The left side of the slide features a decorative vertical bar with a grid pattern, a solid orange vertical bar, and a thin orange vertical line. To the right of these bars are several orange circles of varying sizes, including a large one at the top and a smaller one containing the number '1'.

PARALLEL PROCESSING

UNIT 2

Dr. Ahmed Sallam

1



OUTLINES

- Threads communication patterns
- GPU hardware: Processing model
- GPU hardware: Memory model
- Efficient GPU programming

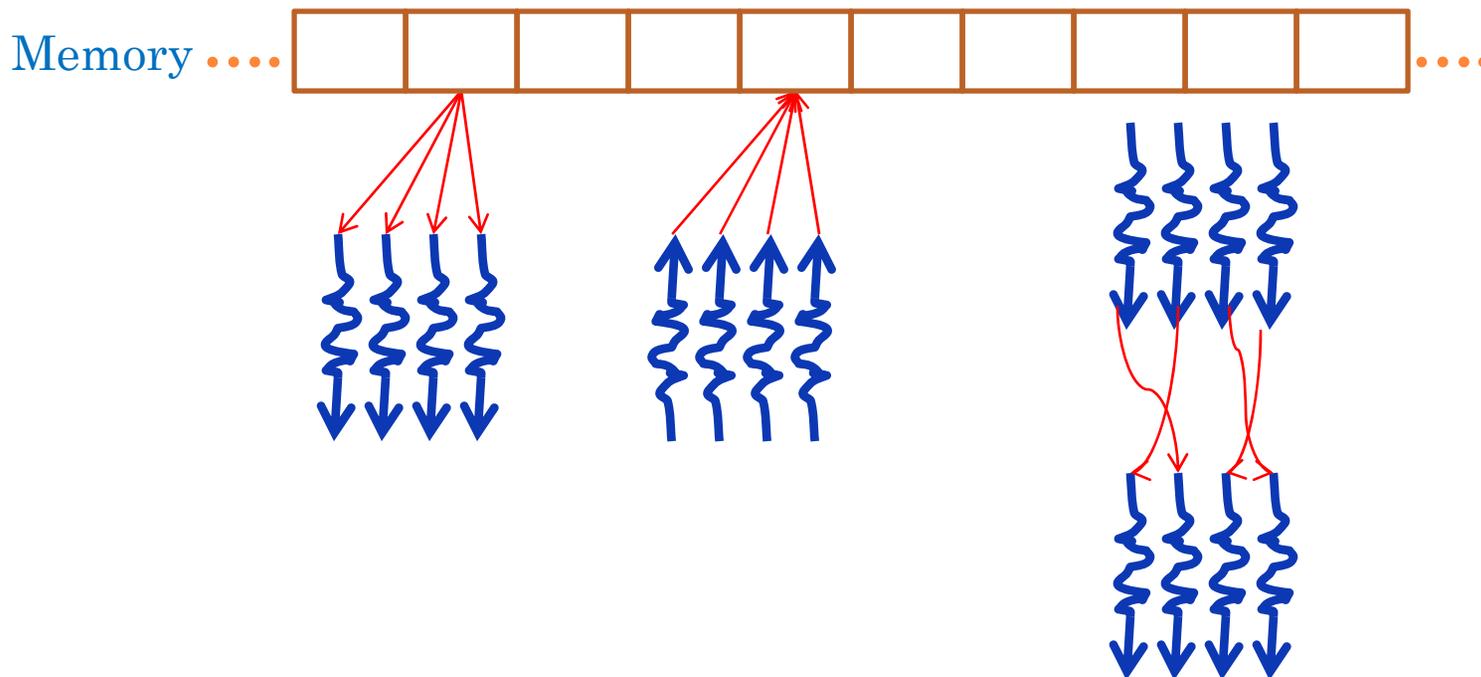
The background features a black field with vertical stripes on the left and right sides. The stripes consist of a thin black line, a wider grey textured band, and a thin orange line. Several orange circles of varying sizes are scattered across the page, with one large circle on the left and another on the right. The text is centered in the middle-right area.

THREADS COMMUNICATION PATTERNS

3

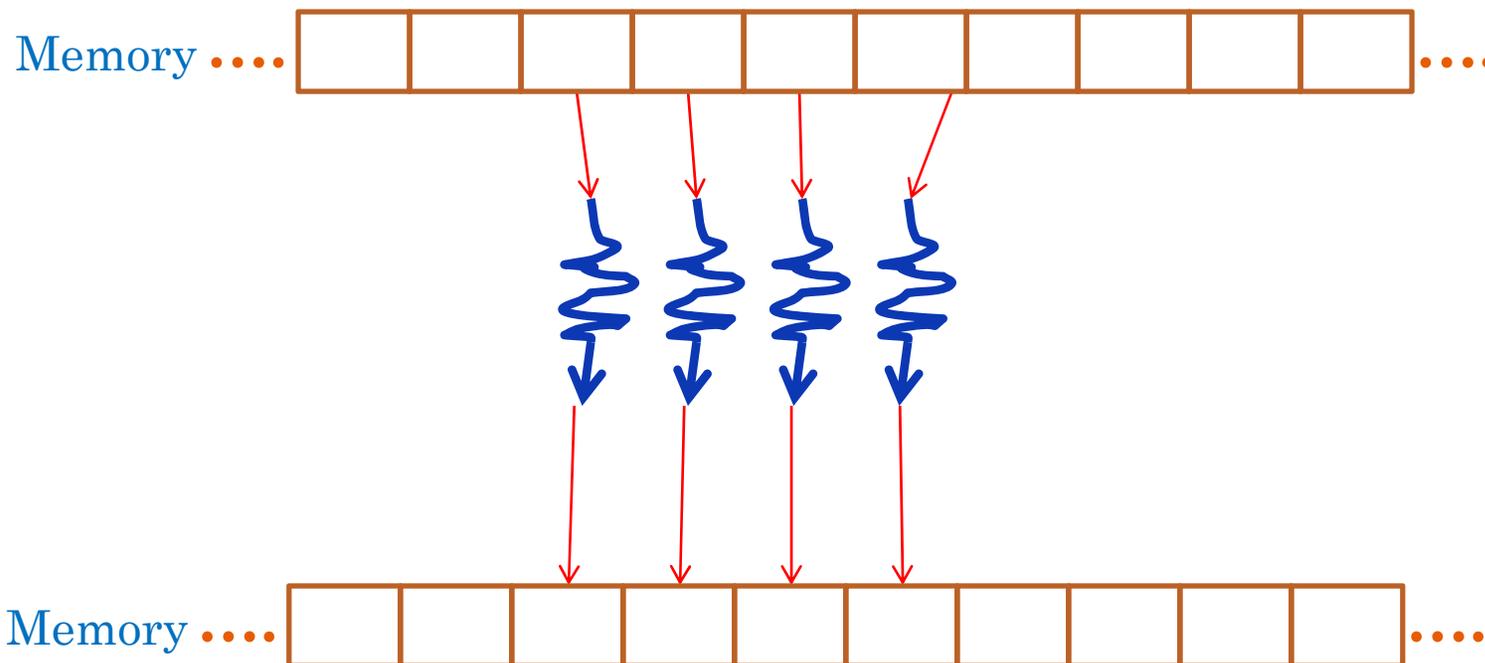
COMMUNICATION PATTERNS

- Parallel Programming: threads solving problem by working together (Communicating through memory)



MAP

- Tasks read from and write to specific data elements



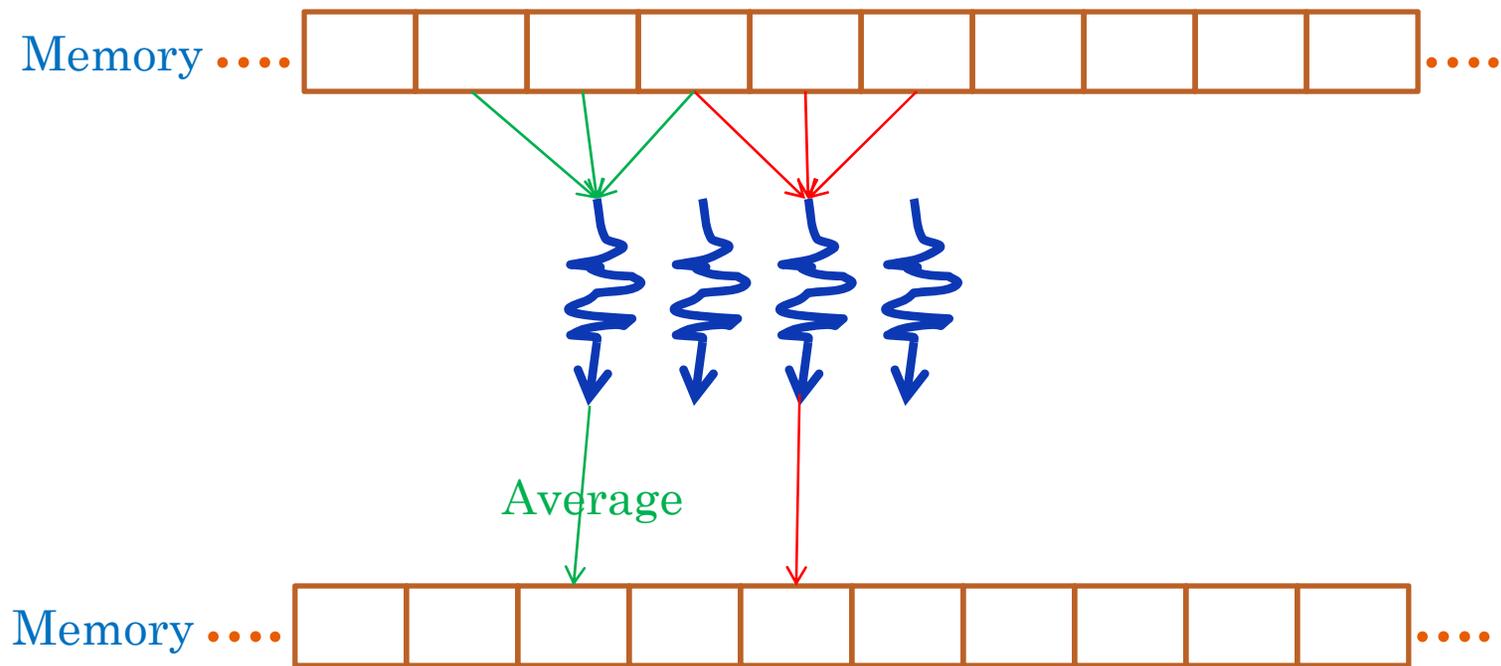
MAP APPLICATION

- Color image to Gray image



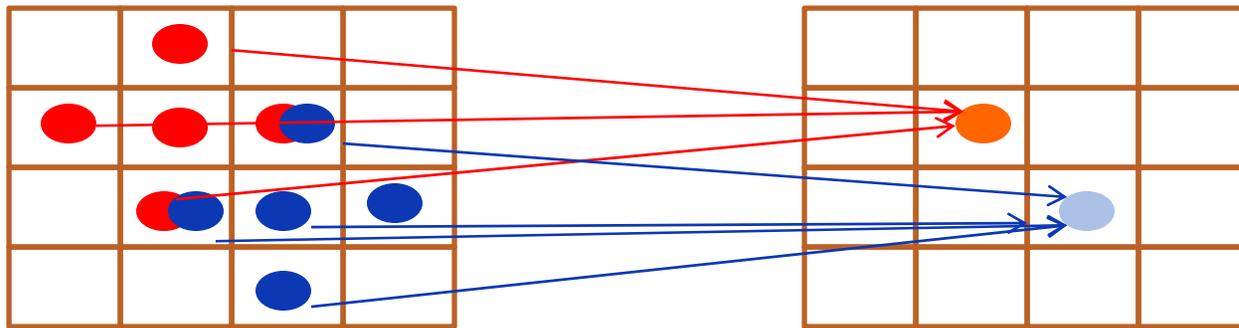
GATHER

- Problem: Many threads read from the same place (Data reuse)



GATHER APPLICATION

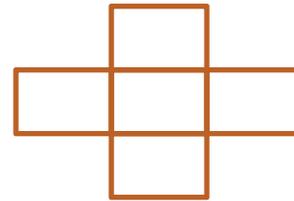
- Image Blurring



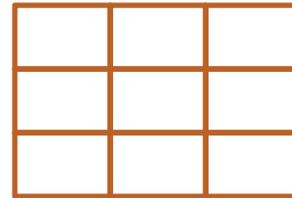
GATHER (STENCIL)

- Tasks gather input from a fixed neighbors.

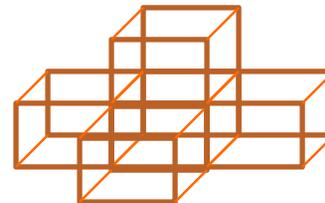
2D Neumann



2D Moore

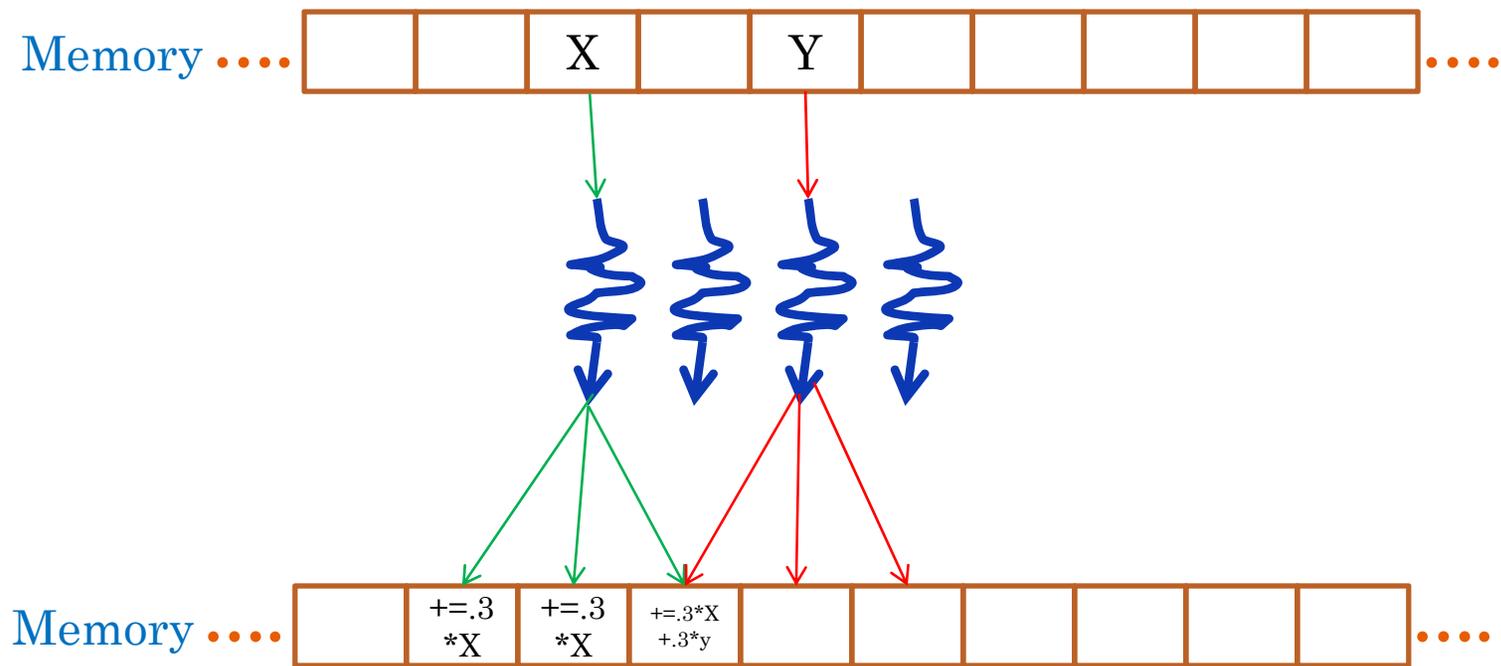


3D Neumann



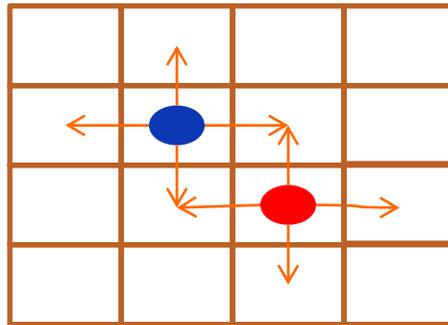
SCATTER

- Tasks compute where to write output
- Problem: threads write to the same place



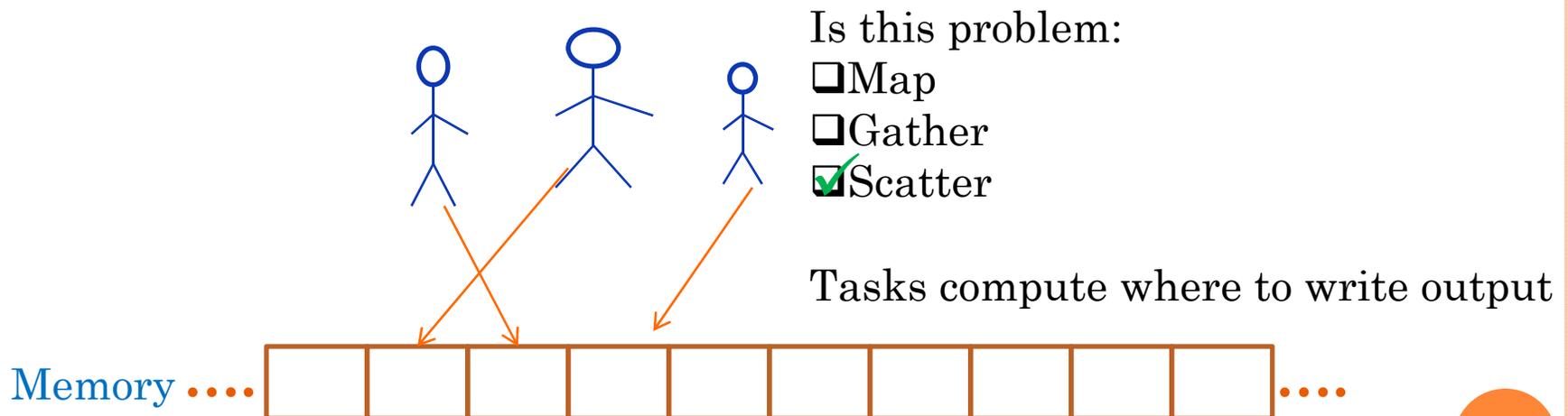
SCATTER APPLICATION

- Image Blurring

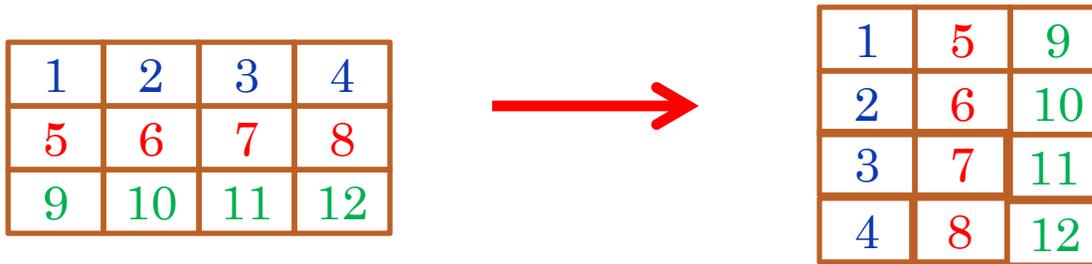


QUIZ

- Given a list of players with the information (name, height, tall-rank)
- Write each player's record in its location in a sorted list according to the tall-rank



TRANSPOSE



Row major order



Column major order



TRANSPOSE APPLICATIONS

- Array
- Matrix
- Image
- Data Structures (AOS to SOA)

```
Struct foo{  
    float f;  
    int i;  
};
```

AOS



SOA



QUIZ

○ Label code snippets by problem:

```
float out[], in [];  
int i= threadIdx.x;  
int j= threadIdx.y;  
  
const float pi= 3.1415;
```

1. Map
2. Gather
3. Scatter
4. Stencil
5. Transpose

1 out[i]= pi* in[i];

5 out[i + j*128]=in[j + i*128];

3 if(i%2){
out[i-1] +=pi * in[i]; out[i+1] += pi * in[i];

2 out[i] = (in[i] + in[i-1]+ in[i+1] * pi/3.0f)
}

NOTES

- Map one-to-one
- Transpose one-to-one
- Gather many-to-one
- Scatter one-to-many
- Stencil several-to-one

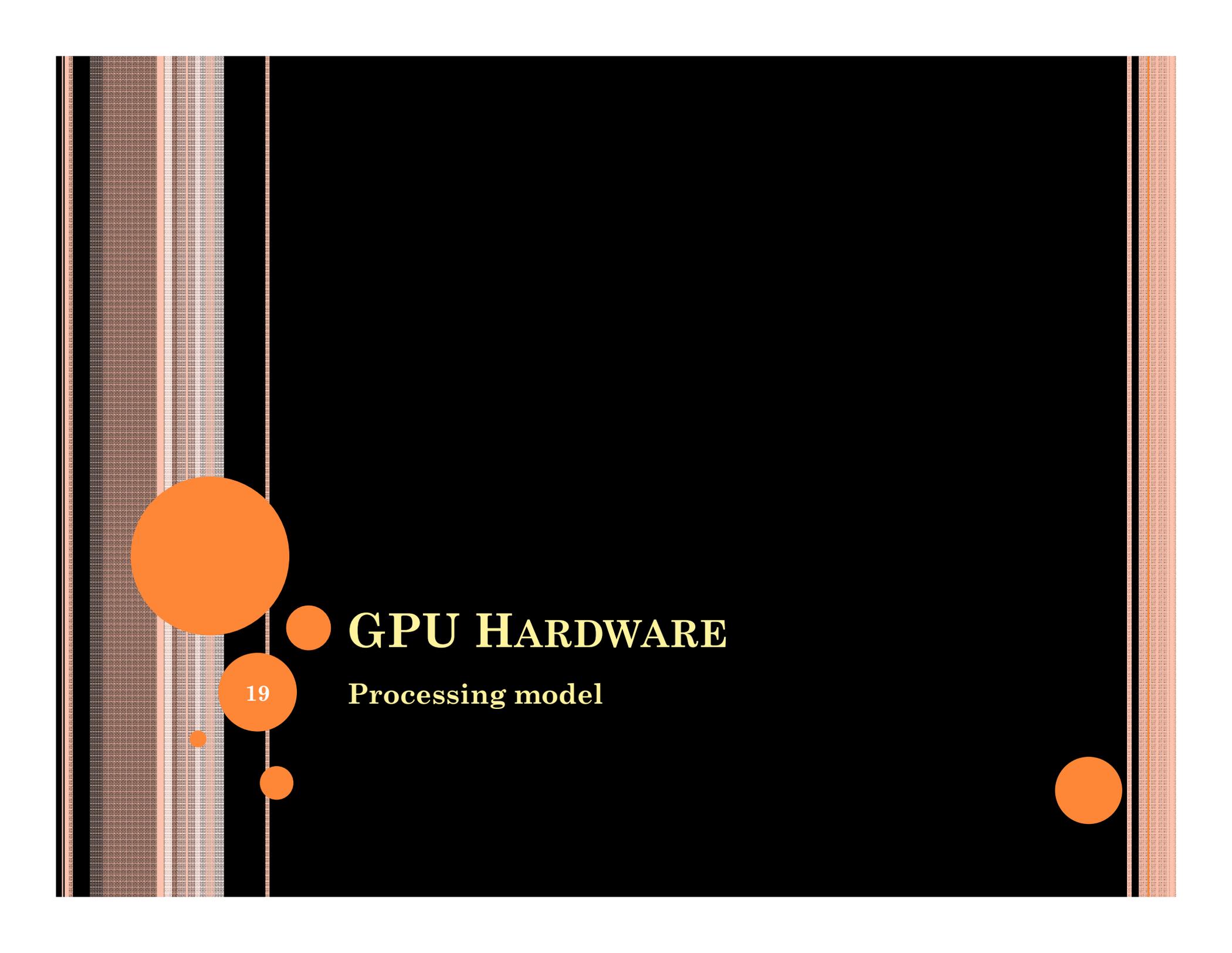
- Reduce all-to-one
- Scan/sort all-to-all

NOTES (CONT.)

- How can threads efficiently access memory in concert?
 - How to exploit data reuse?
- How can threads communicate partial results by sharing memory?
 - safely

OUTLINES

- Threads communication patterns
- GPU hardware: Processing model
- GPU hardware: Memory model
- Efficient GPU programming



GPU HARDWARE

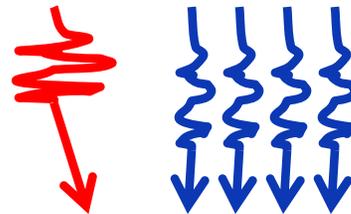
Processing model

19

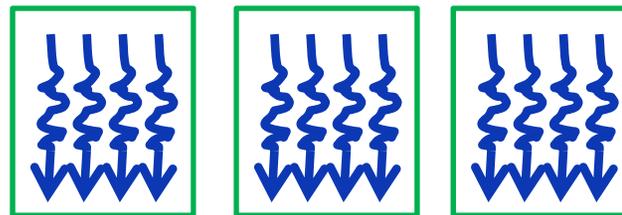
PROGRAMMER VIEW OF THE GPU

- Kernels C/C++ function

- Threads

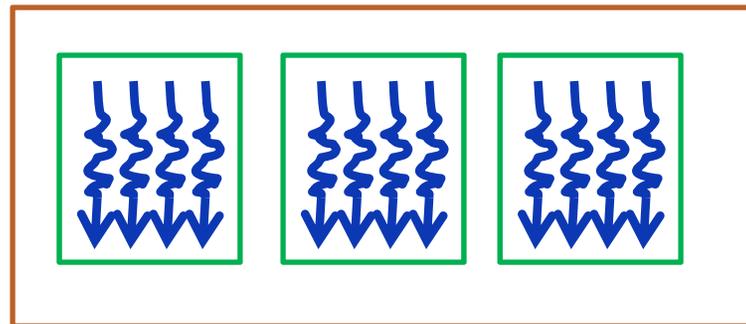


- Blocks: Group of threads that cooperate to solve subproblem.

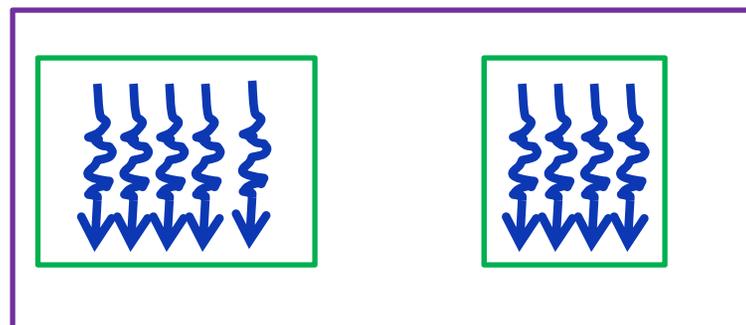


PROGRAMMER VIEW OF THE GPU (CONT.)

- Kernel

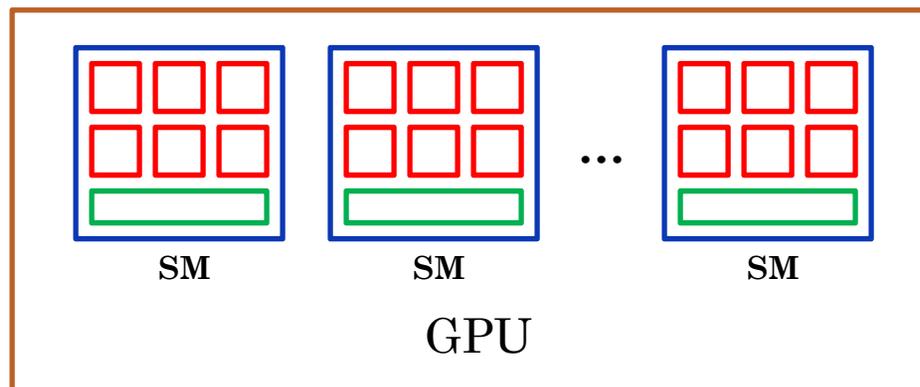


- Another Kernel



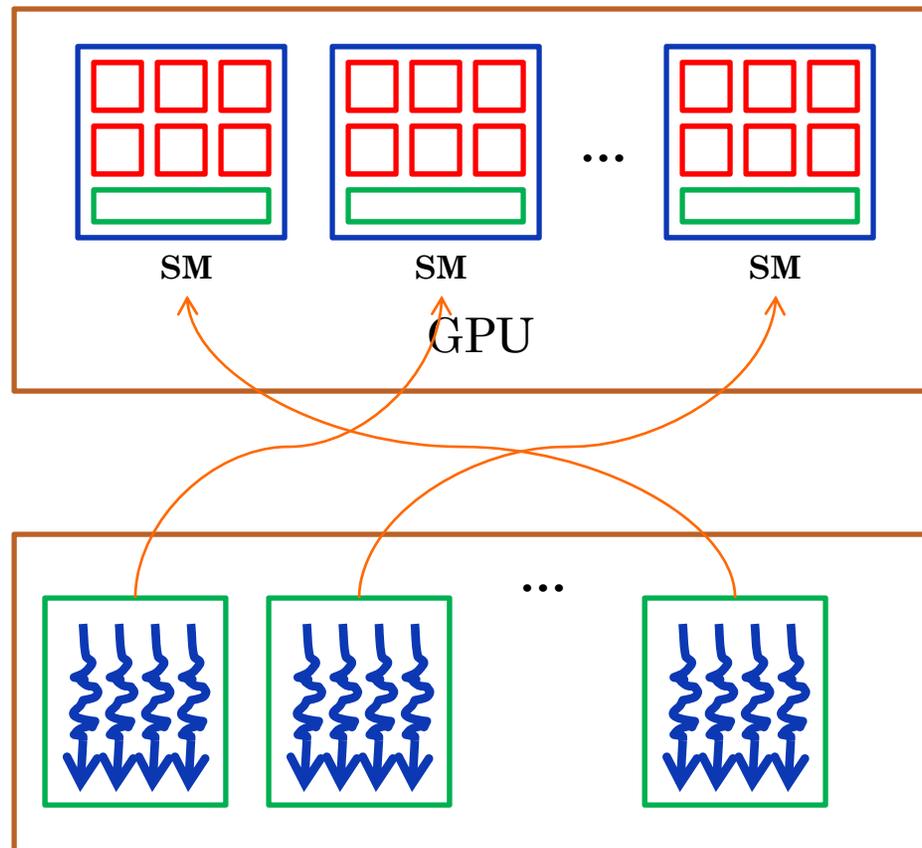
GPU HARDWARE

- GPU consists of Streaming Multiprocessors (SMs)
- SM consists of many cores each has local memory (Registers) and shared memory (L1,L2 cache)



BLOCKS GPU HARDWARE

- GPU is responsible for allocating blocks to SMs



QUIZ

- A thread block contains many threads
- An SM may run more than one block
- A block may run on more than SM
- All the threads in a thread block may cooperate to solve a subproblem
- All the threads that run on a given SM may cooperate to solve a subproblem

- The programmer is responsible for defining thread blocks in software.
- The programmer is responsible for allocating thread blocks to hardware streaming multiprocessors (SMs)

QUIZ

Given a single kernel that is launched on many thread blocks including X and Y, the programmer can specify:

- ❑ That block X will run at the same time as block Y
- ❑ That block X will run after block Y
- ❑ That block X will run on SM Z

CUDA GUARANTEES

CUDA makes few guarantees:

- All threads in the same block can run on the same SM at the same time
- All blocks in the kernel finish before any blocks in another kernel

Advantages:

- Hardware can run things efficiently (Flexibility)
- No waiting on slow blocks if we have free SMs
- Scalability
 - From cell phones to super computers
 - From current to future GPUs

Consequences :

- No assumptions about what SM run what Block.
- No communication between blocks
 - If Block Y needs result from Block X, we could have a dead lock state because Block X will evaporate after completion.

QUIZ

- How many different outputs can different run of this program produce?

```
#include <stdio.h>

#define NUM_BLOCKS 16
#define BLOCK_WIDTH 1
__global__ void hello()
{
    printf("Hello world! I'm a thread in block %d\n", blockIdx.x);
}

int main(int argc, char **argv)
{
    // launch the kernel
    hello<<<NUM_BLOCKS, BLOCK_WIDTH>>>();

    // force the printf()s to flush
    cudaDeviceSynchronize();

    printf("That's all!\n");

    return 0;
}
```

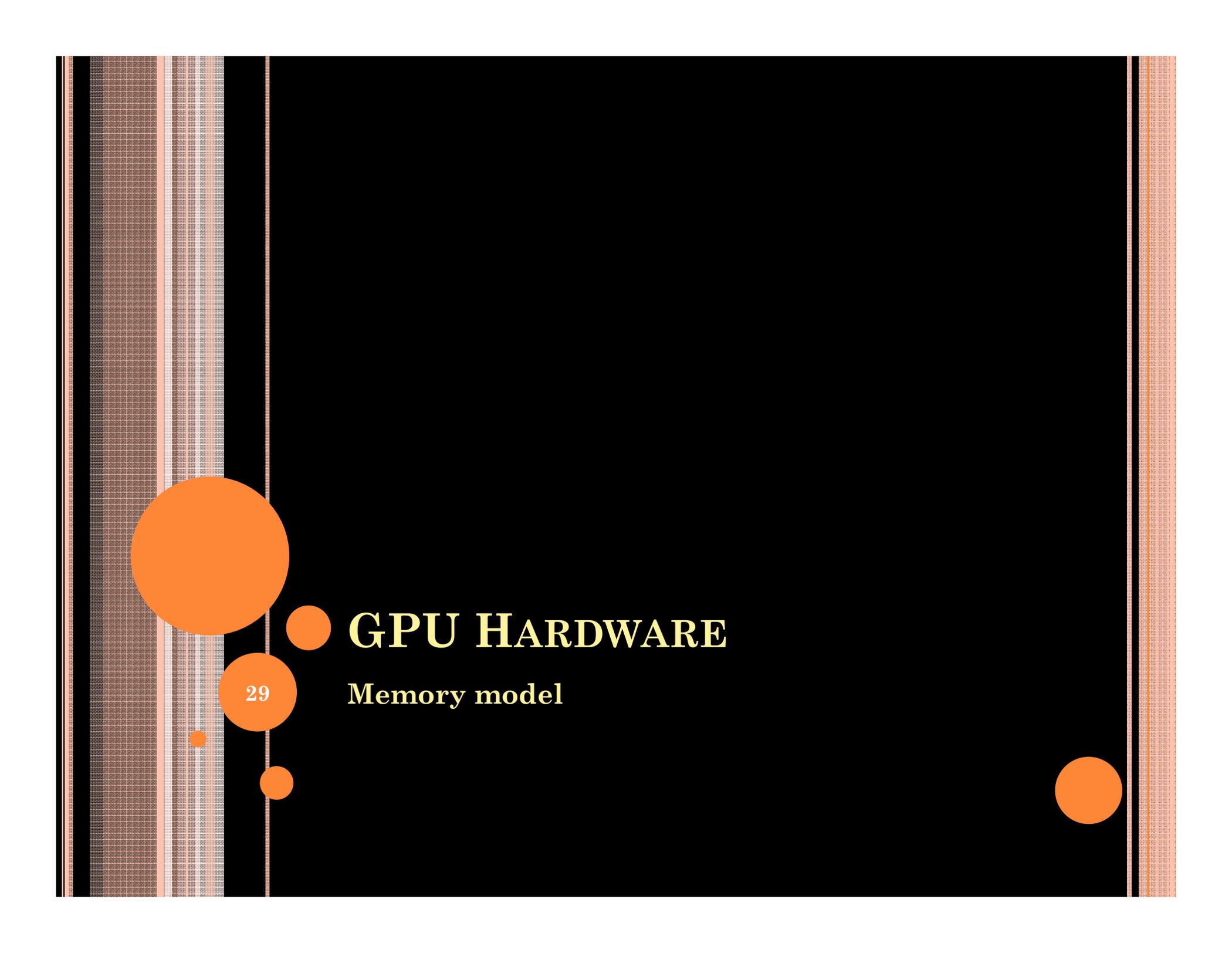
Is this problem:

- 1
- 16
- 65,536
- 21 trillion

=16!

OUTLINES

- Threads communication patterns
- GPU hardware: Processing model
- GPU hardware: Memory model
- Efficient GPU programming

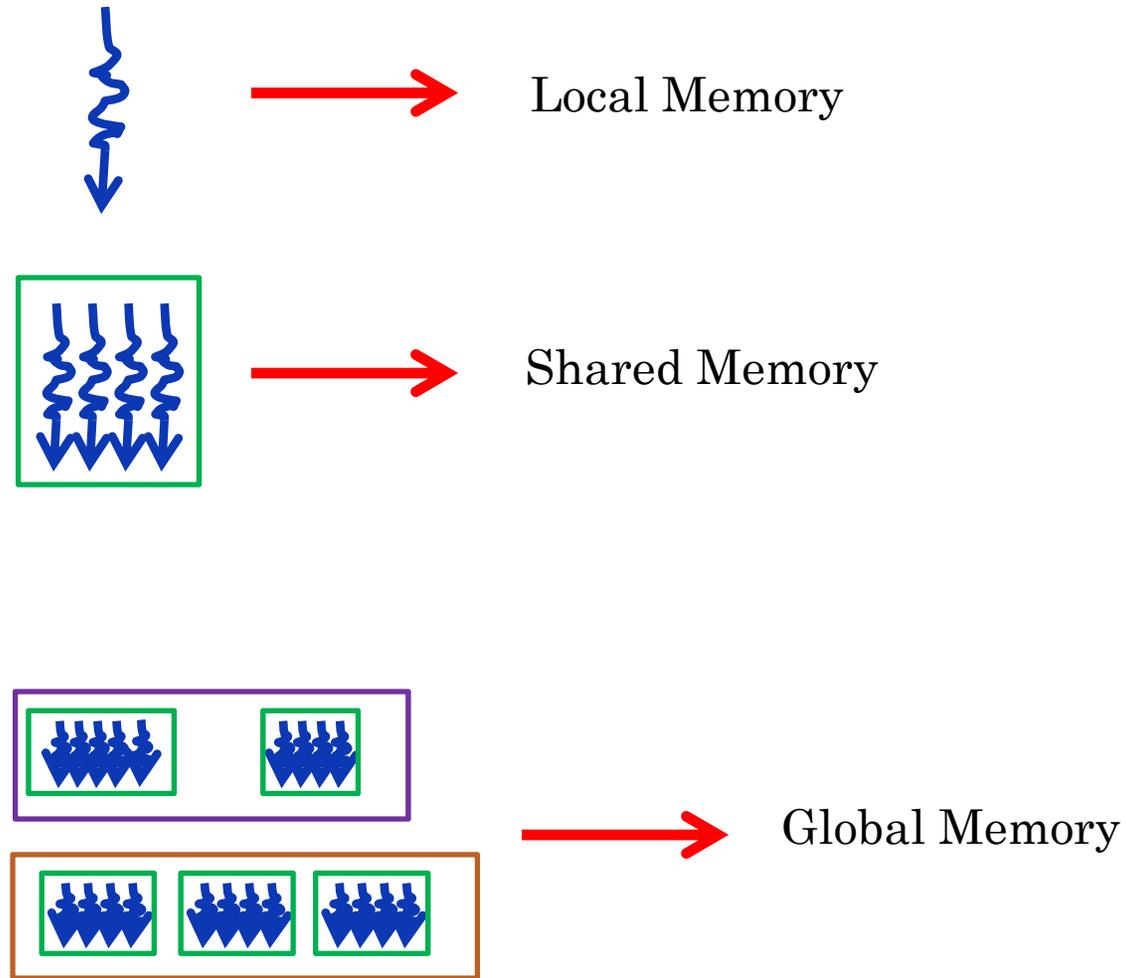


GPU HARDWARE

Memory model

29

GPU MEMORY



GPU MEMORY (CONT.)

```
__global__ void GPU_memory(float *out, float *in )  
{  
  
    int i= threadIdx.x;  
    int j= threadIdx.y;  
  
    const float pi= 3.1415;  
  
    out[i]= pi* in[i];  
  
    __shared__ float sh_arr[128];  
  
}
```

QUIZ

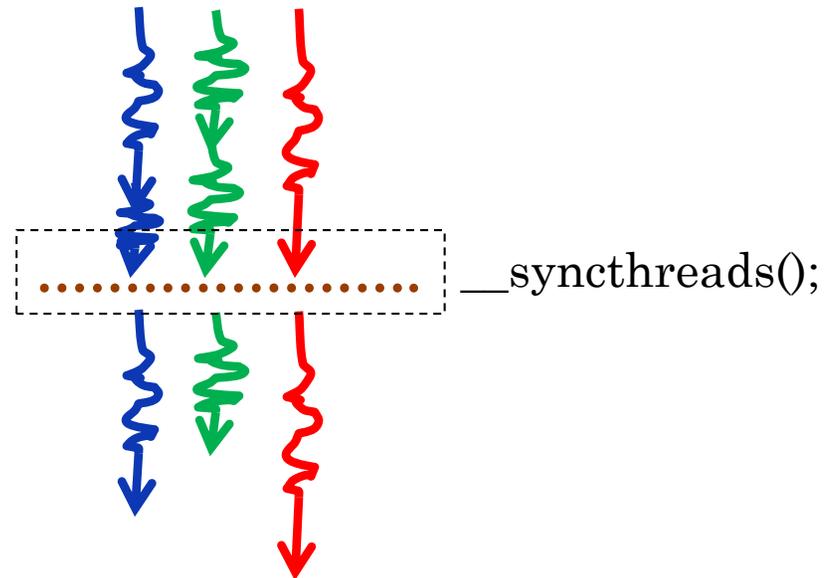
- ✓ All threads from a block can access the same variable in that blocks shared memory
- ✓ Threads from two different blocks can access the same variable in global memory
- ✓ Threads from different blocks have their own copy of local variables in local memory
- ✓ Threads from the same block have their own copy of local variables in local memory

SYNCHRONIZATION

- Threads can access each other's results through shared and global memory. (Work together)
- **Problem:** what if a thread reads a result before another thread writes it?

SYNCHRONIZATION BARRIER

- Barrier: is a point in the program where threads stop and wait when all threads have reached the barrier, then they can proceed.



QUIZ

- Are the following code snippet correct?

```
__global__ void foo( ... )  
{  
    __shared__ int s[128];  
    int i= threadIdx.x;
```

x s[i]=i;
__syncthreads();

x s[i]= s[i+1]; → temp=s[i];
__syncthreads();
s[i]=s[i+1];
__syncthreads();

✓ if(i%2)
 s[i]=s[i-1];
}

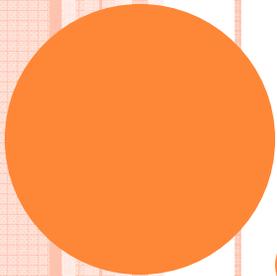
OUTLINES

- Threads communication patterns
- GPU hardware: Processing model
- GPU hardware: Memory model
- Efficient GPU programming

The background features a black field with vertical stripes of varying widths and colors, including a wide grey textured stripe on the left and a thin white stripe. Several orange circles of different sizes are scattered across the page, with one large circle on the left and another on the bottom right.

WRITING EFFICIENT PROGRAM

37



38



EFFICIENT MEMORY ACCESS



MAXIMIZE ARITHMETIC INTENSITY

A high-end GPU can do 3 trillion Floating point Operation Per Second (3 FLOPS)

Math

—
memory

- Maximize compute operations per thread.
- Minimize time spent on memory per thread.
 - (The access time not the amount)

MINIMIZE TIME SPENT ON MEMORY PER THREAD.

- The access time not the amount

Local memory > Shared memory >> Global memory

QUIZ

- Sort following variable based on their memory speed

```
__global__ void GPU_memory(float *out, float *in )
{

    int i= threadIdx.x;
    int j= threadIdx.y;

    const float pi= 3.1415;

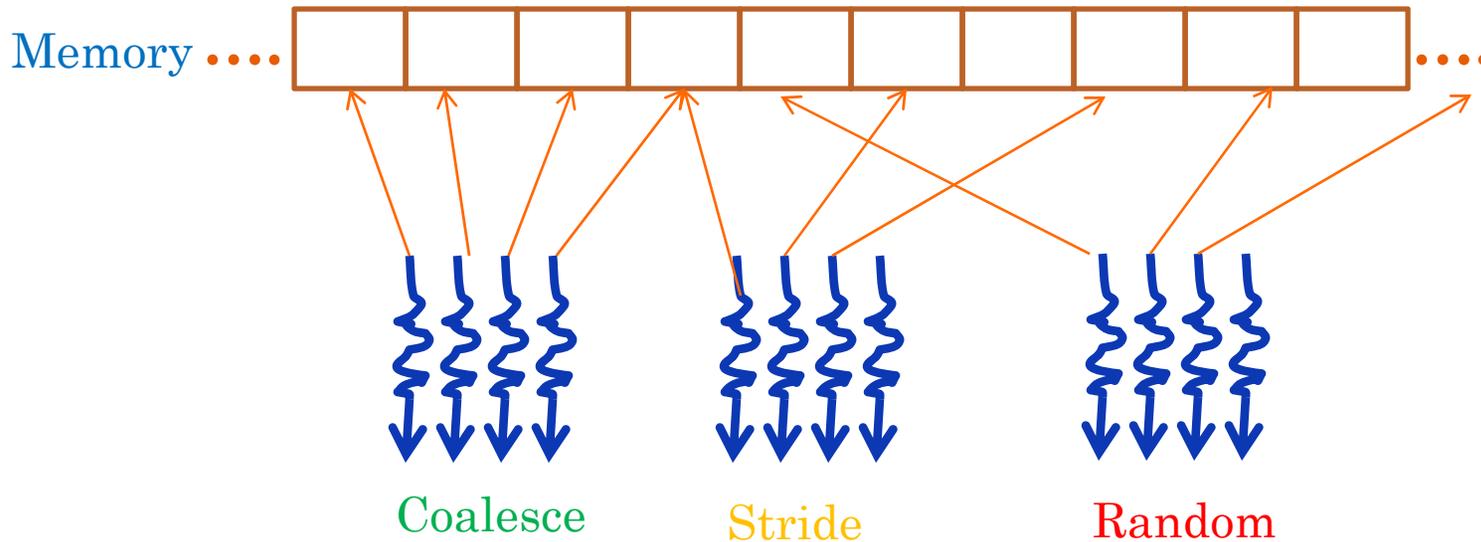
    out[i]= pi* in[i];

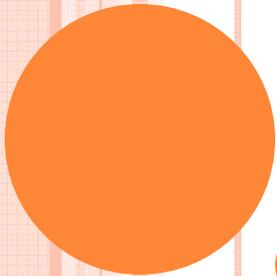
    __shared__ float sh_arr[128];

}
```

USING COALESCED GLOBAL MEMORY ACCESS

- We can read memory by block not by unites





43



ATOMIC OPERATION



ATOMIC OPERATIONS

- Special hardware to prevent more than one thread access the same memory
- Has specific operations (add, sub, xor) and works with specific data types (int)

Example

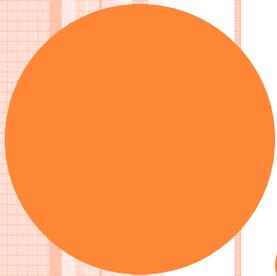
ATOMIC OPERATIONS LIMITATIONS

- Only certain operations
 - `atomicCAS()` has a broader use “Compare and swap”
- Still no ordering constraints
- Serializes access to memory (Slow)

QUIZ

State which sentence give the correct results and sort them based on memory access time :

- 1 million threads incrementing 1 million elements
- 1 million threads atomically incrementing 1 million elements
- 1 million threads incrementing 100 elements
- 1 million threads atomically incrementing 100 elements
- 10 million threads atomically incrementing 100 elements



THREAD DIVERGENCE

47



- Running threads with loops or conditional statement which diverse the time of each thread.



OUTLINES

- Threads communication patterns
- GPU hardware: Processing model
- GPU hardware: Memory model
- Efficient GPU programming