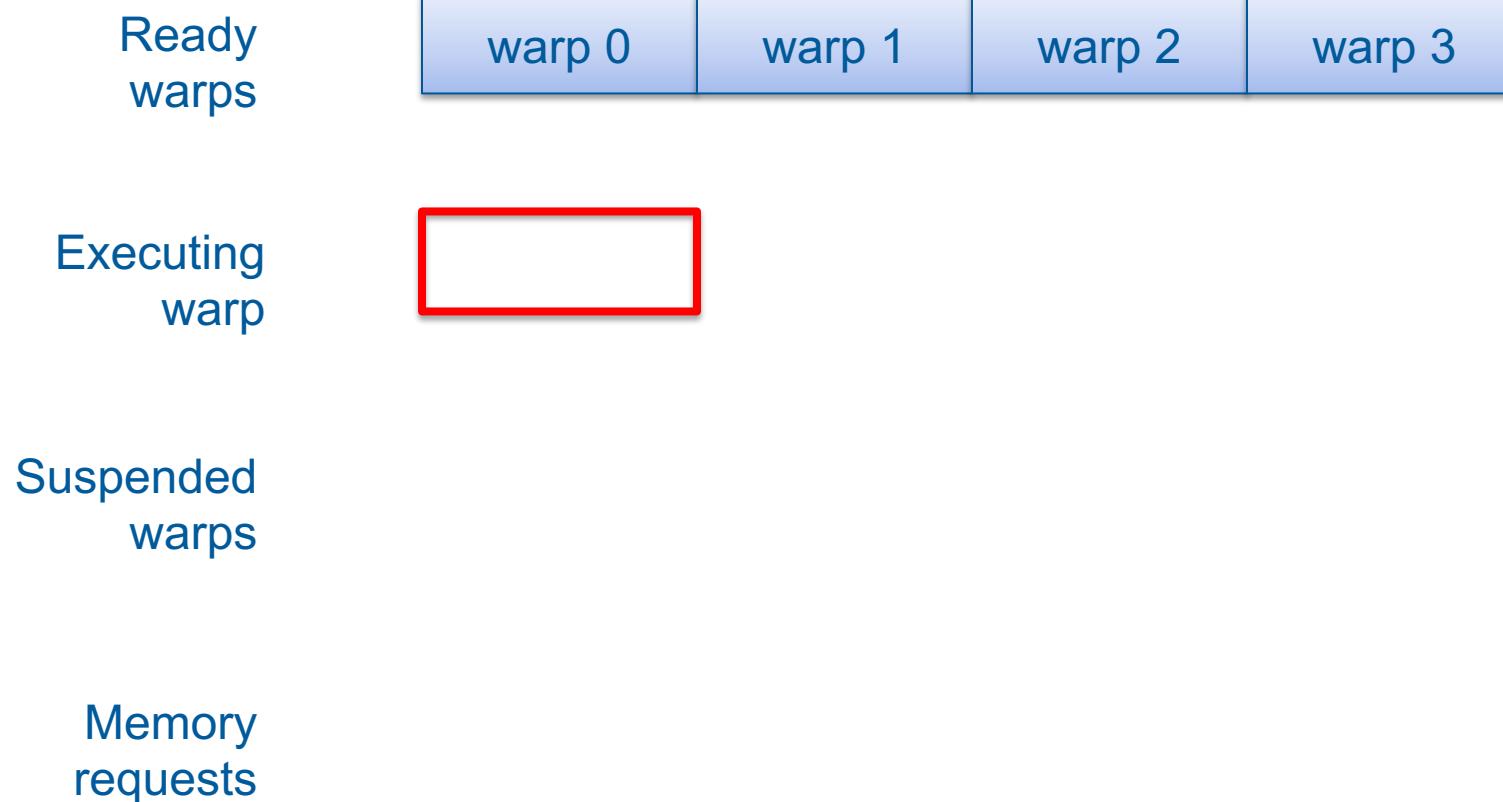


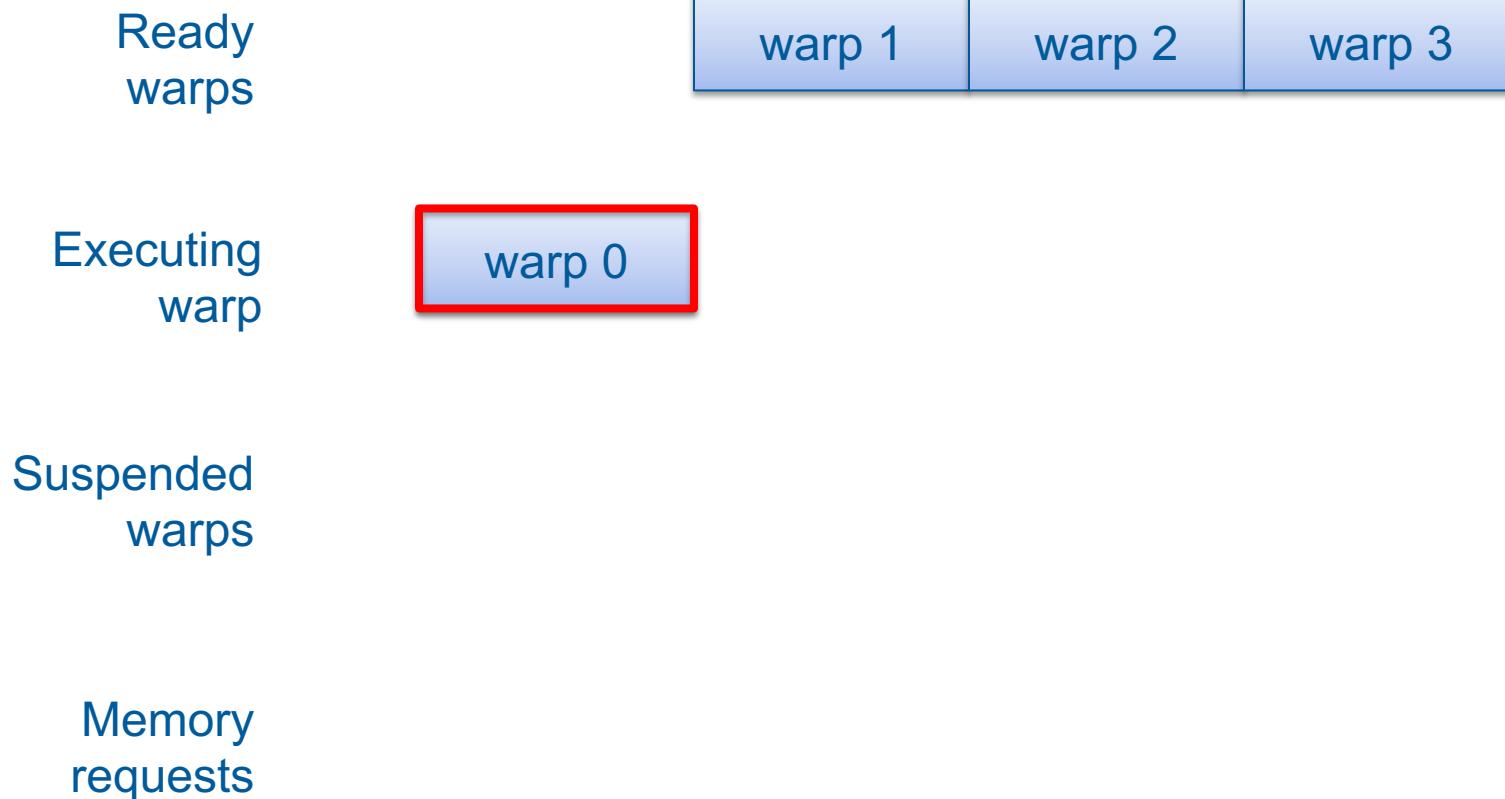
- SIMD hardware
- 32 threads form a *warp*
- Each thread within a warp must execute the same instruction (or be deactivated)
- **1 instruction → 32 values computed**
- Instruction / memory latency

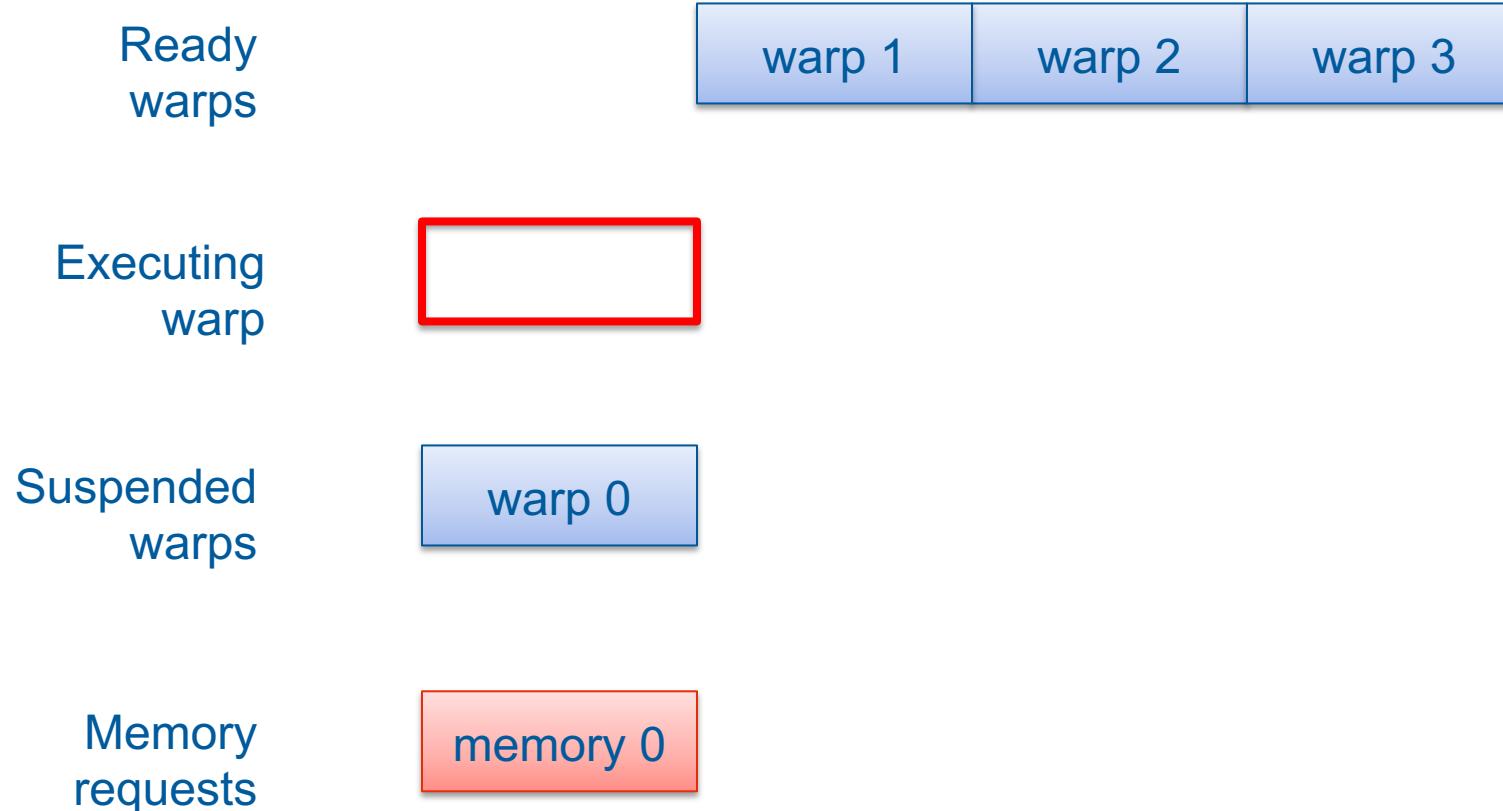


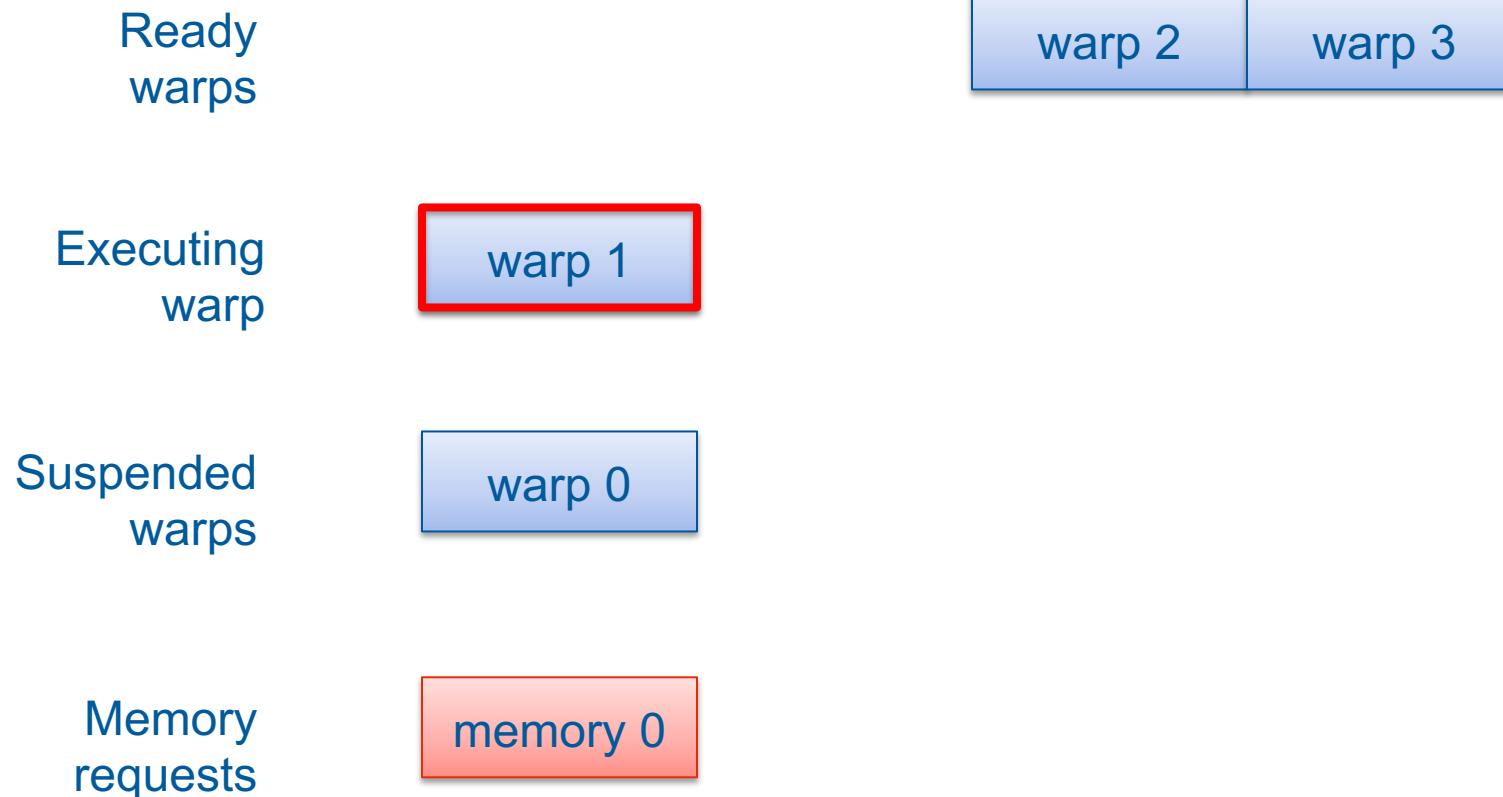
	Volta (in cycles)
L1 hit	28
L1 miss – L2 hit	193
L2 miss – TLB hit	375
TLB miss	1029

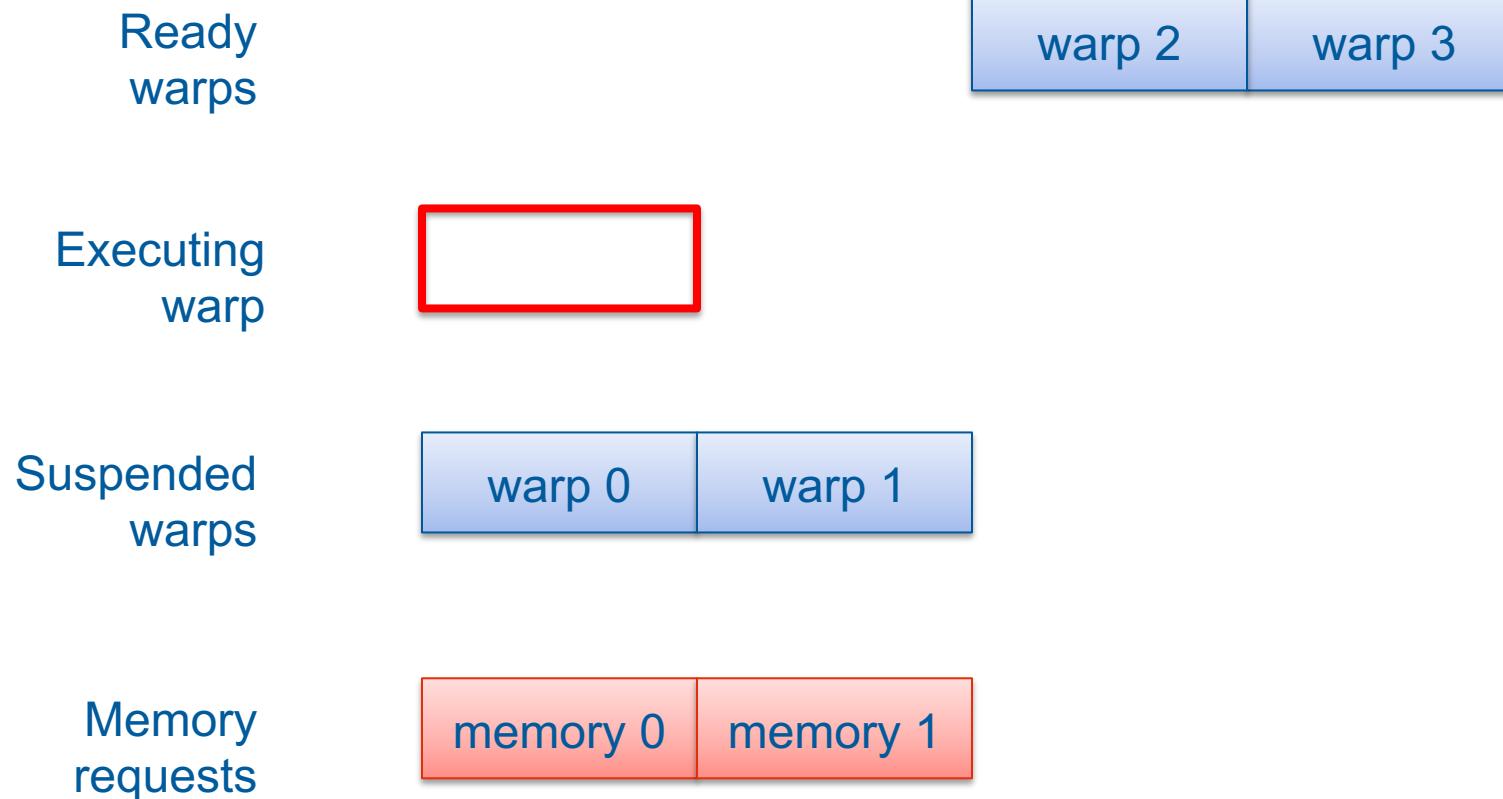
→ handle more warps than cores to hide latency

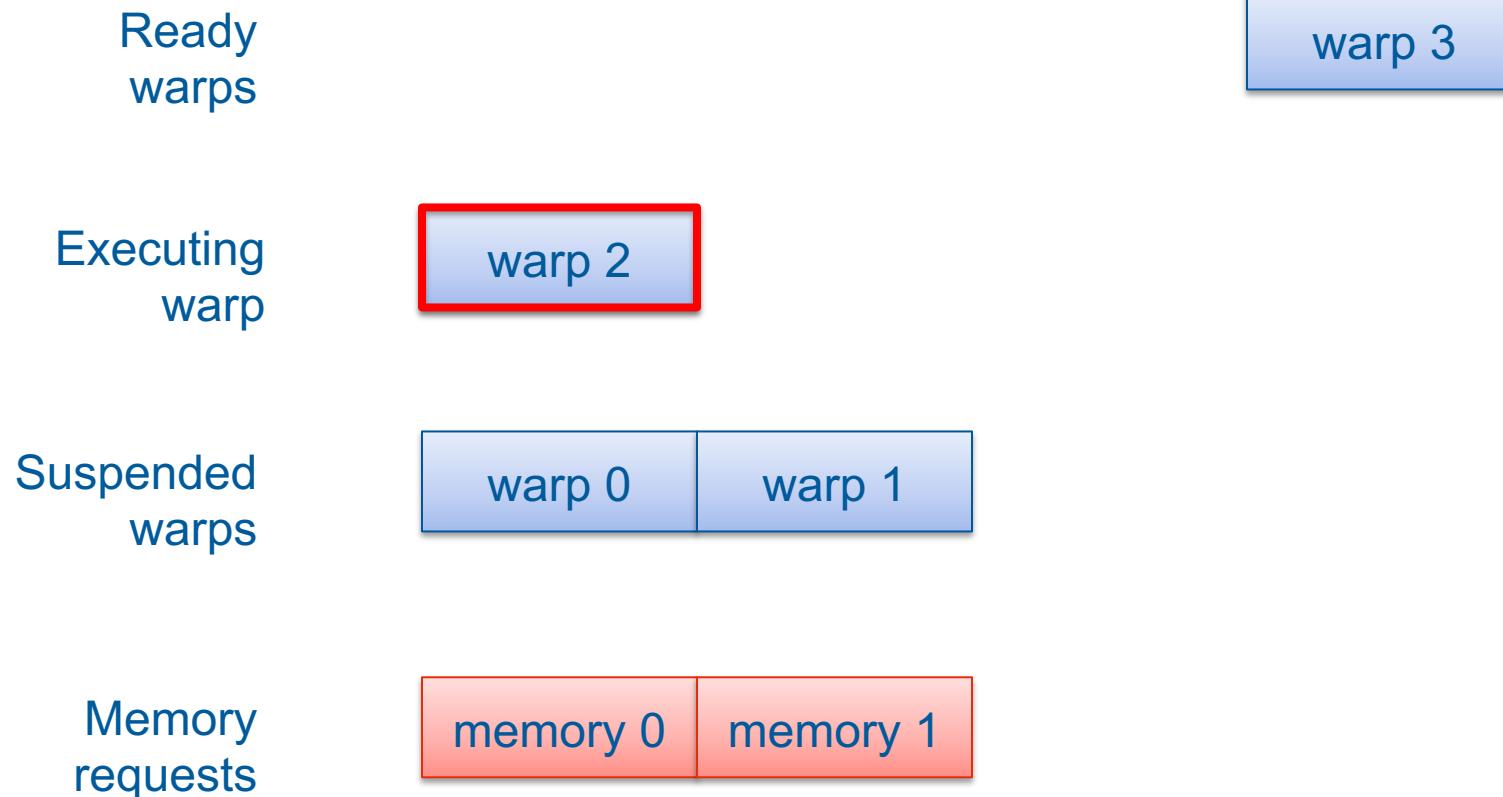


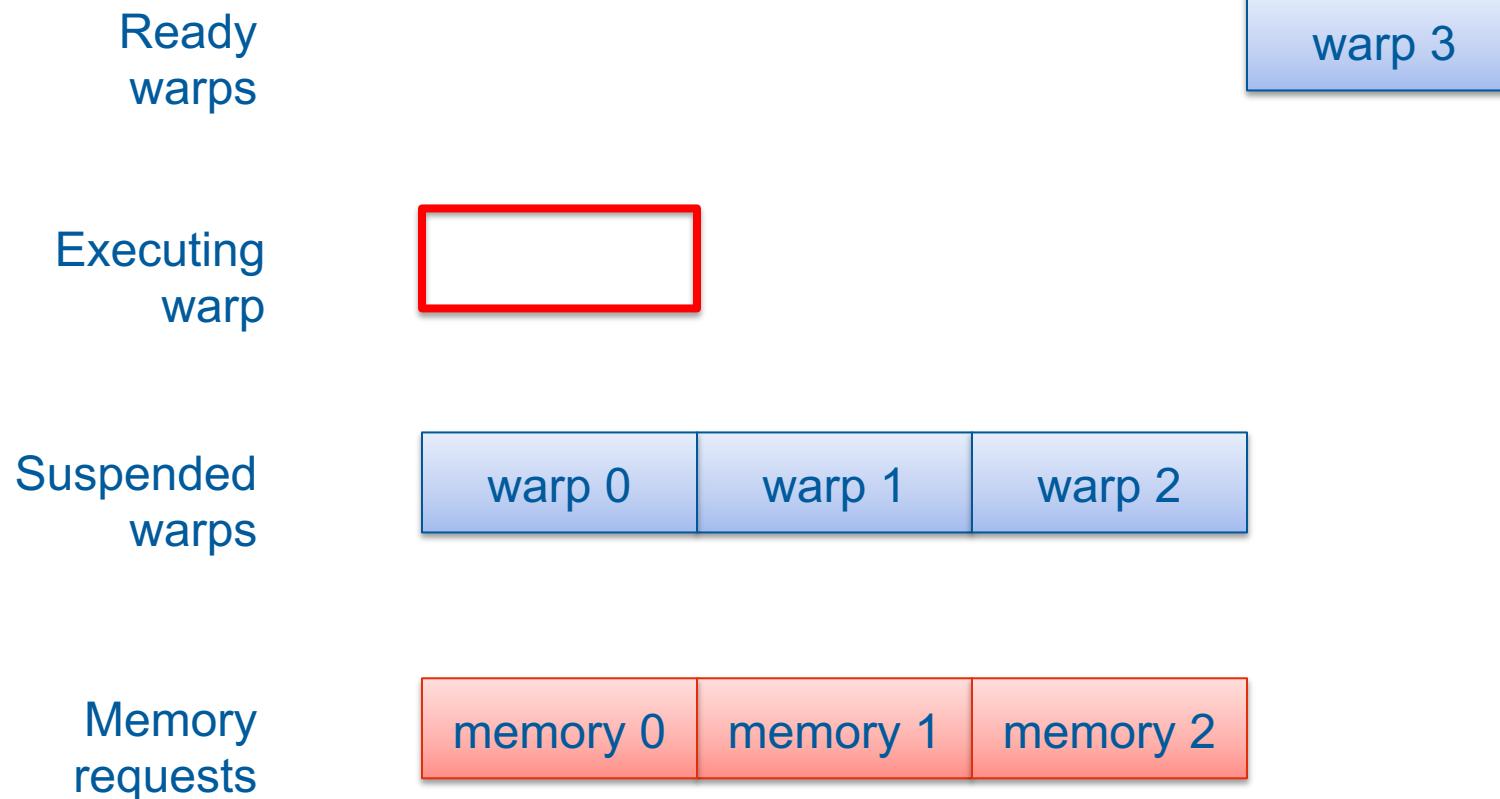


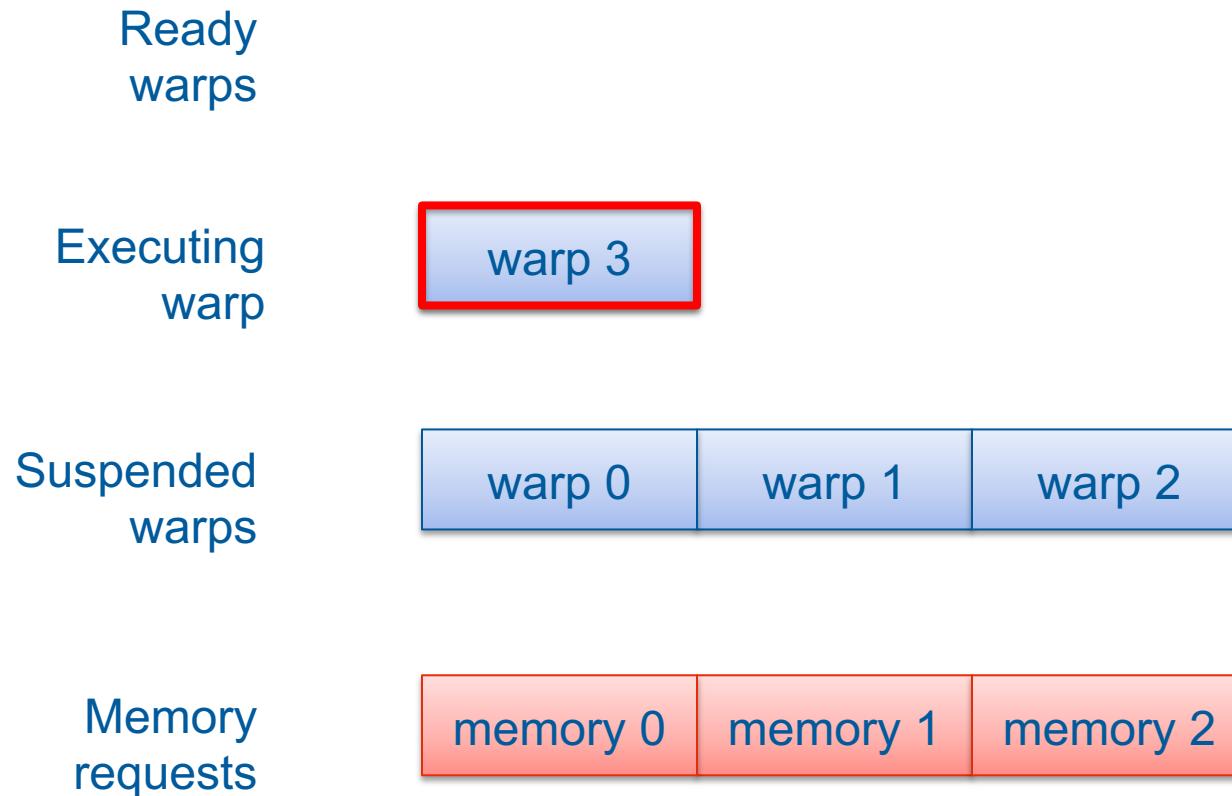


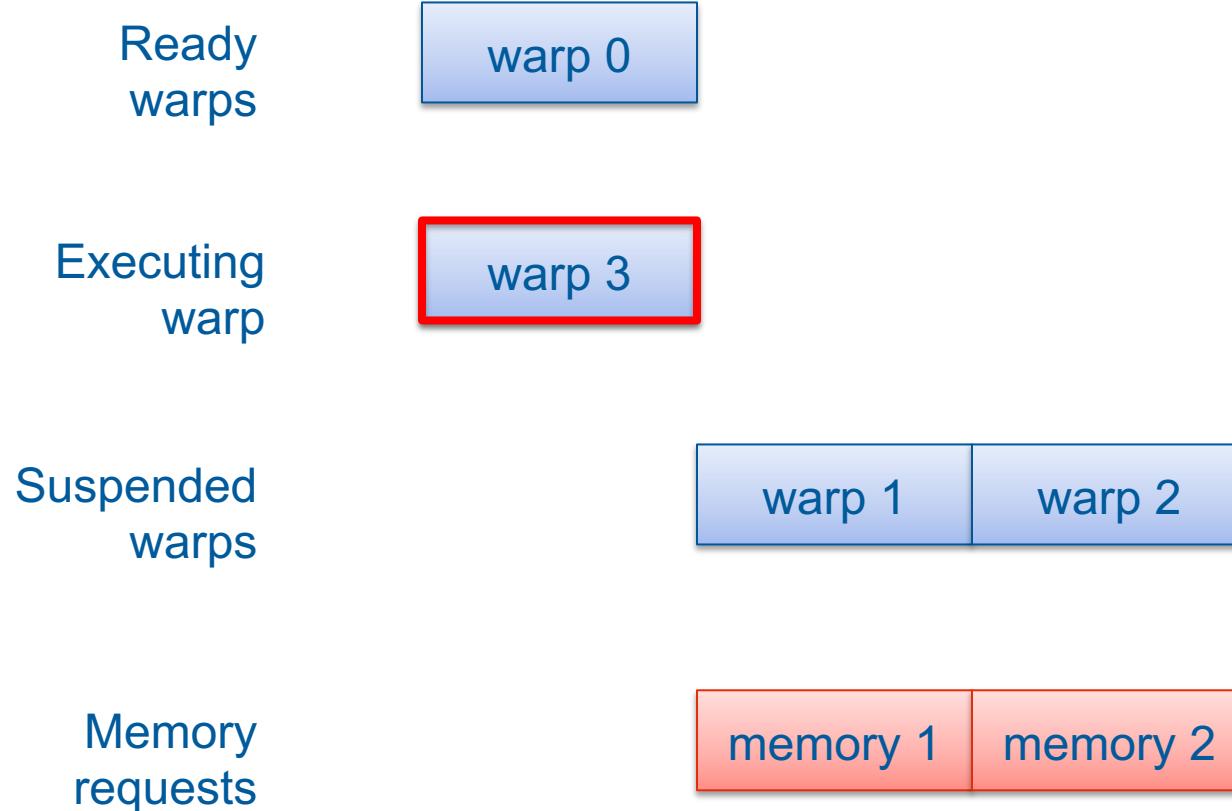


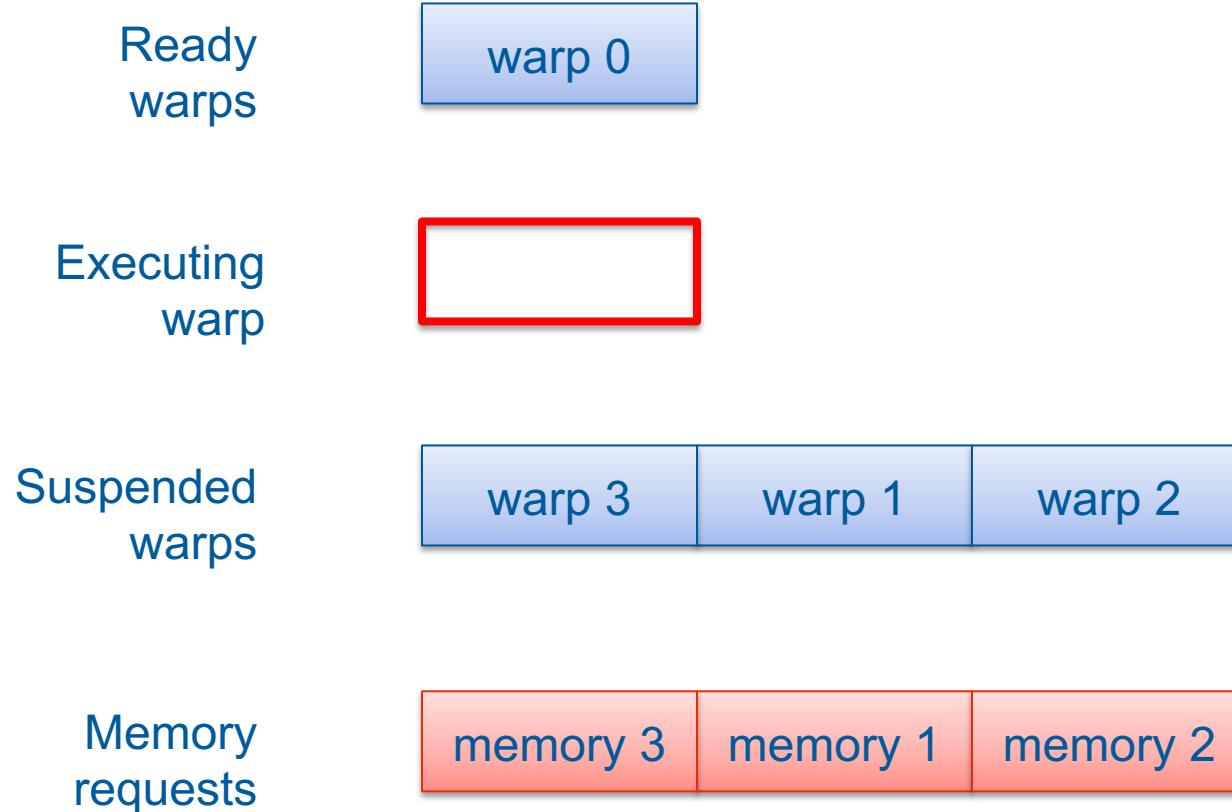


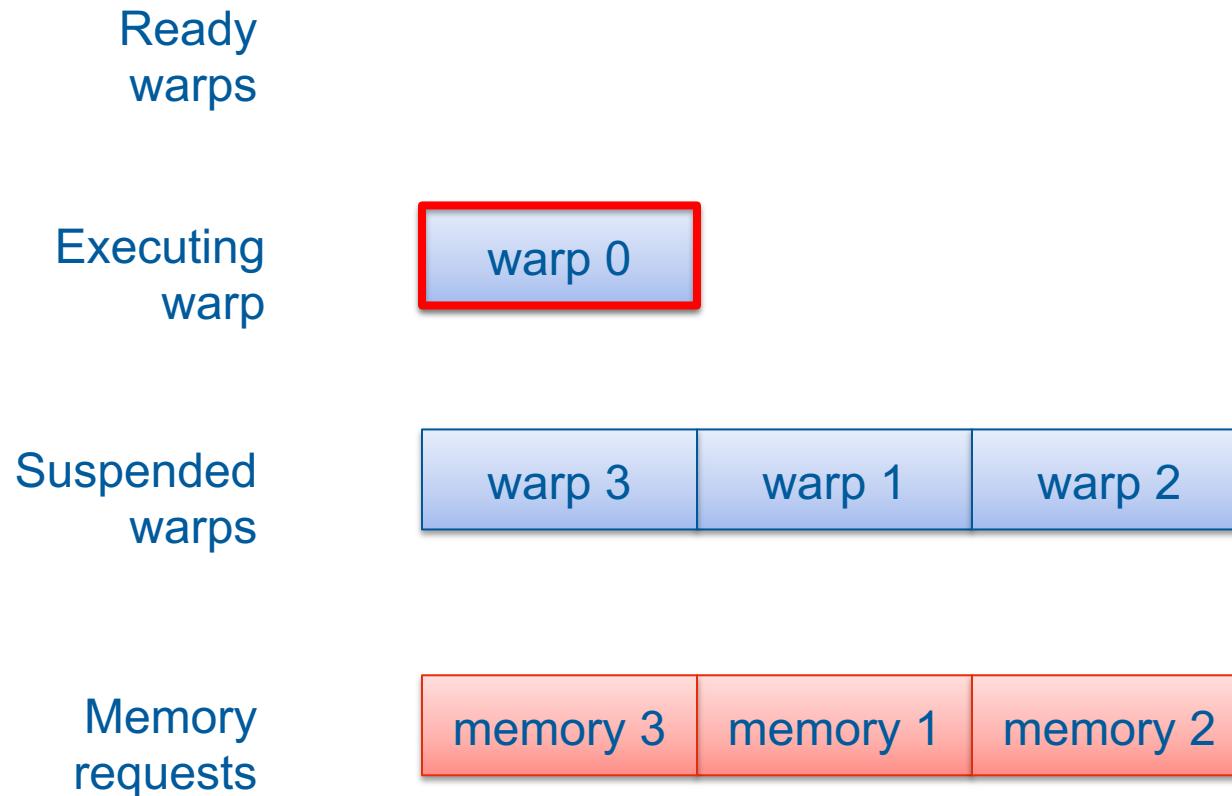














The CUDA Programming Model



- Current GPUs have 10496 cores (RTX 3090 TITAN V)
- Need more threads than cores (warp scheduler)
- Writing code for 40000 threads / 1200 warps?
 - Single-program, multiple-data (SPMD) model
- Write one program that is executed by all threads

- CUDA C++ is C++ with additional keywords to control parallel execution
- Single-source language for host and device code

- Type qualifiers
- Builtins
- Intrinsics
- Runtime API
- GPU kernel

```
_device_ float x;
_global_ void func(int* mem) {
    _shared_ int y[32];
    // ...
    y[threadIdx.x] = blockIdx.x;
    __syncthreads();
}

// ...
cudaMalloc(&d_mem, bytes);
func<<<10, 10>>>(d_mem);
```

launches

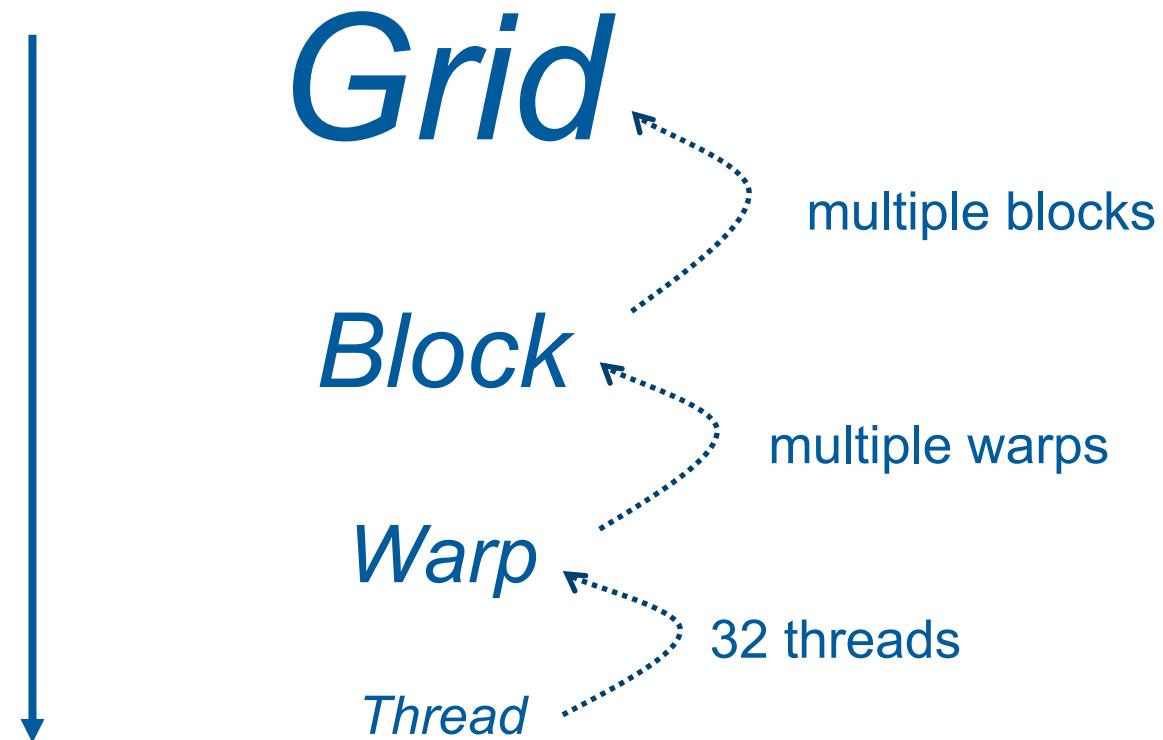


- Write **host** (CPU) and **device** (GPU) code into the same file
- nvcc compiler driver takes care of
 - Separating host and device code
 - Generating additional host code
 - For example, for kernel launches
 - Sending each part through the respective toolchain

- A function that is executed on the GPU
- Each thread is executing the same function

```
__global__ void myfunction(float* input, float* output) {  
    *output = *input;  
}
```

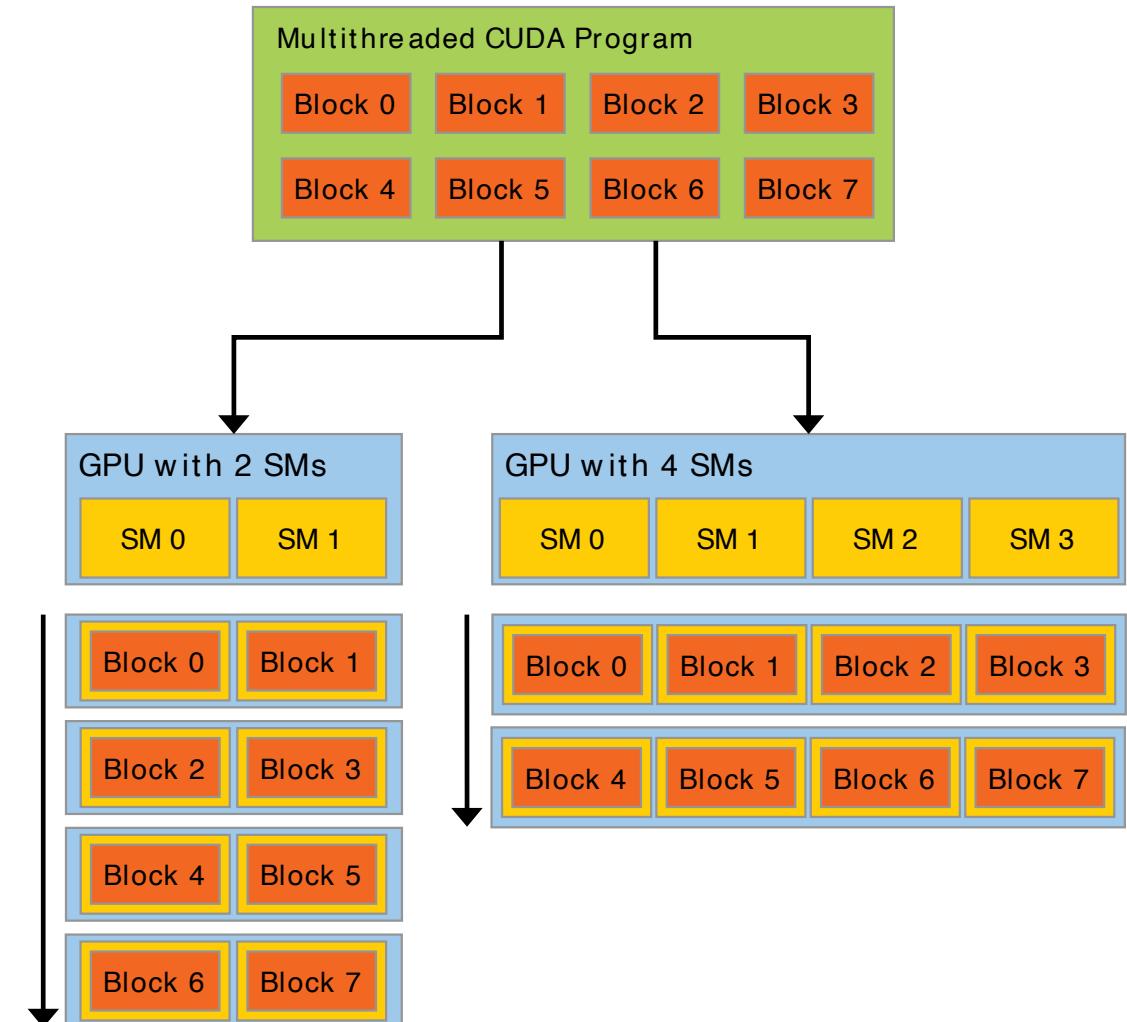
- Indicated by **__global__**
- Must have return value **void**





- Threads within one block ...
 - Are executed together on the same SM
 - Can synchronize
 - Can communicate efficiently
 - Share the same local cache
- Can be used to cooperatively compute some result
- Threads of different blocks ...
 - May be executed one after another on different SMs
 - Cannot synchronize
 - Can only communicate inefficiently
- Should work independently of other blocks

- Block queue feeds multiprocessors
- Number of available multiprocessors determines number of concurrently executed blocks



- Kernel is split up in blocks of threads





- Requires threads per block, number of blocks
- Launch is asynchronous

```
dim3 block_size(128, 1, 1);
dim3 grid_size(12, 1, 1);
myfunction<<<grid_size, block_size>>>(input, output);
```

- Block and thread setup specified within <<< , >>>
- Input parameters are copied to the GPU

■ blockIdx and threadIdx

0,0	1,0	2,0	3,0,0,0 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,0 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,1 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,1 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,2 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,2 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,3 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,3 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,4 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,4 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,5 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,5 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,6 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,6 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,7 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,7 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,8 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,8 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,9 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,9 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,10 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,10 ^{4,0}	5,0	6,0	7,0
0,0	1,0	2,0	3,0,0,11 ^{4,0}	5,0	6,0	7,0	0,0	1,0	2,0	3,0,1,11 ^{4,0}	5,0	6,0	7,0



- Using **threadIdx** and **blockIdx** execution paths are chosen
- With **blockDim** and **gridDim** number of threads can be determined

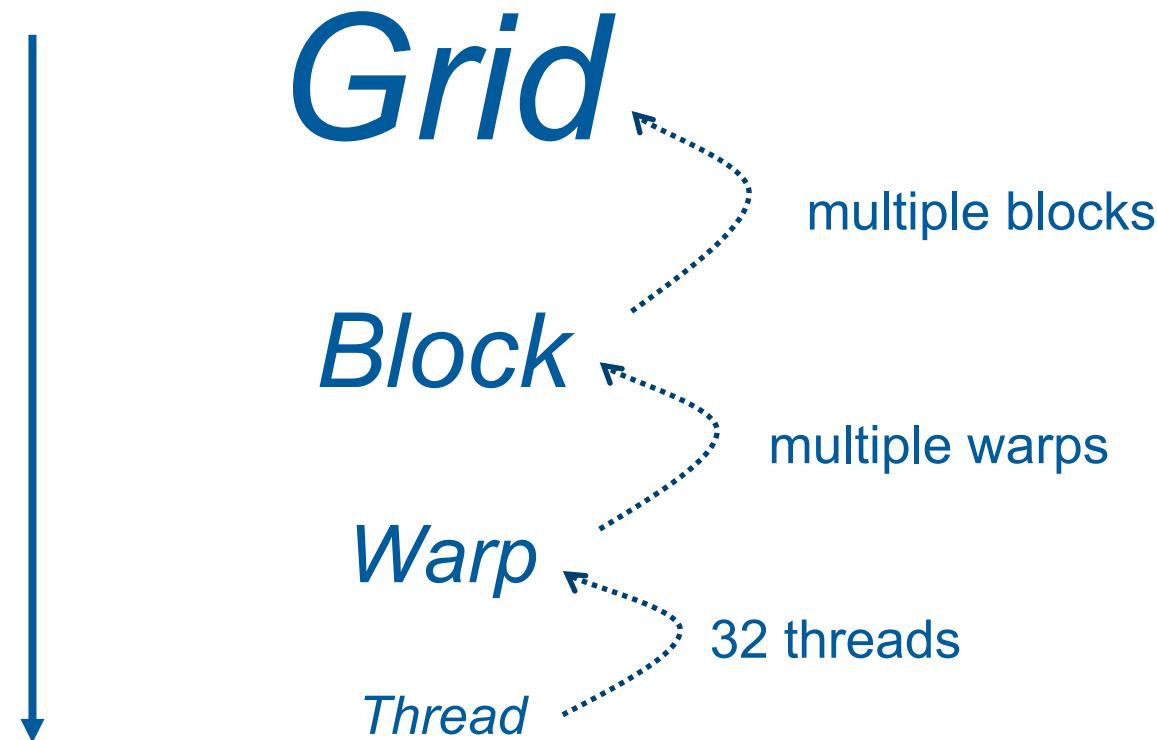
```
__global__ void myfunction(float* input, float* output) {
    uint id = threadIdx.x + blockIdx.x * blockDim.x;
    output[id] = input[id];
}
```

- On each multiprocessor each block is split up into warps
- Threads with the lowest id map to the first warp

0,0	1,0	2,0	3,0	4,0	5,0	6,0	7,0	8,0	9,0	10,0	11,0	12,0	13,0	14,0	15,0
0,1	1,1	2,1	3,1	4,1	5,1	6,1	7,1	8,1	9,1	10,1	11,1	12,1	13,1	14,1	15,1
0,2	1,2	2,2	3,2	4,2	5,2	6,2	7,2	8,2	9,2	10,2	11,2	12,2	13,2	14,2	15,2
0,3	1,3	2,3	3,3	4,3	5,3	6,3	7,3	8,3	9,3	10,3	11,3	12,3	13,3	14,3	15,3
0,4	1,4	2,4	3,4	4,4	5,4	6,4	7,4	8,4	9,4	10,4	11,4	12,4	13,4	14,4	15,4
0,5	1,5	2,5	3,5	4,5	5,5	6,5	7,5	8,5	9,5	10,5	11,5	12,5	13,5	14,5	15,5
0,6	1,6	2,6	3,6	4,6	5,6	6,6	7,6	8,6	9,6	10,6	11,6	12,6	13,6	14,6	15,6
0,7	1,7	2,7	3,7	4,7	5,7	6,7	7,7	8,7	9,7	10,7	11,7	12,7	13,7	14,7	15,7



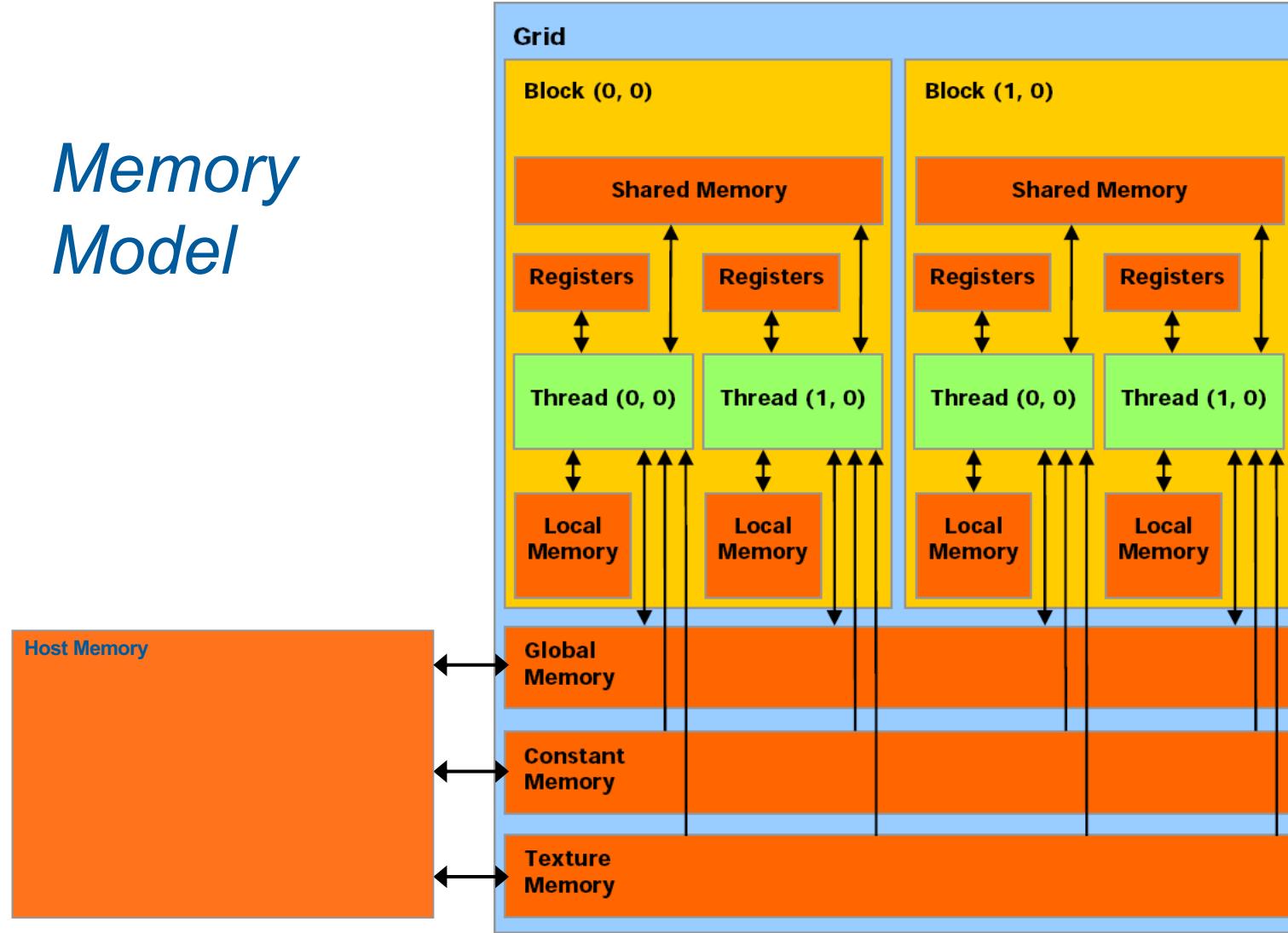
- What block size to choose?
- $$occupancy = \frac{active\ warps}{maximum\ active\ warps}$$
- Potential limits
 - Register usage
 - Shared memory usage
 - Block size
- Start with 128 / 256 threads
- Maximum size 1024 threads





Memory Spaces

Memory Model



- Private for each thread
- Fastest memory
- Automatically allocated from per SM register file (number of registers per block)
- No keyword

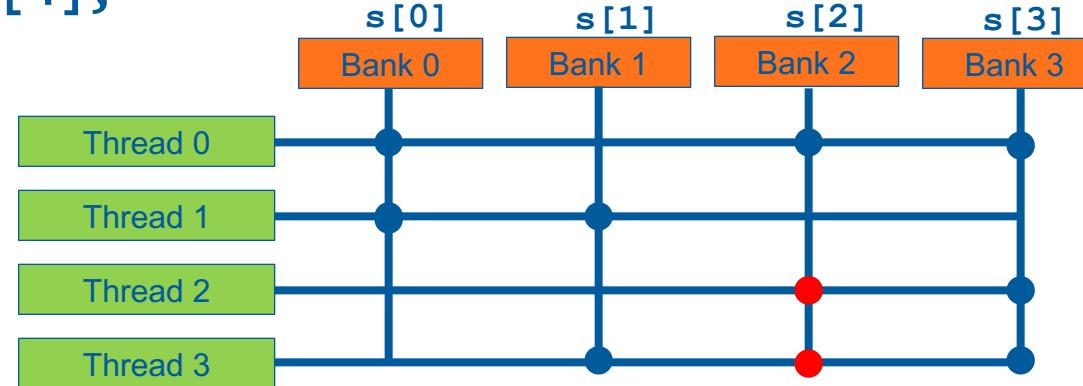
- Main GPU memory
- Accessible from all blocks/threads
- Access pattern matters
 - Older devices coalescing!
 - Newer devices: close together
- Slowest on device memory
 - (up to 1000x times slower than registers)
- Cached on newer devices
- **`__device__` / `cudaMalloc`**



- Accessible from all blocks/threads
- Cached on all devices
 - Special layout for 2D and 3D textures
- Hardware interpolation (same as OpenGL)
- Automatic handling of boundary cases
- Normalized texture coordinates
- Conversion to float
- No concurrent read and write
- **cudaMallocArray + cudaMemcpyToArray**

- Read-only
- Ideal for coefficients and other data that is read uniformly by warps
- Cached on all devices
- `__constant__ + cudaMemcpyToSymbol`

- Shared access within one block (lifetime: block)
- Located on multiprocessor → very fast
- Limited available size on multiprocessor
- Crossbar: simultaneous access to distinct banks
- `__shared__ int s[4];`

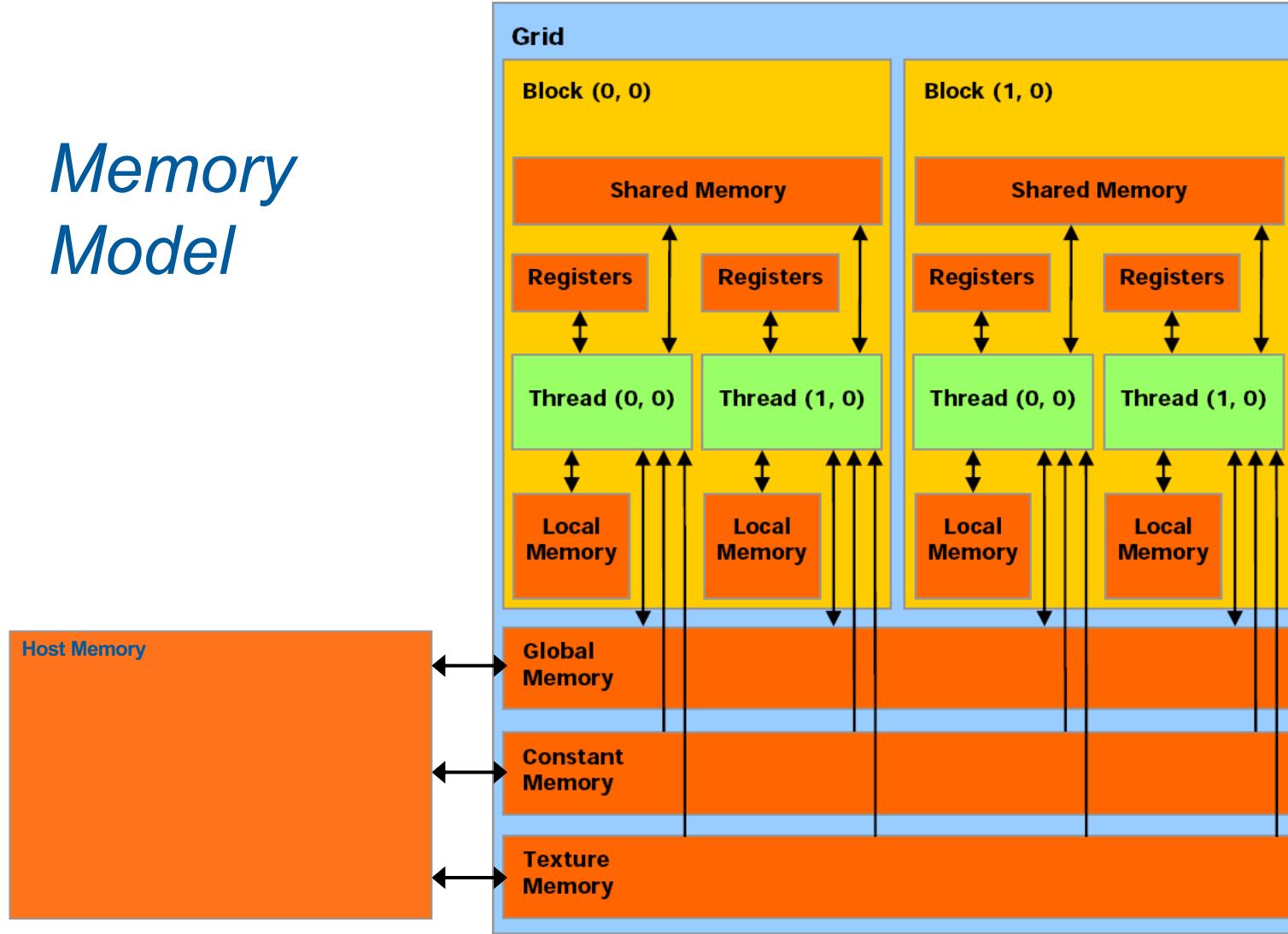


- If above register limited, registered are spilled to global
- Very slow in comparison to registers
- Besides that, same behavior as registers
- Always perfect access pattern



- Can be mapped to device memory space
- Same data can be accessed from host and device concurrently
- Even slower than global memory
- Mapped memory:
cudaHostAlloc with **cudaHostAllocMapped**

Memory Model



	Global	Constant	Texture	Shared	Local	Registers
Host Code	Dynamic allocation R/W	No allocation R/W	Dynamic allocation R/W	Dynamic allocation No access	No allocation No access	No allocation No access
Device Code	Static allocation (Dynamic allocation) R/W	Static allocation Read-only	Static allocation Read-only	Static allocation R/W	Static allocation R/W	Static allocation R/W