

Introduction to CUDA II

... curtsey of Massimo Bernaschi (CNR - <http://www.iac.cnr.it/~massimo>) &
John E. Stone (Univ. of Illinois at Urbana-Champaign - <http://www.ks.uiuc.edu/~johns/>)

Ivan Girotto – igirotto@ictp.it

Information & Communication Technology Section (ICTS)

International Centre for Theoretical Physics (ICTP)



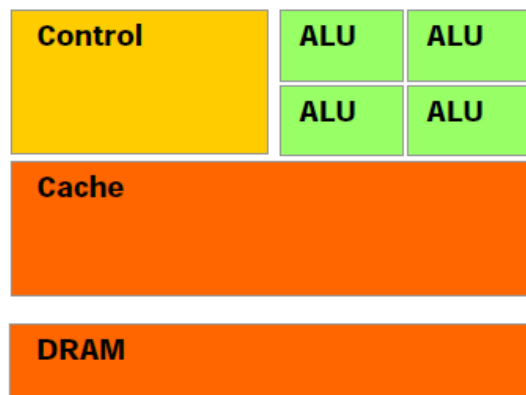
Scuola Internazionale Superiore
di Studi Avanzati



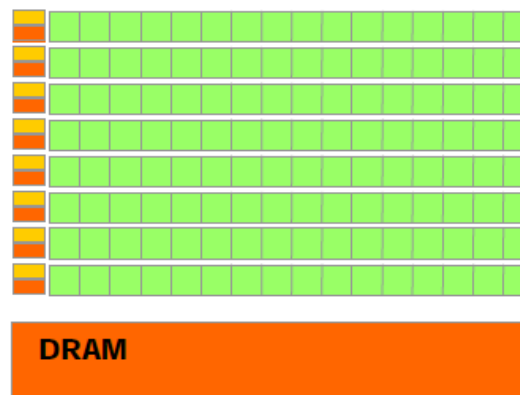
GPU Work Abstraction

- CUDA Kernels can be thought of as telling a GPU to compute all iterations of a set of nested loops concurrently
- Threads are dynamically scheduled onto hardware according to a hierarchy of thread groupings

CPU: Cache heavy, focused on individual thread performance



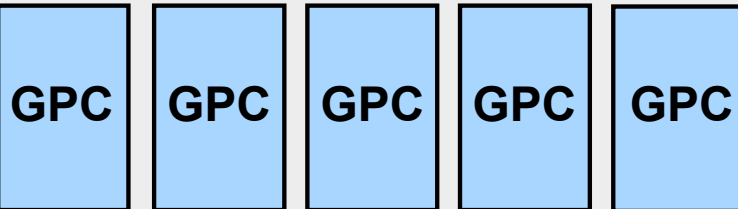
GPU: ALU heavy, massively parallel, throughput oriented



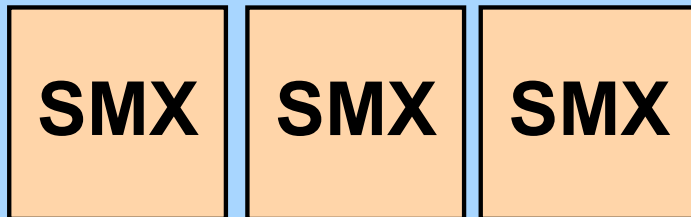
NVIDIA K20 GPU

3-12 GB DRAM Memory w/ ECC

1280KB - Level 2 - Cache



Graphics Processor Cluster

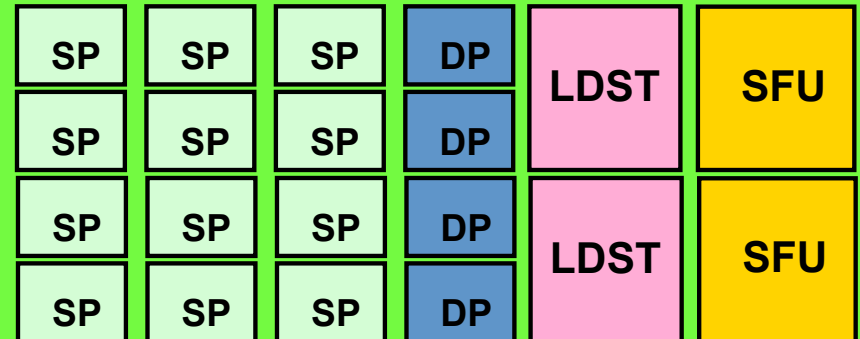


Streaming Multiprocessor - SMX

64 KB Constant Cache

64 KB L1 Cache / Shared Memory

48 KB Tex + Read-only Data Cache

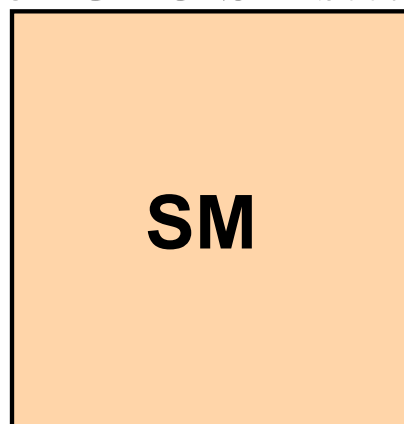


Tex Unit

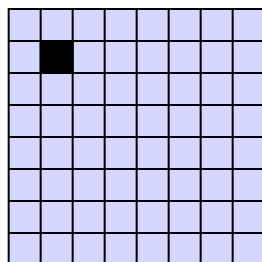
16 × Execution block =
192 SP, 64 DP,
32 SFU, 32 LDST

Grids, Thread Blocks, Threads

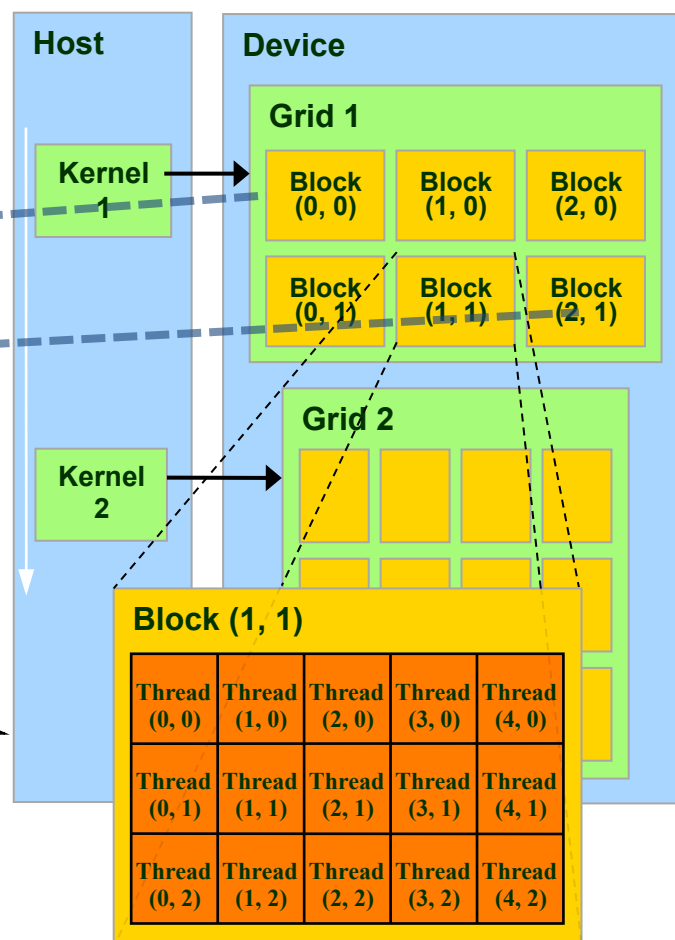
Thread blocks are scheduled onto pool of GPU SMs...



1-D, 2-D, 3-D thread block:

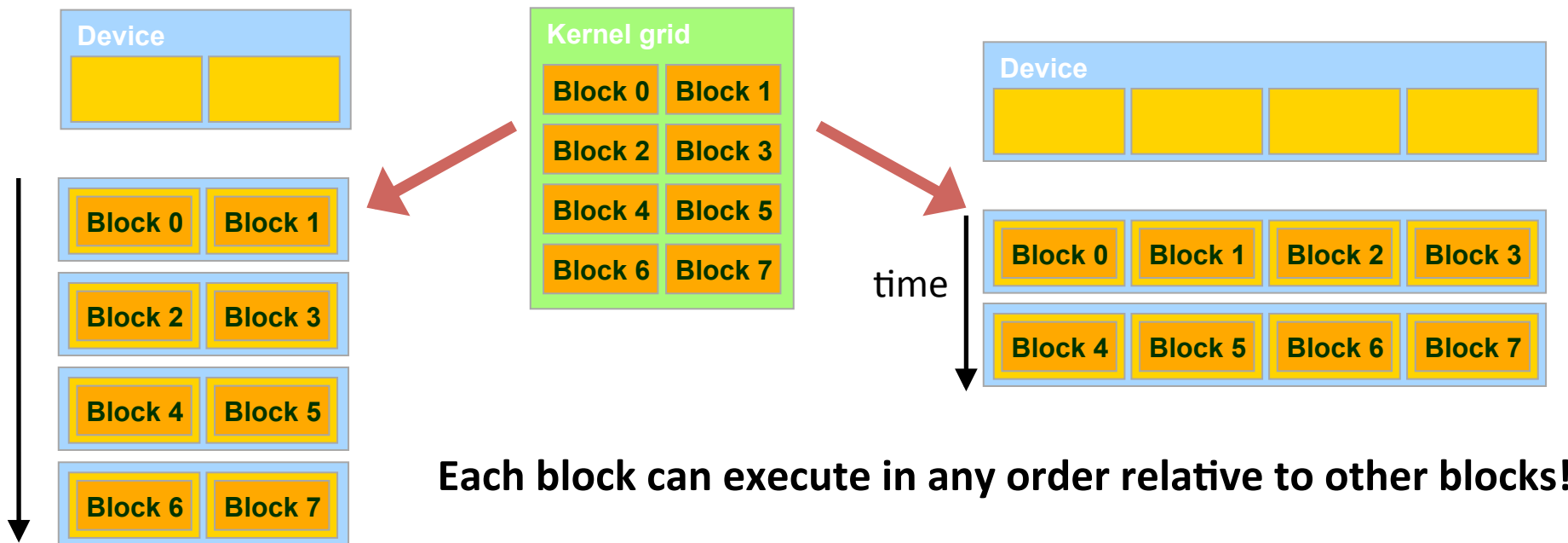


1-D, 2-D, or 3-D Grid of Ths blocks:



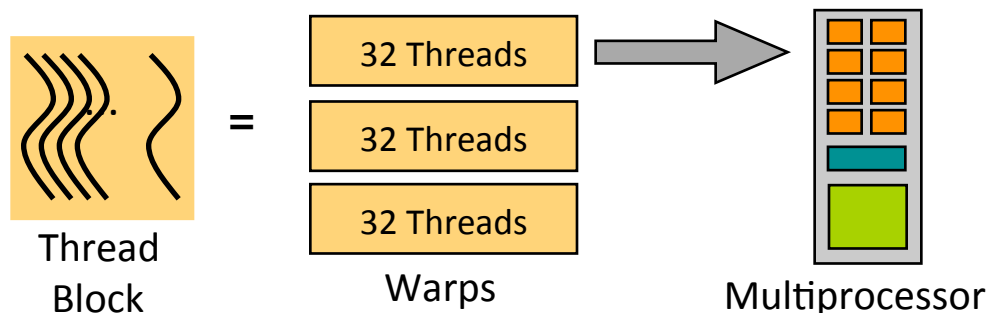
Transparent Scalability

- Hardware is free to assign blocks to any processor at any time
 - A kernel scales across any number of parallel processors

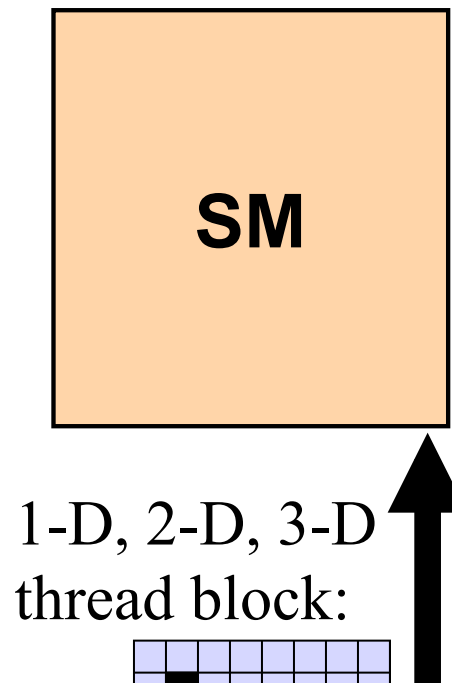


GPU Thread Block Execution

- Thread blocks are decomposed onto hardware in **32-thread “warps”**
- Hardware execution is scheduled in units of **warps**
 - an SM can execute warps from several thread blocks
- **Warps** run in SIMD-style execution:
 - All threads execute the same instruction in lock-step
 - If one thread stalls, the entire warp stalls...
 - A branch taken by a thread has to be taken by all threads...
(divergence is bad)



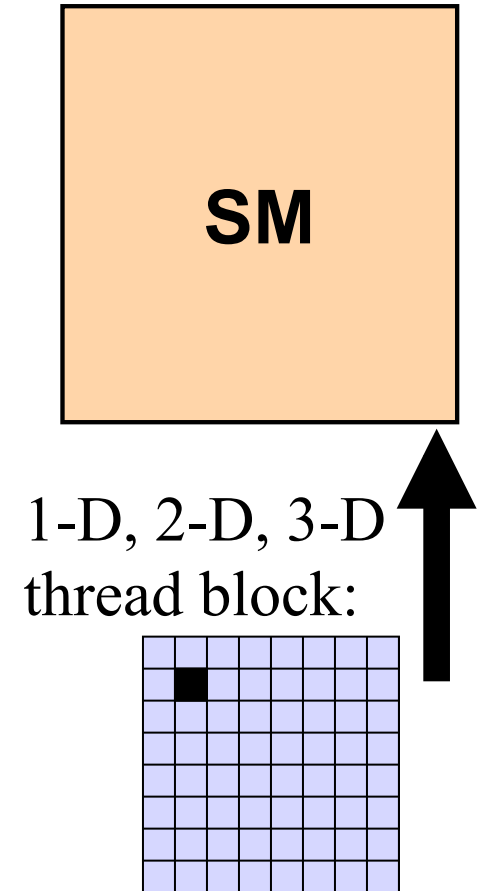
Thread blocks are multiplexed onto pool of GPU SMs...



GPU Warp Branch Divergence

- Branch divergence: when not all threads take the same branch, the entire warp has to **execute both sides of the branch**
- Branch divergence issue not unique to GPUs, affects **all** SIMD hardware platforms...
- On GPUs, we get fast **hardware-based** implementation of predication/masking/etc...
- GPU blocks memory writes from disabled threads in the “if then” branch, then inverts all thread enable states and runs the “else” branch
- GPU hardware detects warp re-convergence and then runs with all threads enabled...

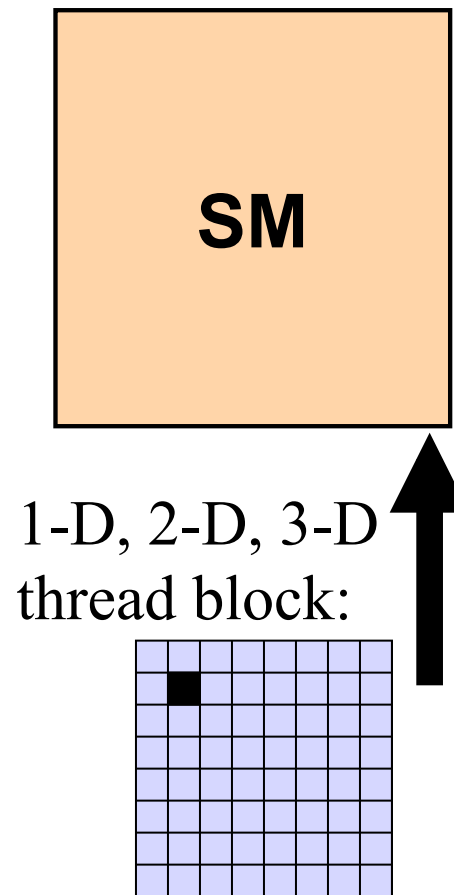
Thread blocks are multiplexed onto pool of GPU SMs...



GPU Warp Branch Divergence

- Threads within the same thread block can communicate with each other in fast on-chip shared memory
- Once scheduled on an SM, thread blocks run until completion
- Because the order of thread block execution is arbitrary and blocks cannot be stopped, they cannot communicate or synchronize with other thread blocks (*)
- (*) Atomic memory ops are an exception wrt/communication

Thread blocks are multiplexed onto pool of GPU SMs...



Execution Model

Software

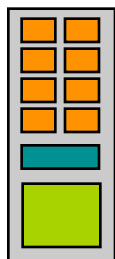
Thread



Thread
Block

Hardware

Scalar
Processor



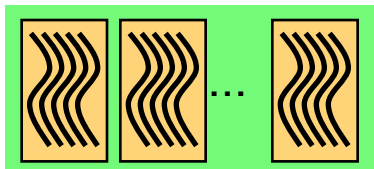
Multiprocessor

Threads are executed by scalar processors

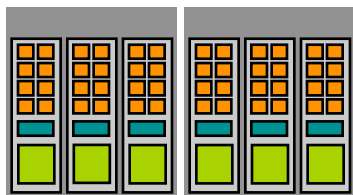
Thread blocks are executed on multiprocessors

Thread blocks **do not migrate**

Several concurrent thread blocks can reside on one multiprocessor - limited by multiprocessor resources (shared memory and register file)



Grid

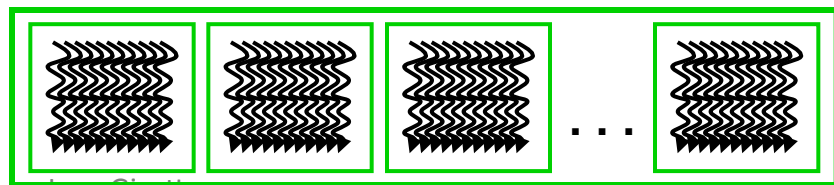
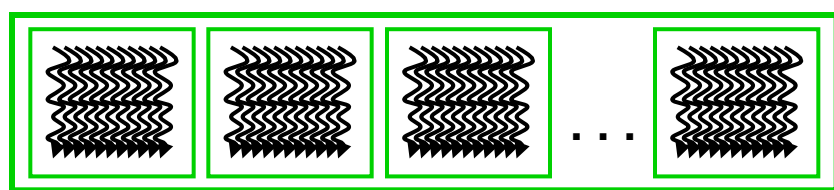
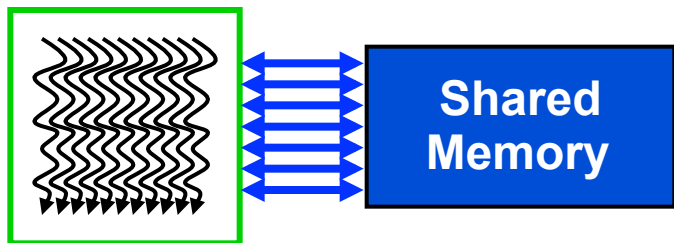


Device

A kernel is launched as a grid of thread blocks

Memory Hierarchy

- Registers (fast up to availability)
- Local Memory: per-thread
 - Private per thread but *slow*!
 - Auto variables, register spill
- Shared Memory: per-block
 - Shared by threads of the same block
 - *Fast* inter-thread communication
- Global Memory: per-application
 - Shared by all threads
 - Inter-Grid communication



**Sequential Grids
Execution in Time**



Atomic Operations

- Terminology: Read-modify-write uninterruptible when *atomic*
- Many *atomic operations* on memory available with CUDA C
 - `atomicAdd()` ▪ `atomicInc()`
 - `atomicSub()` ▪ `atomicDec()`
 - `atomicMin()` ▪ `atomicExch()`
 - `atomicMax()` ▪ `atomicCAS() old == compare ? val : old`
- Predictable result when simultaneous access to memory required

Multiblock Dot Product: dot ()

```
__global__ void dot( int *a, int *b, int *c ) {  
    __shared__ int temp[THREADS_PER_BLOCK];  
  
    int index = threadIdx.x + blockIdx.x * blockDim.x;  
  
    temp[threadIdx.x] = a[index] * b[index];  
  
    __syncthreads();  
  
    if( 0 == threadIdx.x ) {  
        int sum = 0;  
        for( int i = 0; i < THREADS_PER_BLOCK; i++ ) sum += temp[i];  
        atomicAdd( c , sum );  
    }  
}
```

- We need to atomically add **sum** to **c** in our multiblock dot product

Built-in Variables to manage grids and blocks

dim3 => a new datatype defined by CUDA:

- **struct dim3 { unsigned int x, y, z };**
- three unsigned ints where any unspecified component defaults to 1.

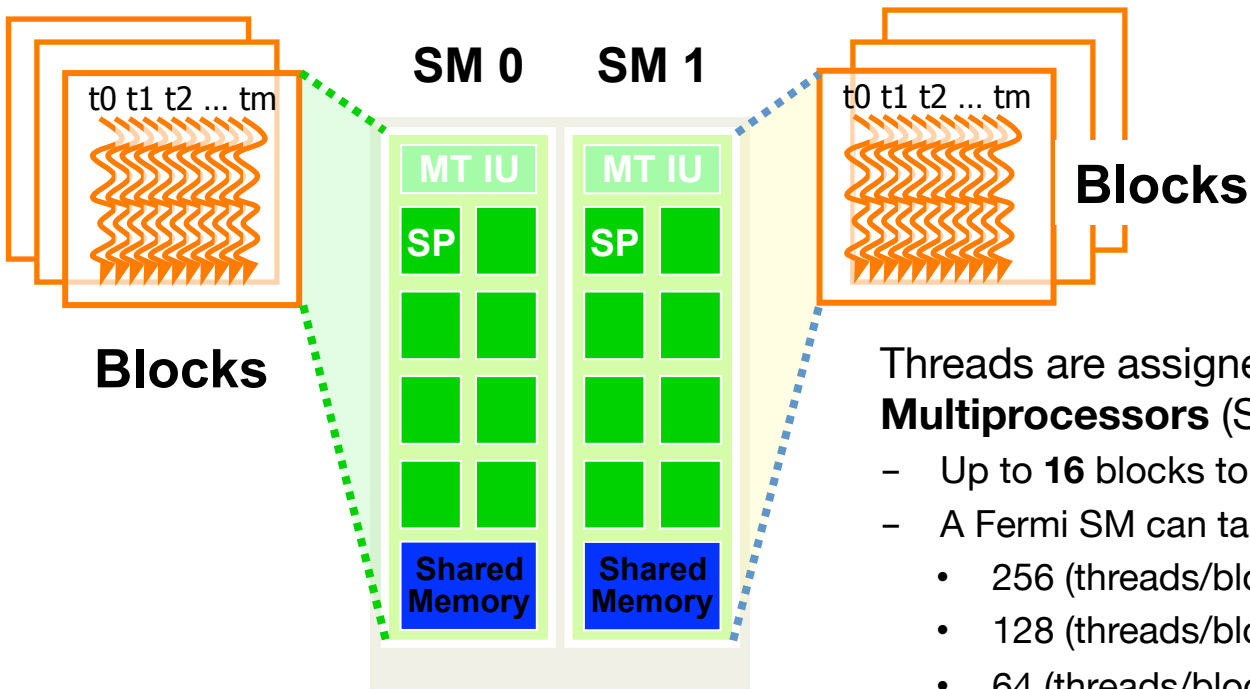
- **dim3 gridDim;**
 - Dimensions of the grid in blocks
- **dim3 blockDim;**
 - Dimensions of the block in threads
- **dim3 blockIdx;**
 - Block index within the grid
- **dim3 threadIdx;**
 - Thread index within the block

Bi-dimensional threads configuration: set the elements of a square matrix

```
__global__ void kernel( int *a, int dimx, int dimy ) {  
    int ix  = blockIdx.x*blockDim.x + threadIdx.x;  
    int iy  = blockIdx.y*blockDim.y + threadIdx.y;  
    int idx = iy*dimx + ix;  
  
    a[idx] = idx+1;  
}
```

```
int main() {  
    int dimx = 16;  
    int dimy = 16;  
    int num_bytes = dimx*dimy*sizeof(int);  
  
    int *d_a=0, *h_a=0; // device and host pointers  
  
    h_a = (int*)malloc(num_bytes);  
    cudaMalloc( (void**)&d_a, num_bytes );  
  
    dim3 grid, block;  
    block.x = 4;  
    block.y = 4;  
    grid.x = dimx / block.x;  
    grid.y = dimy / block.y;  
  
    kernel<<<grid, block>>>( d_a, dimx, dimy );  
  
    cudaMemcpy(h_a,d_a,num_bytes,  
               cudaMemcpyDeviceToHost);  
  
    for(int row=0; row<dimy; row++) {  
        for(int col=0; col<dimx; col++)  
            printf("%d ", h_a[row*dimx+col] );  
        printf("\n");  
    }  
  
    free( h_a );  
    cudaFree( d_a );  
    return 0;  
}
```

Executing Thread Blocks



Threads are assigned to **Streaming Multiprocessors** (SM) in block granularity

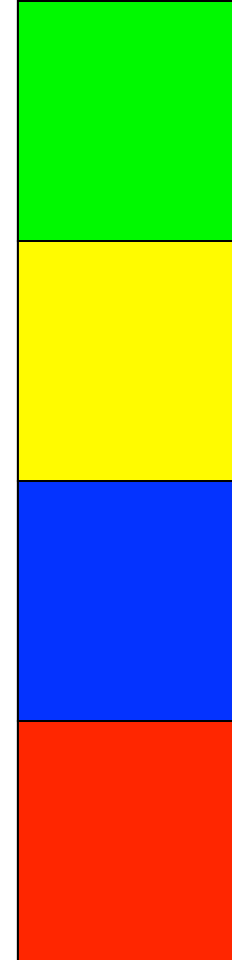
- Up to **16** blocks to each SM on K20.
- A Fermi SM can take up to 2048 threads (examples):
 - 256 (threads/block) * 8 blocks
 - 128 (threads/block) * 16 blocks
 - 64 (threads/block) * 32 blocks, **not allowed!**

- Threads run concurrently SM manages/schedules thread execution
- The number of threads in a block depends on the capability => **1024** threads on K20.

Programmer View of Register File

- There are up to **65536** (32 bit) registers in each Kepler SM
 - This is an implementation decision, not part of CUDA
 - Registers are dynamically partitioned across all blocks assigned to the SM
 - Once assigned to a block, the register is NOT accessible by threads in other blocks
 - Each thread in the same block only access registers assigned to itself

4 blocks



3 blocks



GPU On-Board Global Memory

GPU arithmetic rates dwarf memory bandwidth

For Kepler K40 hardware:

~4.3 SP TFLOPS vs. ~288 GB/sec

The ratio is roughly **60 FLOPS per memory reference**
for single-precision floating point

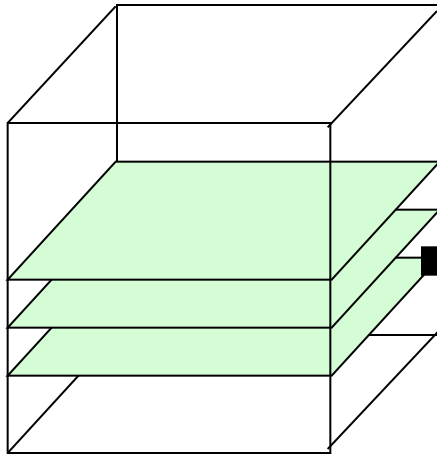
Peak performance achieved with **“coalesced”** memory access patterns –
patterns that result in a single hardware memory transaction for a SIMD
“warp” – a contiguous group of 32 threads

Memory Coalescing (Oversimplified explanation)

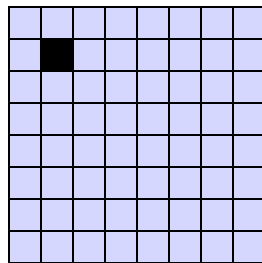
- Threads in a warp perform a read/write operation that can be serviced in a single hardware transaction
- Rules vary slightly between hardware generations, but new GPUs are much more flexible than old ones
- If all threads in a warp read from a contiguous region that's 32 items of 4, 8, or 16 bytes in size, that's an example of a coalesced access
- Multiple threads reading the same data are handled by a hardware broadcast
- Writes are similar, but multiple writes to the same location yields undefined results

CUDA Grid/Block/Thread Decomposition

**1-D, 2-D, or 3-D
Computational Domain**

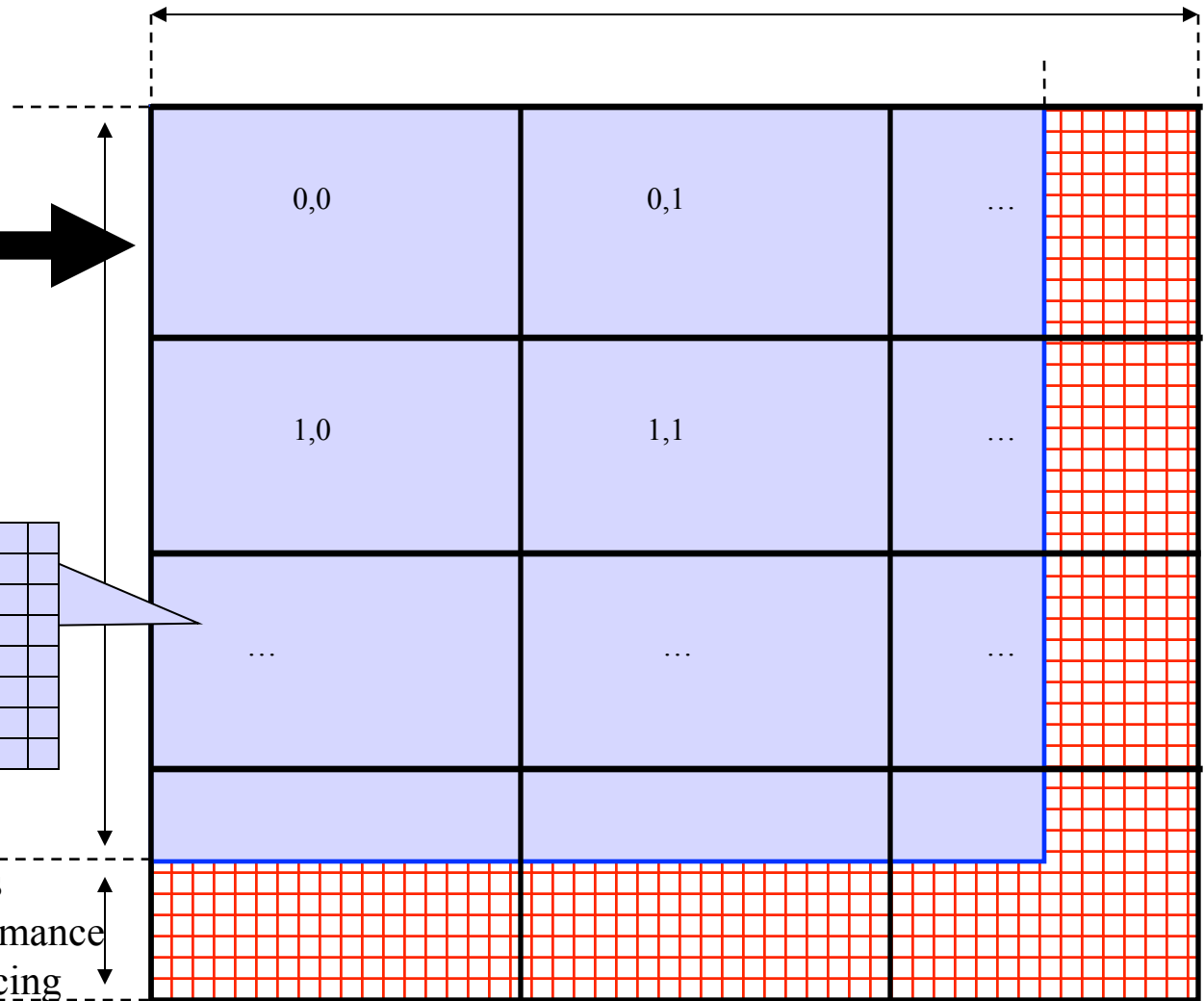


**1-D, 2-D, 3-D
thread block:**



Padding arrays out to full blocks
optimizes global memory performance
by guaranteeing memory coalescing

**1-D, 2-D, or 3-D (SM \geq 2.x)
Grid of thread blocks:**



Optimizing threads per block

- ✓ Choose threads per block as a multiple of warp size (32)
 - ✓ Avoid wasting computation on under-populated warps
 - ✓ Facilitates efficient memory access (*coalescing*)
- ✓ Run as many warps as possible per multiprocessor (hide latency)
 - ✓ SM can run up to 16 (on Kepler) blocks at a time
- ✓ Heuristics
 - ✓ Minimum: 64 threads per block
 - ✓ 192 or 256 threads a better choice
 - ✓ Usually still enough registers to compile and invoke successfully
- ✓ The right tradeoff depends on your computation, so **experiment, experiment, experiment!!!**

CUDA Compiler: nvcc basic options

- `-arch=sm_35` → enable code for a given capability
- `-G` → enable debug for device code
- `--ptxas-options=-v` → show register and memory usage
- `--maxrregcount <N>` → limit the number of registers
- `-use_fast_math` → use fast math library
- `-O3` → Enables compiler optimization
- `-ccbin compiler_path` → use a different C compiler
- `--compiler-options` → Specify options directly to the compiler/preprocessor.

Getting Performance From GPUs

- Don't worry (much) about counting arithmetic operations...at least until you have nothing else left to do
- GPUs provide tremendous memory bandwidth, but even so, **memory bandwidth often ends up being the performance limiter**
- Keep/reuse data in **registers** as long as possible
- The main consideration when programming GPUs is **accessing memory efficiently**, and storing operands in the **most appropriate memory system** according to data size and access pattern

References

- <http://www.ks.uiuc.edu/Research/gpu/>
- <http://indico.ictp.it/event/a14302/other-view?view=ictp timetable> (ICTP - SMR2760)
- <http://www.iac.rm.cnr.it/~massimo/PMC.html>
- CUDA Zone: <https://developer.nvidia.com/cuda-zone>
- CUDA by Example

