

# GPU programming



## 27-th Summer School on Parallel Computing

*Rome edition*

*9 - 20 July 2018*

***Sergio Orlandini***

s.orlandini@ Cineca.it

# ■ Memory Hierarchy on CUDA

- *Global Memory*
  - *caches*
  - *type of global memory accesses*
- *Shared Memory*
  - Matrix-Matrix Product using *Shared Memory*
- *Constant Memory*
- *Texture Memory*
- *Registers and Local Memory*



# Memory Hierarchy

All CUDA threads in a block have access to:

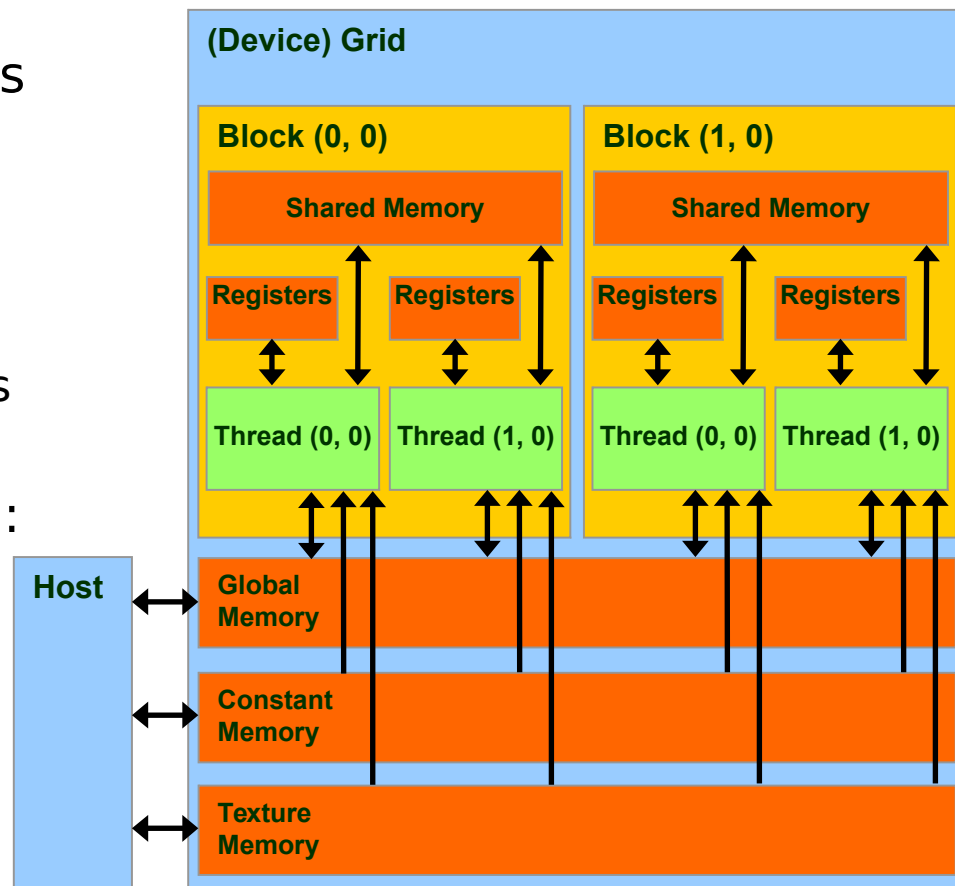
- resources of the SM assigned to its block:
  - **Registers**
  - **Shared Memory**

NB: thread belonging to different blocks cannot share these resources

- all memory type available on GPU:
  - **Global Memory**
  - **Constant Memory** (read only)
  - **Texture Memory** (read only)

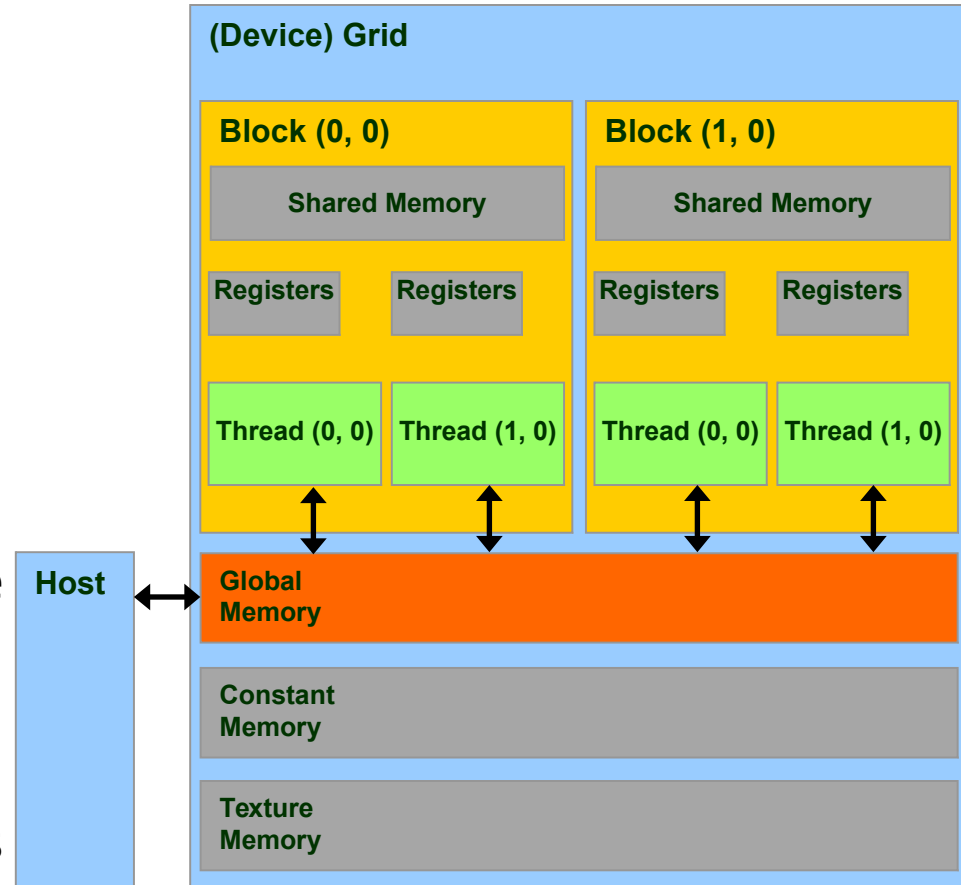
NB: CPU can access and initialize both constant and texture memory

NB: global, constant and texture memory have persistent storage duration



# Global Memory

- **Global Memory** is the larger memory available on a *device*
  - Comparable to a RAM for CPU
  - Its status is maintained among different kernel launches
  - Can be access both read/write from all threads of the kernel grid
  - Unique memory that can be use in read/write access from the CPU
  - **Very high bandwidth**  
Throughput up to 240-760 GB/s
  - **Very high latency**  
about 400-800 clock cycles



# Declare Variable in *Global Memory*

- How to allocate a variable in Global Memory:

```
__device__ type variable_name; // static
```

*or dynamic allocation*

```
type *pointer_to_variable;  
cudaMalloc((void **) &pointer_to_variable, size);  
cudaFree(pointer_to_variable);
```

```
type, device :: variable_name
```

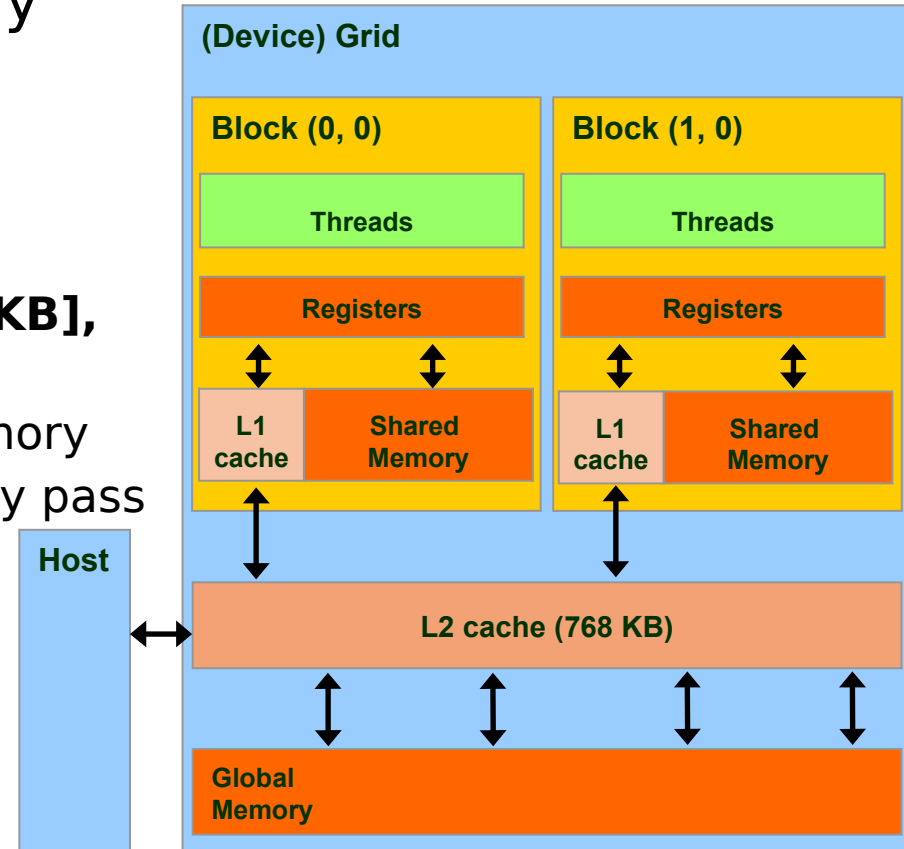
*or dynamic allocation*

```
type, device, allocatable :: variable_name  
allocate(variable_name, size)  
deallocate(variable_name)
```

- Lifetime of the application
- Accessible by all threads of a CUDA grid and by the host

# Cache Hierarchy for *Global Memory*

- Starting with the Fermi architecture, a cache hierarchy has been introduced
- 2 Levels of cache:
  - L2** : share among all SM
    - Fermi [768 KB], Kepler [1536 KB], Pascal [4MB]**
    - 25% less latency than Global Memory
    - NB : all accesses to global memory pass through L2 cache, also H2D/D2H memory transfers
  - L1** : private to each SM
    - [16/48 KB]** configurable
    - L1 + Shared Memory = 64 KB
    - Kepler/Pascal**: configurable at **32 KB**



```
cudaFuncSetCacheConfig(kernel1, cudaFuncCachePreferL1);          // 48KB L1 / 16KB ShMem  
cudaFuncSetCacheConfig(kernel2, cudaFuncCachePreferShared);      // 16KB L1 / 48KB ShMem
```

# Cache Hierarchy for *Global Memory*

Two different types of **load** operations:

- **Caching (default mode)**

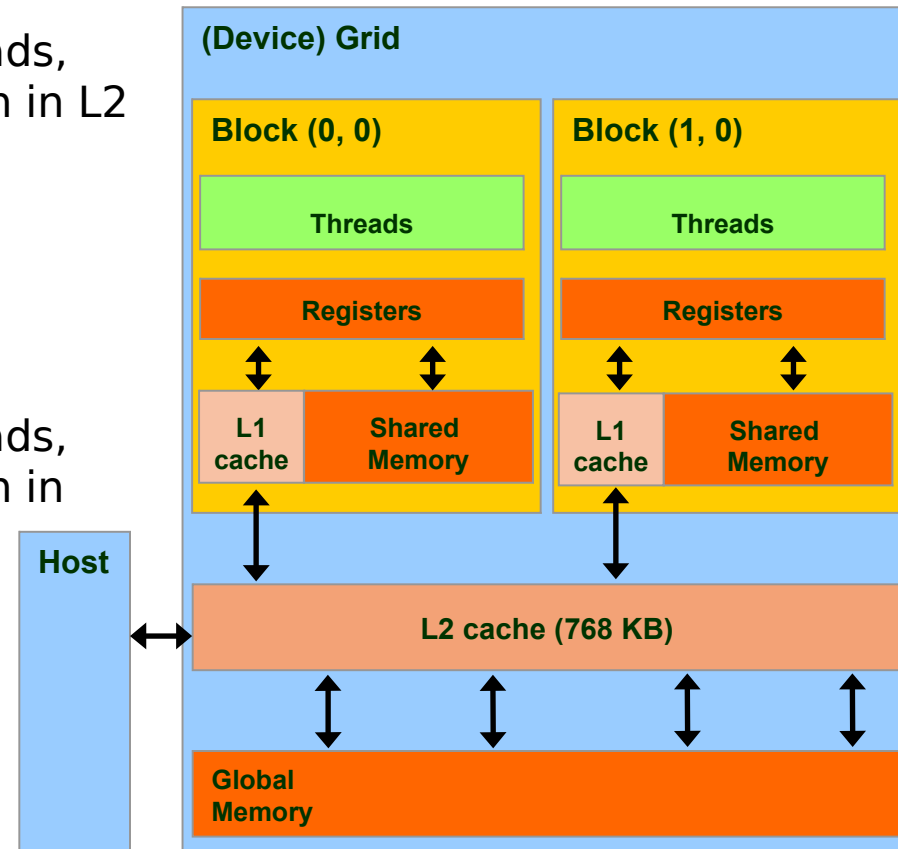
- when data is requested by some threads, data is first searched in L1 cache, then in L2 cache, then in global memory
- cache line length is **128-byte**

- **Non-caching**

- L1 cache is disabled
- when data is requested by some threads, data is first searched in L2 cache, then in global memory
- cache line length is **32-bytes**
- Activated at *compile time* with option:  
-Xptxas -dlcm=cg

Just one type of **store** operation:

- when data should be store in global memory, its L1 copy is invalidated and L2 cache value is updated



# Global Memory Load/Store

```
// strided data copy
__global__ void strideCopy (float *odata, float* idata, int stride) {
    int xid = (blockIdx.x*blockDim.x + threadIdx.x) * stride;
    odata[xid] = idata[xid];
}
```

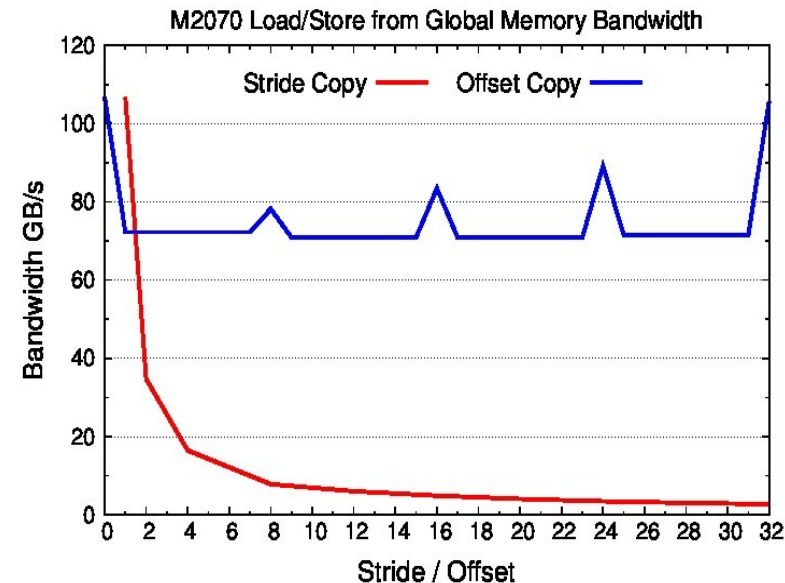
```
// offset data copy
__global__ void offsetCopy(float *odata, float* idata, int offset) {
    int xid = blockIdx.x * blockDim.x + threadIdx.x + offset;
    odata[xid] = idata[xid];
}
```

## Strided copy

Stride	Bandwidth GB/s
1	106.6
2	34.8
8	7.9
16	4.9
32	2.7

## Offset copy

Offset	Bandwidth GB/s
0	106.6
1	72.2
8	78.2
16	83.4
32	105.7



Measured on a M2070; Total elements = 16776960; Num. Blocks = 65535; Block length = 256



# Loads from *Global Memory*

- All load/store request in global memory are issued per *warp* (as all other instructions)
  1. each *thread* in a *warp* compute the address to access
  2. *load/store* units calculate in which memory segments data resides
  3. *load/store* units start up requests for segment to transfer

**Warp requires 32 consecutive 4-byte word aligned to segment (total 128 bytes)**

## Caching Load

addresses fall within 1 cache line

128 bytes are moved across the bus

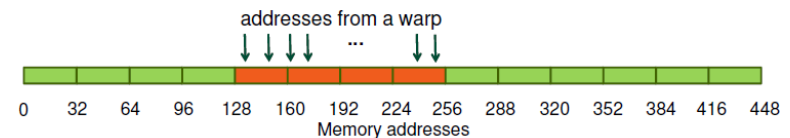
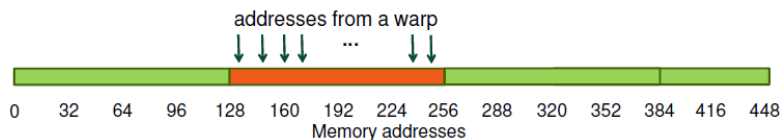
bus utilization: **100%**

## Non-caching Load

addresses fall within 4 cache segments

128 bytes are moved across the bus

bus utilization: **100%**



# Loads from *Global Memory*

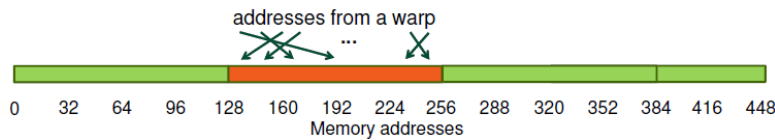
**Warp requests 32 permuted 4-byte words aligned to a segment (total 128 bytes)**

## Caching Load

addresses fall within 1 cache line

128 bytes are moved across the bus

bus utilization: **100%**

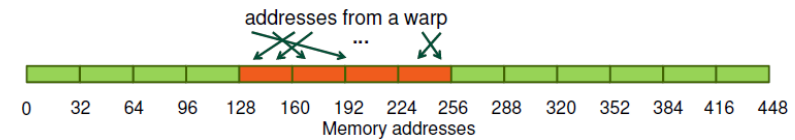


## Non-caching Load

addresses fall within 4 cache segments

128 bytes are moved across the bus

bus utilization: **100%**



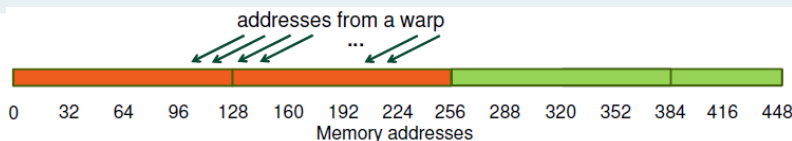
**Warp requests 32 consecutive 4-bytes words not aligned to a segment (total 128 bytes)**

## Caching Load

addresses fall within 2 cache lines

256 bytes are moved across the bus

bus utilization: **50%**

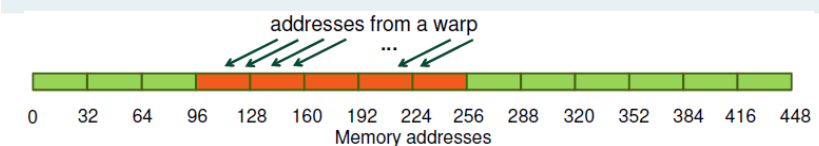


## Non-caching Load

addresses fall within at most 5 segments

256 bytes are moved across the bus

bus utilization: at least **80%**



# Loads from *Global Memory*

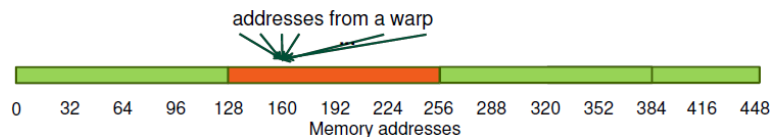
**All threads in a warp request the same 4-byte word (total 4 bytes)**

## Caching Load

addresses fall within a single cache line

128 bytes are moved across the bus

bus utilization: **3.125%**

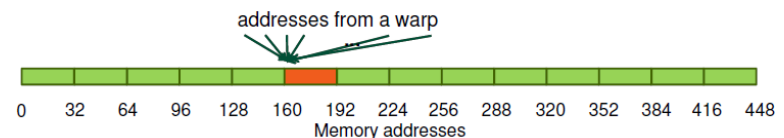


## Non-caching Load

addresses fall within a single segment

32 bytes are moved over the bus

bus utilization: **12.5%**



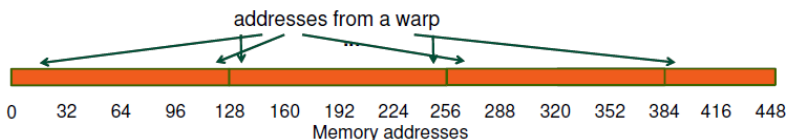
**Warp requests 32 not contiguous 4-bytes words (total 128 bytes)**

## Caching Load

addresses fall within N different cache lines

$N \times 128$  bytes are moved across the bus

bus utilization:  **$128 / (N \times 128)$**

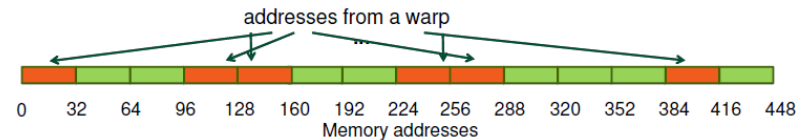


## Non-caching Load

addresses fall within N different segments

$N \times 32$  bytes are moved across the bus

bus utilization:  **$128 / (N \times 32)$**



# Data alignment in *Global Memory*

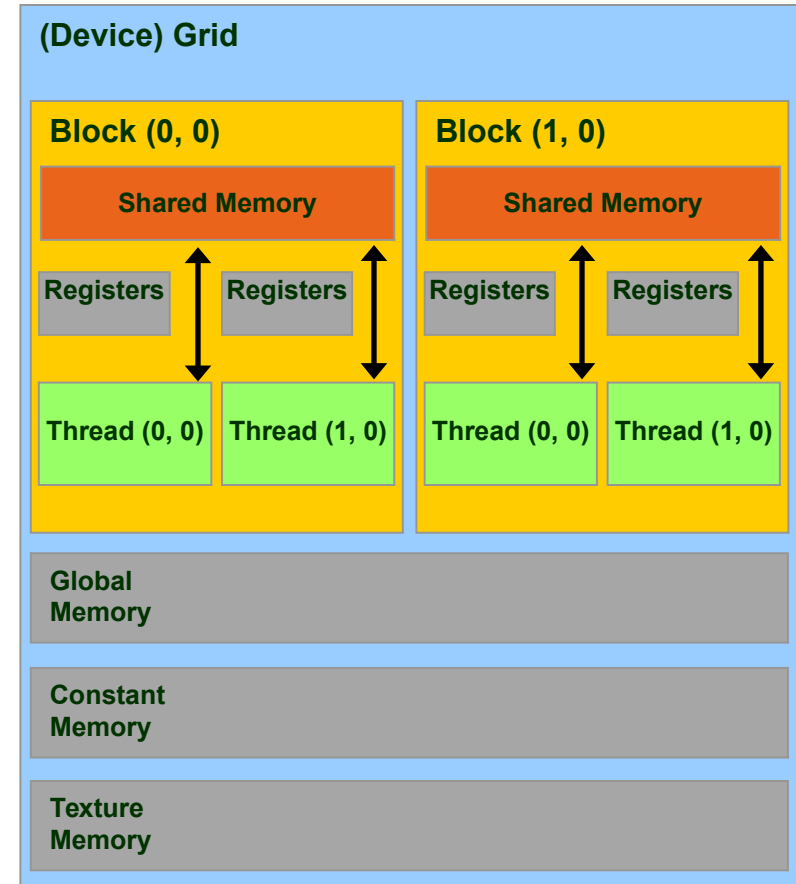
- It is very important to align data in memory so to have aligned accesses (*coalesced*) during load/store operation in global memory, reducing the number of bytes moved across the bus
  - **cudaMalloc()** grants the alignment of first element in global memory, useful for one dimensional arrays
  - **cudaMallocPitch()** must be used to allocate 2D buffers
    - elements are padded so each row is aligned for coalescing accesses
    - returns an integer (*pitch*) which can be used as a stride to access row elements

```
// host code
int width = 64, height = 64;
float *devPtr;
int pitch;
cudaMallocPitch(&devPtr, &pitch, width * sizeof(float), height);

// device code
__global__ myKernel(float *devPtr, int pitch, int width, int height)
{
    for (int r = 0; r < height; r++) {
        float *row = devPtr + r * pitch;
        for (int c = 0; c < width; c++)
            float element = row[c];
    }
    ...
}
```

# Shared Memory

- The **Shared Memory** is a small, but quite fast memory mounted on each SM
  - Accessible in read/write mode for only threads of a block
  - Alike a cache memory under the direct control of the programmer
  - Its status is not maintained among different kernel calls
- Specifications:
  - **Very low latency**: 2 clock cycles
  - Throughput: 32 bit every 2 cycles
  - Dimension : **48 KB** [default]  
(Configurable : 16/48 KB)
  - **Kepler** : also **32 KB**



# Shared Memory Allocation

```
// statically inside the kernel
__global__ myKernelOnGPU (...) {
    ...
    __shared__ type shmem[MEMSZ];
    ...
}
```

or dynamic allocation

```
// dynamically sized
extern __shared__ type *dynshmem;

__global__ myKernelOnGPU (...) {
    ...
    dynshmem[i] = ... ;
    ...
}
```

```
void myHostFunction() {
    ...
    myKernelOnGPU<<<gs,bs,MEMSZ>>>();
}
```

```
! statically inside the kernel
attribute(global)
subroutine myKernel(...)
    ...
    type, shared:: variable_name
    ...
end subroutine
```

or dynamic allocation

```
! dynamically sized
type, shared:: dynshmem(*)

attribute(global)
subroutine myKernel(...)
    ...
    dynshmem(i) = ...
    ...
end subroutine
```

- Lifetime of CUDA block of threads (NOT persistent along kernel launch!)
- Accessible only by threads of the same block

# Thread Block Synchronization

- All threads in the same block can be synchronized using the CUDA runtime API:

`__syncthreads()` | `call syncthreads()`

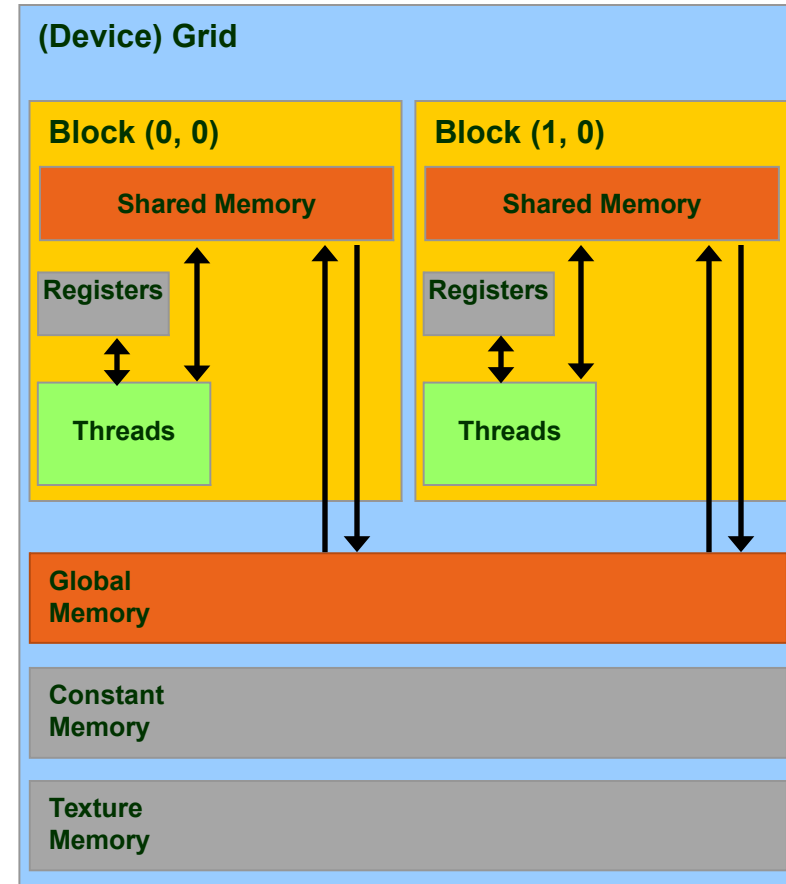
which blocks execution until all other threads reach the same call location

- NB: can be used in conditional too, but only if all thread in the block reach the same synchronization call

*“... otherwise the code execution is likely to hang or produce unintended side effects”*

# Shared Memory - Thread Cooperation

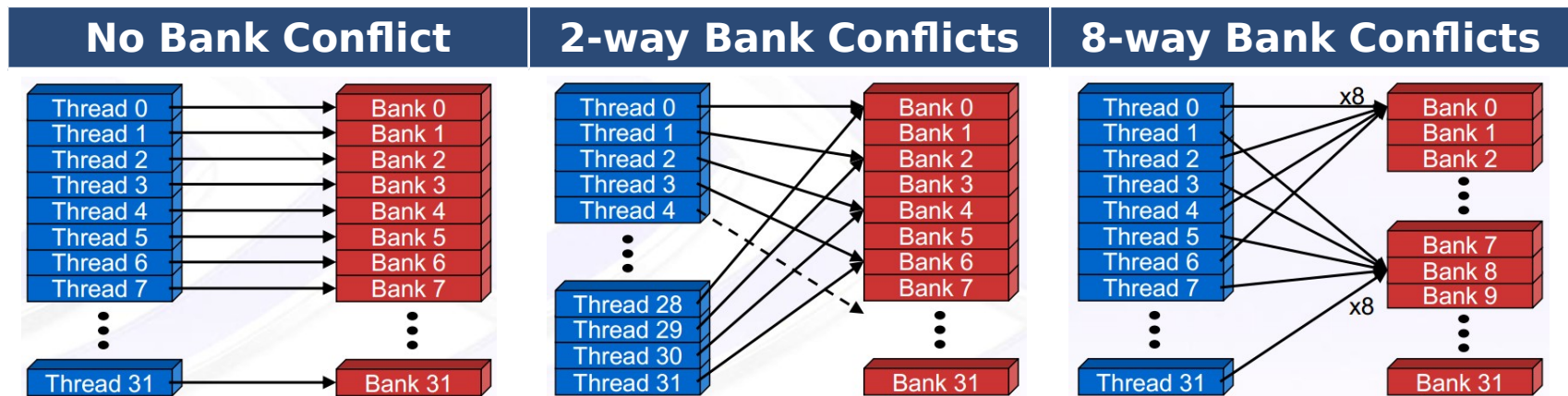
- Threads belonging to the same block can cooperate together using the shared memory to share data
  - if a thread needs some data which has been already retrieved by another thread in the same block, this data can be shared using the shared memory
- Typical Shared Memory usage:
  - declare a buffer residing on shared memory (this buffer is per block)
  - load data into shared memory buffer
  - synchronize threads so to make sure all needed data is present in the buffer
  - perform operation on data
  - synchronize threads so all operations have been performed
  - write back results to global memory





# Shared Memory and Bank Accesses

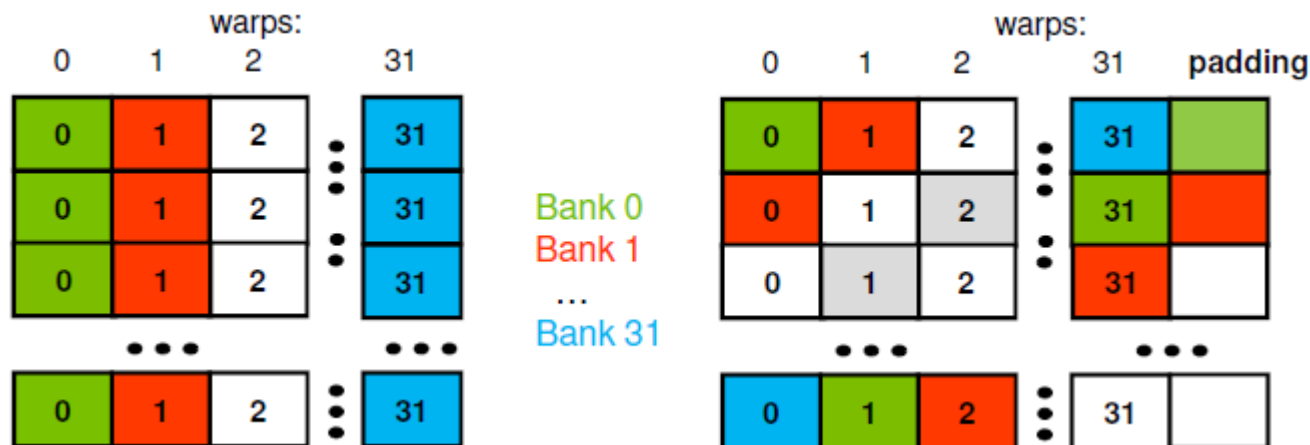
- Shared memory has 32 banks organized such that 32-bit words map a banks
  - Data are distributed every 4-bytes cycling over successive banks
  - Shared memory accesses are per warp
  - Multicast** : if N threads of the same warp request the same element, access is executed with only one transaction
  - Broadcast** : if ALL threads of the same warp request the same element, access is executed with only one transaction
  - Bank Conflict**: if two or more threads requests different data belonging to the same bank, each access is serialized



# Avoid Bank Conflict

- A naive implementation of CUDA kernels using shared memory would use a tile of size 32x32 floats
  - each element resides on a single bank (4-byte)
  - data are on the same bank every 32 floats
  - so read/write by columns will turn into the worst type of bank conflict
- Use a common trick: let's size the tile using 33 elements
  - now all elements belonging to the same column reside on different banks

```
__shared__ float tile[TILE_DIM][TILE_DIM+1];
```



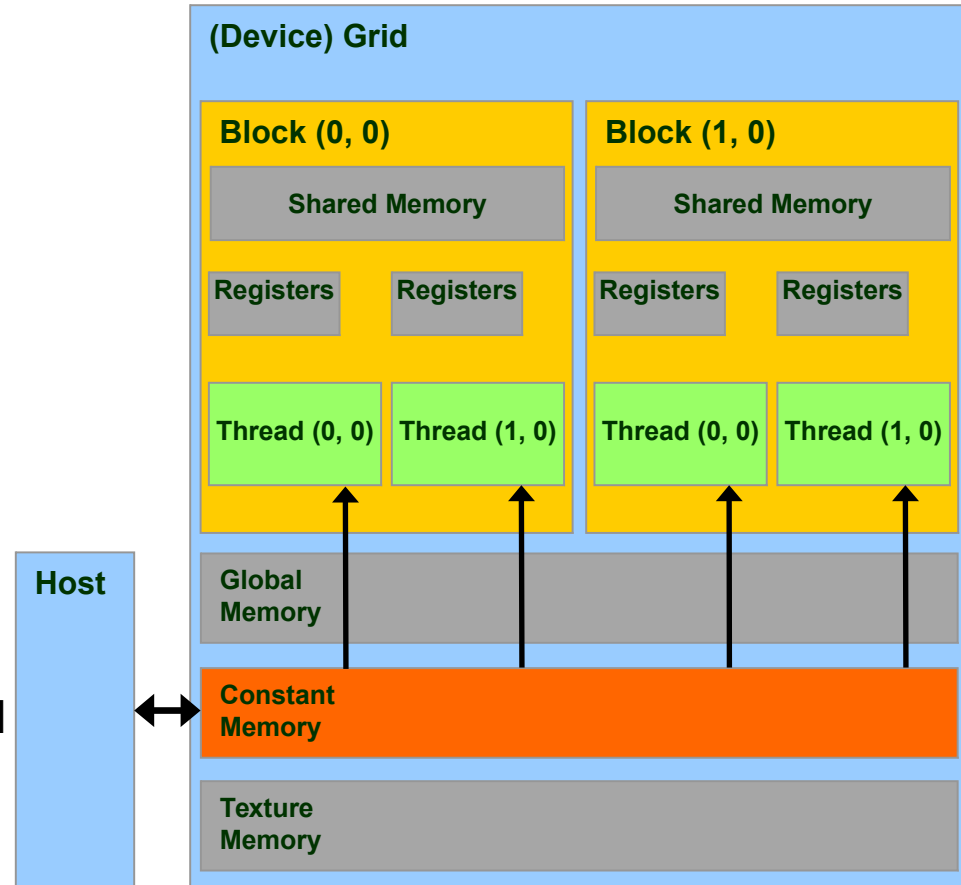
# Constant Memory

▪ **Constant Memory** is the ideal place to store constant data in **read-only** access from all threads

- constant memory data actually reside in the global memory, but fetched data is moved into a dedicated *constant-cache*
- very efficient when all *thread* of a *warp* request the same memory address
- Constant memory is initialized from host code using a special CUDA API

▪ Specifications:

- Dimension : **64 KB**
- Throughput: 32 bits per warp every 2 clock cycles

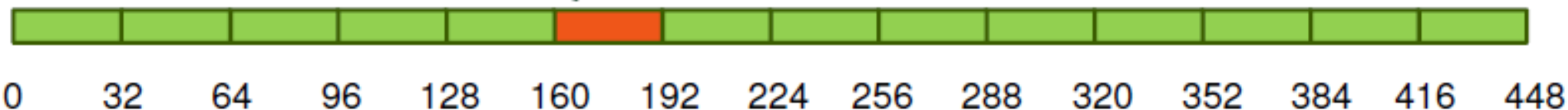


# Accessing *Constant Memory*

Suppose a kernel is launched using 320 warps per SM and all threads requests the same data

- if data is on global memory:
  - all *warp* will request the same segment from global memory
  - the first time segment is copied into L2 cache
  - if other data pass through L2, there are good chances it will be lost
  - there are good chances that data should be requested 320 times
- if data is in constant memory:
  - during first *warp* request, data is copied in *constant-cache*
  - since there is less traffic in *constant-cache* , there are good chances *all other warp will find the data already in cache*, so no more traffic on the BUS

addresses from a warp



# Constant Memory Allocation

```
__constant__  type  variable_name; // static

cudaMemcpyToSymbol(const_mem, &host_src, sizeof(type), cudaMemcpyHostToDevice);

// warning
// cannot be dynamically allocated
```

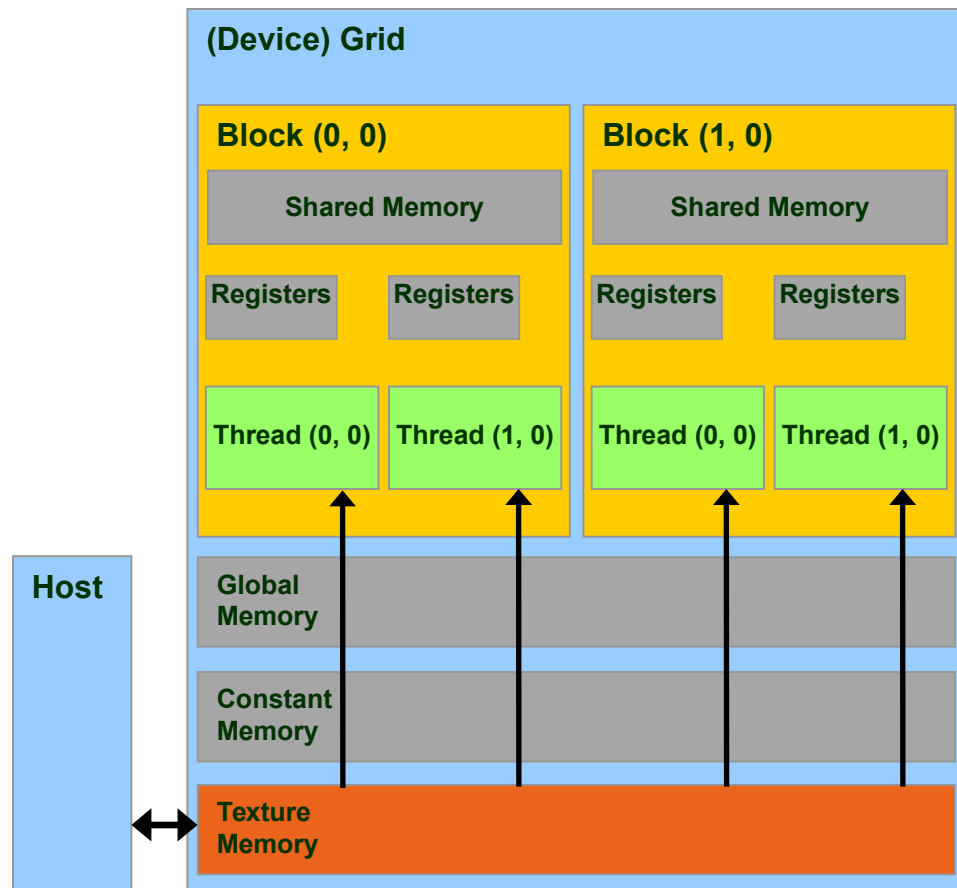
```
type, constant :: variable_name

! warning
! cannot be dynamically allocated
```

- data will reside in the constant memory address space
- has static storage duration (persists until the application ends)
- readable from all threads of a kernel

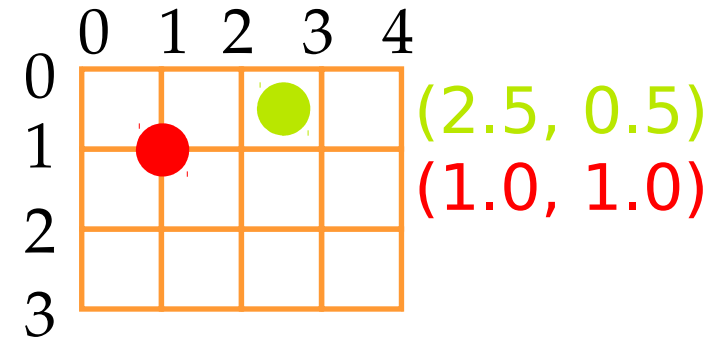
# Texture Memory

- **Texture Memory** is a basic graphic rendering functionality
- as for constant memory, data actually reside in global memory, but is fetched across a dedicated texture-cache
- data is accessed in **read-only** using special CUDA API function, called **texture fetch**
- Specifications:
  - address resolution is more efficient since it is performed on dedicated hardware
- specialized hardware for:
  - out-of-bound address resolution
  - floating-point interpolation
  - type conversion or bit operations

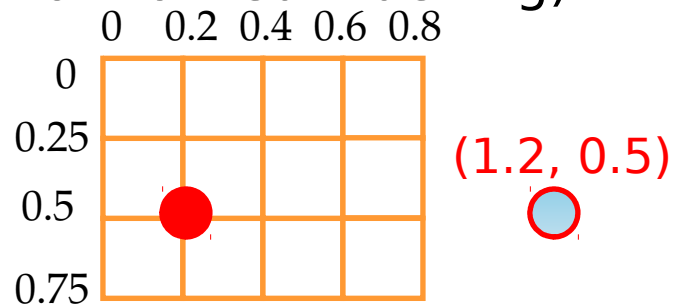


# Texture Memory Addressing Features

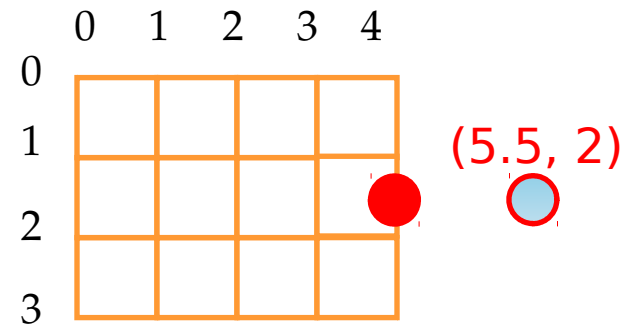
- integer 1D:  $[0, N-1]$
- normalized 1D:  $[0, 1-1/N]$
- available interpolations:
  - floor, linear, bilinear
  - weights are 9 bit



**Wrap:** out-of-border coordinates are replaced in the box using modulus (available only for normalized indexing)



**Clamp:** out-of-border coordinates are clamped to nearest box bound



# Steps for Accessing *Texture Memory*

## CPU

- Allocate global memory on the device (standard, pitched or as cudaArray)  
`cudaMalloc(&d_a, memsize);`
- Create a “texture reference” object at file scope:  
`texture<datatype, dim> d_a_texRef;`  
`datatype` cannot be a double; `dim` can be 1, 2 or 3
- Create a “channel descriptor” object to describe the return type of texture memory load:  
`cudaChannelFormatDesc d_a_desc = cudaCreateChannelDesc<datatype>();`
- Bind the texture reference to memory  
`cudaBindTexture(0, d_a_texRef, d_a, d_a_desc);`
- when finished: unbind the texture reference (there is a maximum number of usable textures):  
`cudaUnbindTexture(d_a_texRef);`

## GPU

- Access data from CUDA kernels through “texture reference”:
  - `tex1Dfetch(d_a_texRef, indirizzo)` - for linear memory
  - `tex1D()`, `tex2D()`, `tex3D()` - for pitched linear texture and cudaArray



# Texture Usage Example

```
__global__ void shiftCopy(int N, int shift, float *odata, float *idata)
{
    int xid = blockIdx.x * blockDim.x + threadIdx.x;
    odata[xid] = idata[xid+shift];
}

texture<float, 1> texRef; // TEXTURE creation

__global__ void textureShiftCopy(int N, int shift, float *odata)
{
    int xid = blockIdx.x * blockDim.x + threadIdx.x;
    odata[xid] = tex1Dfetch(texRef, xid+shift); // TEXTURE FETCHING
}

...

ShiftCopy<<<nBlocks, NUM_THREADS>>>>(N, shift, d_out, d_inp);

cudaChannelFormatDesc d_a_desc = cudaCreateChannelDesc<float>(); // CREATE DESC
cudaBindTexture(0, texRef, d_a, d_a_desc); // BIND TEXTURE MEMORY
textureShiftCopy<<<nBlocks, NUM_THREADS>>>>(N, shift, d_out);
```

# Texture Memory in Kepler: aka **Read-only Cache**

- Starting from Kepler architecture (cc 3.5) constant memory loads from global memory can pass thorough the *texture cache* :
  - without using a explicit texture *binding*
  - without limits on the maximum allowed number of texture

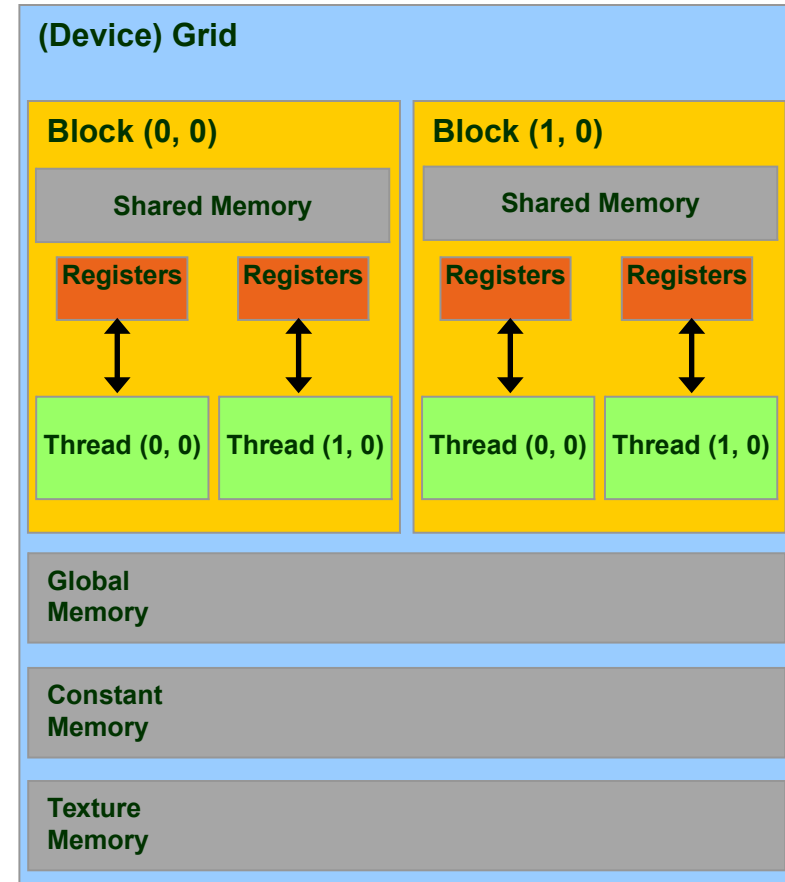
```
__global__ void kernel_copy (float *odata, float *idata) {  
    int index = blockIdx.x * blockDim.x + threadIdx.x;  
    odata[index] = __ldg(idata[index]);  
}
```

```
__global__ void kernel_copy (float *odata, const __restrict__  
float *idata) {  
    int index = blockIdx.x * blockDim.x + threadIdx.x;  
    odata[index] = idata[index];  
}
```

# Registers

- **Registers** are used to store scalars or small array variables with frequent access by each thread
  - **Fermi**: 63 registers per thread / 32 KB
  - **Kepler**: 255 registers per thread / 64 KB
  - **Pascal**: same as Kepler
- **WARNING:**
  - Less registers a kernel needs, more blocks can be assigned to a SM
  - Attention to *Register Pressure*: can be a limiting factor
  - Number of registers per kernel can be limited during *compile time*:  
`--maxregcount max_registers`
  - Number of active blocks per kernel can be forced using the CUDA special qualifier  
`__launch_bounds__`

```
__global__ void  
__launch_bounds__(maxThreadsPerBlock,  
                  minBlocksPerMultiprocessor)  
my_kernel( ... ) { ... }
```



# Local Memory

- **Local Memory** does not correspond to a real physical memory place
- Automatic variables are often place in local memory by the compiler:
  - large structures or arrays that would consume too much register space
- If a kernel uses more registers than available (register spilling), can move variables into local memory
- Local memory is often mapped to global memory
  - using same *Caching* hierachies (L1 for read-only variables)
  - facing same latency and bandwidth limitation of global memory
- In order to obtain information on how much local, constant, shared memory and registers are required for each kernel, you can provide the following compiler options

**--ptxas-options=-v**

```
$ nvcc -arch=sm_20 -ptxas-options=-v my_kernel.cu
...
ptxas info : Used 34 registers, 60+56 bytes lmem, 44+40 bytes
smem, 20 bytes cmem[1], 12 bytes cmem[14]
...
```

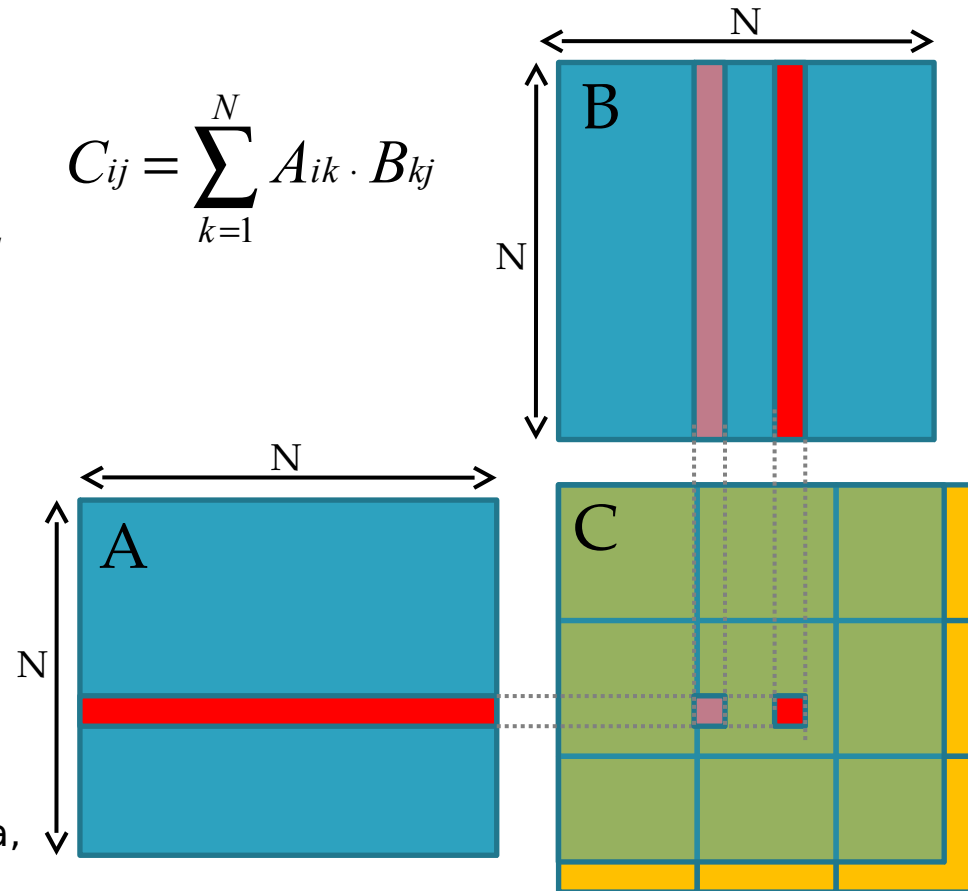
## ■ Matrix-Matrix Product

- limits of global memory implementation
- using shared memory
- implementation guidelines



# Matrix Product using Global Memory

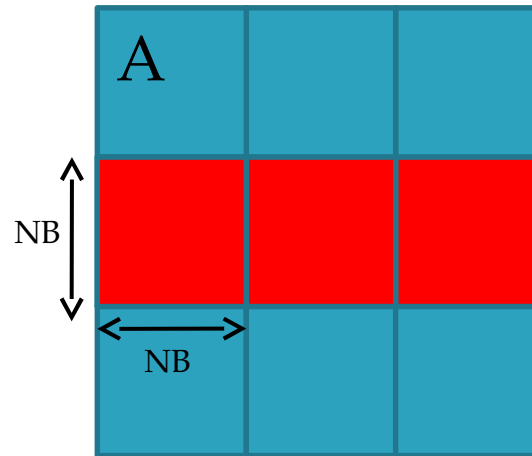
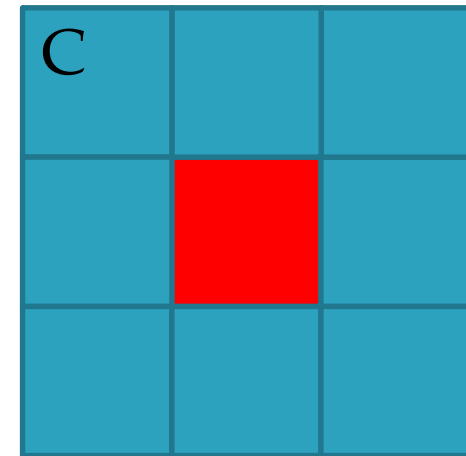
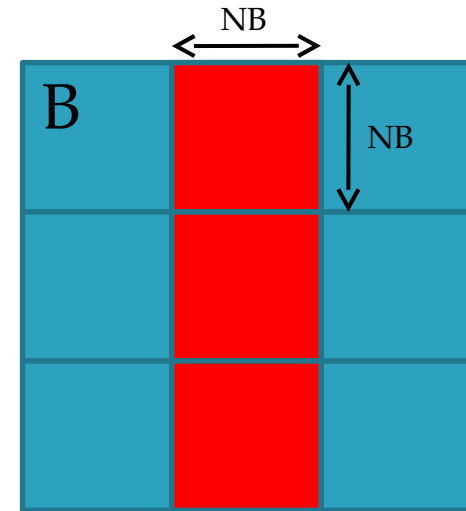
- Each thread compute one element of C, using 2N elements (N from A, N from B) and performing 2N floating-point operations (N add , N mul)
- NB: every element of C shares same row or column retrived N times the same elements from A or B
- This implementation results in  $2N^3$  loads !!!
- We can avoid requesting the same elements many times, sharing them through the shared memory
  - each thread can retrieve just one data element data in parallel and store it into shared memory
  - when all threads have loaded needed data, they can access all the elements by the threads belonging to the same block, for example sharing a full row or column
- Unfortunately shared memory size is small
  - 16/48 KB depending on the compute capability



# Matrix Product using Shared Memory

- Let's solve the problem using blocks of (NB,NB) dimension
  - each CUDA thread block computes the elements of a single matrix block of size (NB·NB) of matrix C
  - each resulting matrix block of matrix C is obtained as the product of all sub-matrices of A and B

$$C_{ij} = \sum_{S=1}^{N/NB} \sum_{k=1}^{NB} A_{Sik} \cdot B_{Skj}$$



- The kernel is divided in two phases:
- threads load a block of A and B from global memory to shared memory
  - threads compute the element of sub-block C reading from shared memory
- Elements of each sub-block C are accumulated using local variables in registers, then stored in global memory
  - Threads synchronizations are required
    - after the load of sub-block of matrix A and B, in order to grant all data is available for sub-block matrix product
    - after the partial sub-block matrix product, in order to grant that next load of other sub-block will not overwrites elements not yet used in current block evaluation

# Matrix Product using Shared Memory: Flow

$C_{ij}=0.$

Cycle on block  
 $kb=0, N/NB$

$As(it,jt) = A(ib*NB + it, kb*NB + jt)$   
 $Bs(it,jt) = B(kb*NB + it, jb*NB + jt)$

Thread Synchronization

Cycle on block:  $k=1,NB$

$C_{ij}=C_{ij}+As(it,k) \cdot Bs(k,jt)$

Thread Synchronization

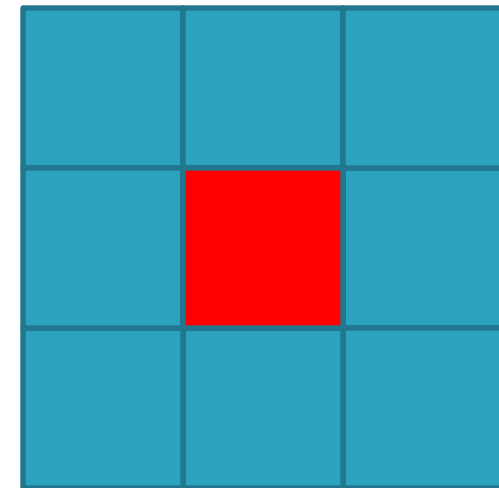
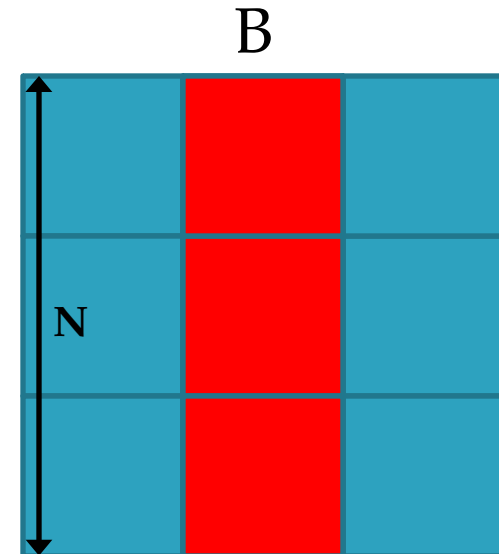
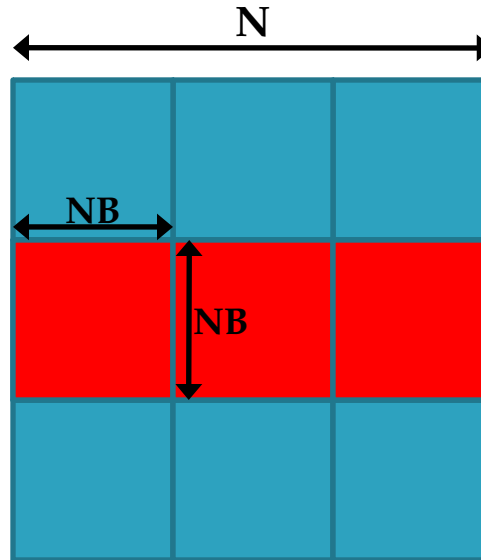
$C(i,j)=C_{ij}$

```
it = threadIdx.y
jt = threadIdx.x

ib = blockIdx.y
jb = blockIdx.x
```

```
it = threadIdx%x
jt = threadIdx%y

ib = blockIdx%x - 1
jb = blockIdx%y - 1
```



A

C



# Matrix Product using Shared Memory: Kernel

```
// Matrix multiplication kernel called by MatMul_gpu()
__global__ void MatMul_kernel (float *A, float *B, float *C, int N)
{

    // Shared memory used to store Asub and Bsub respectively
    __shared__ float Asub[NB][NB];
    __shared__ float Bsub[NB][NB];

    // Block row and column
    int ib = blockIdx.y;
    int jb = blockIdx.x;

    // Thread row and column within Csub
    int it = threadIdx.y;
    int jt = threadIdx.x;

    int a_offset , b_offset, c_offset;

    // Each thread computes one element of Csub
    // by accumulating results into Cvalue
    float Cvalue = 0;

    // Loop over all the sub-matrices of A and B that are
    // required to compute Csub
    // Multiply each pair of sub-matrices together
    // and accumulate the results
```

```
    for (int kb = 0; kb < (A.width / NB); ++kb) {

        // Get the starting address of Asub and Bsub
        a_offset = get_offset (ib, kb, N);
        b_offset = get_offset (kb, jb, N);

        // Load Asub and Bsub from device memory to shared memory
        // Each thread loads one element of each sub-matrix
        Asub[it][jt] = A[a_offset + it*N + jt];
        Bsub[it][jt] = B[b_offset + it*N + jt];

        // Synchronize to make sure the sub-matrices are loaded
        // before starting the computation
        __syncthreads();

        // Multiply Asub and Bsub together
        for (int k = 0; k < NB; ++k) {
            Cvalue += Asub[it][k] * Bsub[k][jt];
        }
        // Synchronize to make sure that the preceding
        // computation is done
        __syncthreads();
    }

    // Get the starting address (c_offset) of Csub
    c_offset = get_offset (ib, jb, N);
    // Each thread block computes one sub-matrix Csub of C
    C[c_offset + it*N + jt] = Cvalue;

}
```

# Rights & Credits

These slides are CINECA 2014 and are released under the Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) Creative Commons license, version 3.0.

Uses not allowed by the above license need explicit, written permission from the copyright owner. For more information see:

<http://creativecommons.org/licenses/by-nc-nd/3.0/>

Slides and examples were authored by:

Isabella Baccarelli, Luca Ferraro, Sergio Orlandini