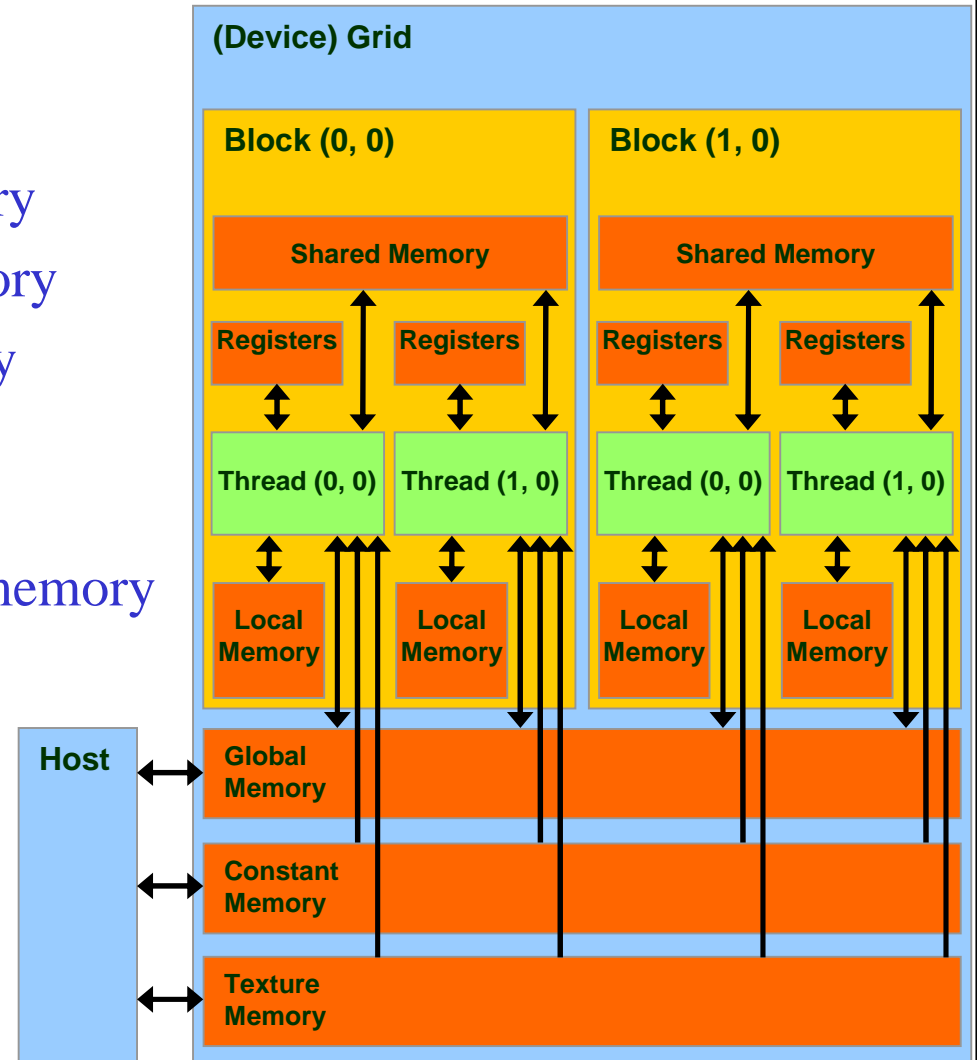# ECE 498AL

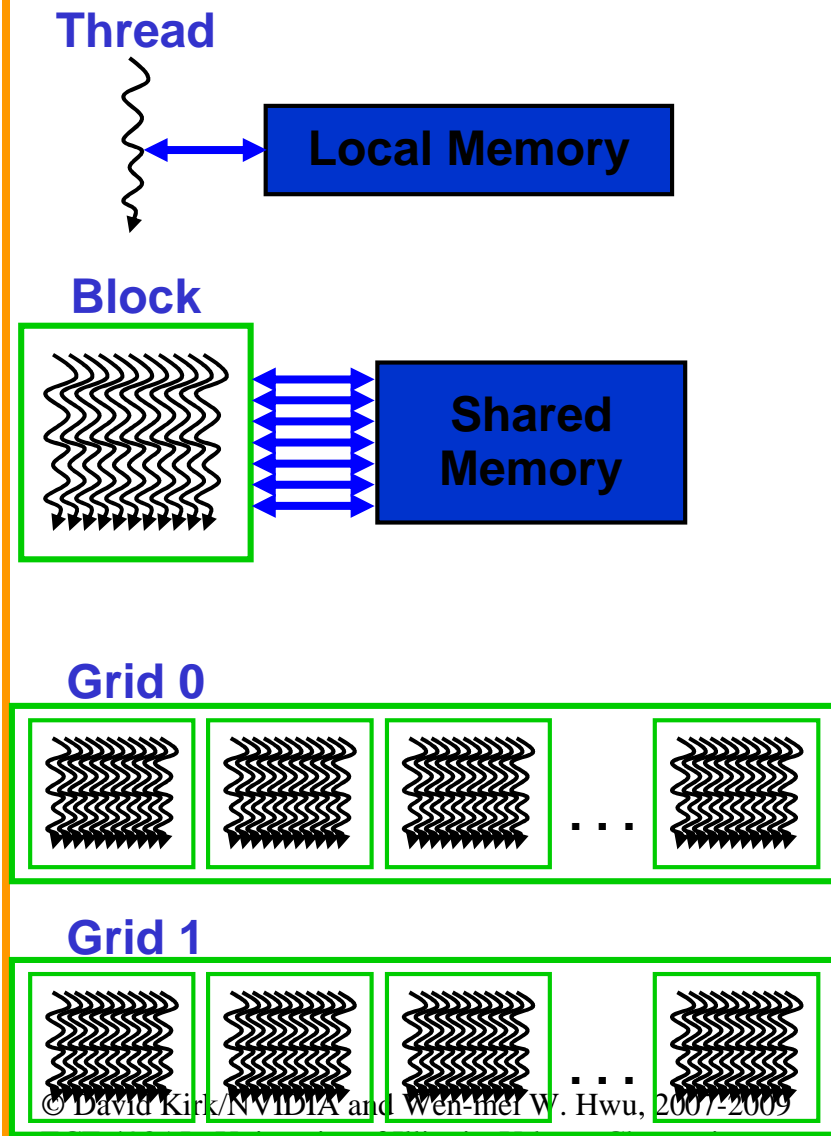# Lectures 9:
# Memory Hardware in G80

# CUDA Device Memory Space: Review

- Each thread can:
    – R/W per-thread registers
    – R/W per-thread local memory
    – R/W per-block shared memory
    – R/W per-grid global memory
    – Read only per-grid constant memory
    – Read only per-grid texture memory

- The host can R/W global, constant, and texture memories



**(Device) Grid**

**Block (0, 0)**
- Shared Memory
- Registers
- Registers
- Thread (0, 0)
- Thread (1, 0)
- Local Memory
- Local Memory

**Block (1, 0)**
- Shared Memory
- Registers
- Registers
- Thread (0, 0)
- Thread (1, 0)
- Local Memory
- Local Memory

**Host**

- Global Memory
- Constant Memory
- Texture Memory

# Parallel Memory Sharing

**Thread**

**Local Memory**

**Block**

**Shared Memory**

**Grid 0**

· · ·

**Grid 1**

· · ·

**Global Memory**

**Sequential Grids in Time**
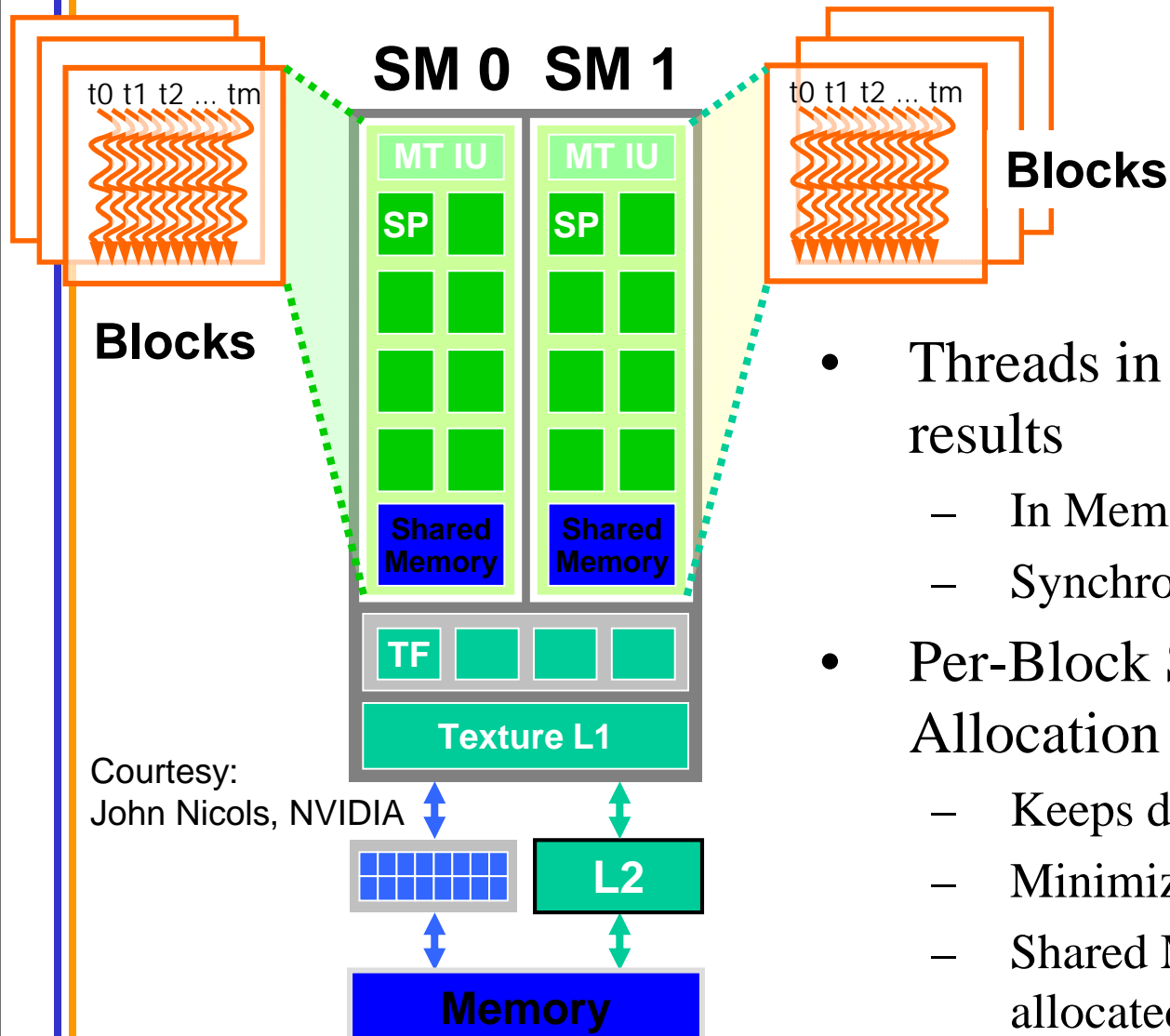
- Local Memory:    per-thread
  - Private per thread
  - Auto variables, register spill
- Shared Memory:  per-Block
  - Shared by threads of the same block
  - Inter-thread communication
- Global Memory:  per-application
  - Shared by all threads
  - Inter-Grid communication

3

# SM Memory Architecture

**SM 0  SM 1**

t0 t1 t2 … tm

t0 t1 t2 … tm

**Blocks**

**Blocks**

MT IU

MT IU

SP

SP

**Shared Memory**

**Shared Memory**

TF

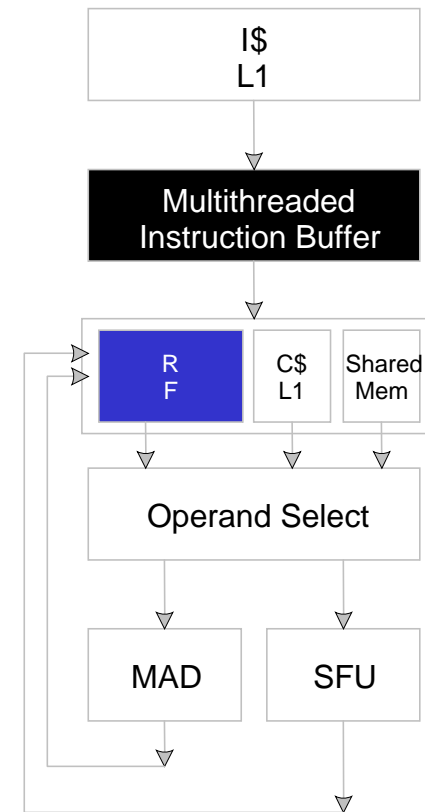**Texture L1**

**L2**

**Memory**

Courtesy:
John Nicols, NVIDIA

- Threads in a block share data & results
  - In Memory and Shared Memory
  - Synchronize at barrier instruction
- Per-Block Shared Memory Allocation
  - Keeps data close to processor
  - Minimize trips to global Memory
  - Shared Memory is dynamically allocated to blocks, one of the limiting resources
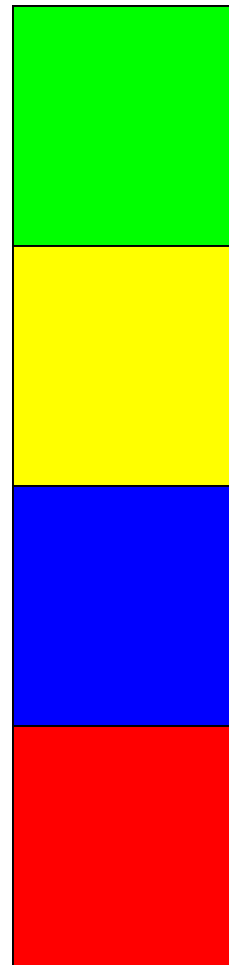
4

# SM Register File

- Register File (RF)
  - 32 KB (8K entries) for each SM in G80
- TEX pipe can also read/write RF
  - 2 SMs share 1 TEX
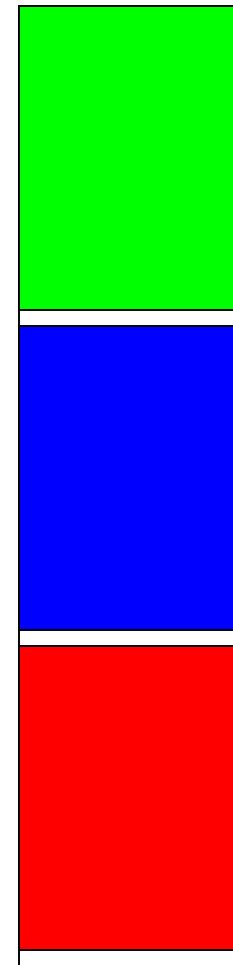- Load/Store pipe can also read/write RF

```
        ┌──────────────┐
        │      I$      │
        │      L1      │
        └──────────────┘
                │
        ┌──────────────┐
        │ Multithreaded│
        │ Instruction  │
        │   Buffer     │
        └──────────────┘
                │
   ┌─────┬──────┬────────┐
   │ R   │ C$   │ Shared │
   │ F   │ L1   │ Mem    │
   └─────┴──────┴────────┘
        │     │      │
   ┌──────────────────────┐
   │   Operand Select     │
   └──────────────────────┘
        │           │
   ┌────────┐   ┌────────┐
   │  MAD   │   │  SFU   │
   └────────┘   └────────┘
```

# Programmer View of Register File

4 blocks          3 blocks

- There are 8192 registers in each SM in G80
  - 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

# Matrix Multiplication Example

- If each Block has 16X16 threads and each thread uses 10 registers, how many thread can run on each SM?
  - Each block requires 10*256 = 2560 registers
  - 8192 = **3** * 2560 + change
  - So, three blocks can run on an SM as far as registers are concerned

- How about if each thread increases the use of registers by 1?
  - Each Block now requires 11*256 = 2816 registers
  - 8192 < 2816 *3
  - Only two Blocks can run on an SM, **1/3 reduction of parallelism**!!!

# More on Dynamic Partitioning

- Dynamic partitioning gives more flexibility to compilers/programmers
  - One can run a smaller number of threads that require many registers each or a large number of threads that require few registers each
    - This allows for finer grain threading than traditional CPU threading models.
  - The compiler can tradeoff between instruction-level parallelism and thread level parallelism

# ILP vs. TLP Example

- Assume that a kernel has 256-thread Blocks, 4 independent instructions for each global memory load in the thread program, and each thread uses 10 registers, global laods have 200 cycles
  - 3 Blocks can run on each SM

- If a compiler can use one more register to change the dependence pattern so that 8 independent instructions exist for each global memory load
  - Only two can run on each SM
  - However, one only needs 200/(8*4) = 7 Warps to tolerate the memory latency
  - Two blocks have 16 Warps. The performance can be actually higher!

# Memory Layout of a Matrix in C

| | | | |
|---|---|---|---|
| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ |
| $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ |
| $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ |
| $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

M

| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Memory Coalescing

- When accessing global memory, peak performance utilization occurs when all threads in a half warp access continuous memory locations.

Not coalesced          coalesced

Md

Thread 1
Thread 2

WIDTH

Nd

WIDTH

# Memory Layout of a Matrix in C



Access direction in Kernel code

$M_{0,0}$ $M_{1,0}$ $M_{2,0}$ $M_{3,0}$
$M_{0,1}$ $M_{1,1}$ $M_{2,1}$ $M_{3,1}$
$M_{0,2}$ $M_{1,2}$ $M_{2,2}$ $M_{3,2}$
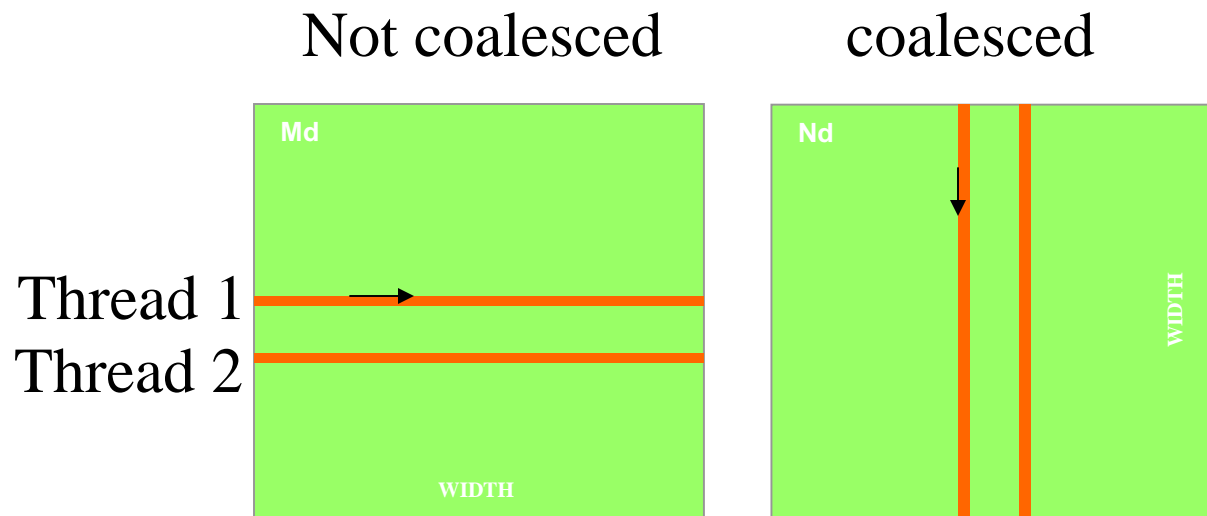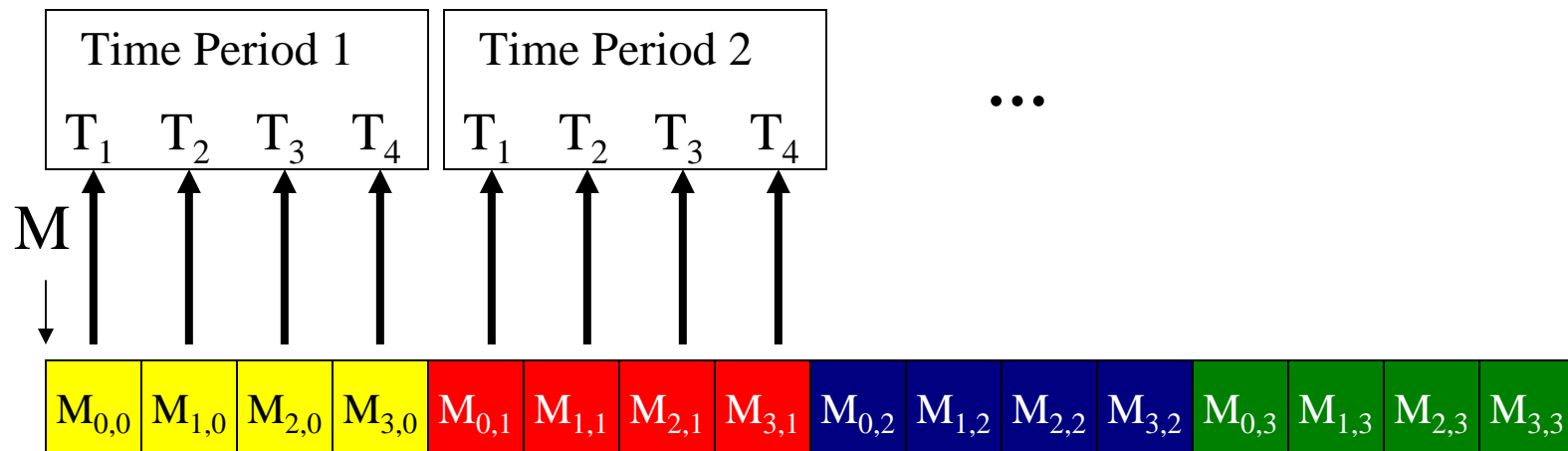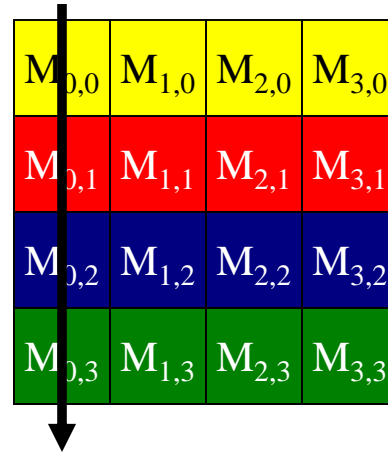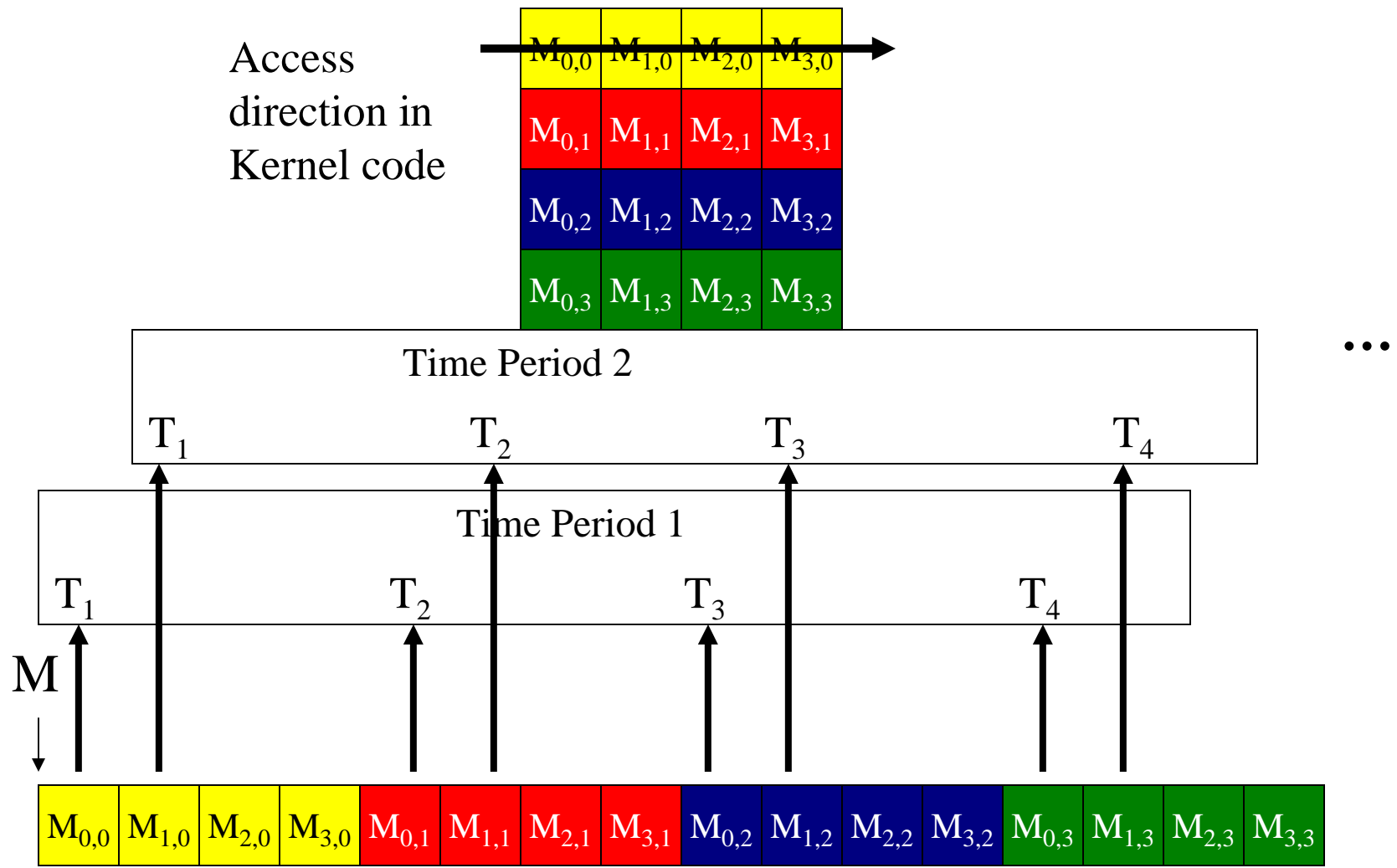$M_{0,3}$ $M_{1,3}$ $M_{2,3}$ $M_{3,3}$

| Time Period 1 | | | | Time Period 2 | | | | ... |
|---|---|---|---|---|---|---|---|---|
| $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_1$ | $T_2$ | $T_3$ | $T_4$ | |

M

$M_{0,0}$ $M_{1,0}$ $M_{2,0}$ $M_{3,0}$ $M_{0,1}$ $M_{1,1}$ $M_{2,1}$ $M_{3,1}$ $M_{0,2}$ $M_{1,2}$ $M_{2,2}$ $M_{3,2}$ $M_{0,3}$ $M_{1,3}$ $M_{2,3}$ $M_{3,3}$

# Memory Layout of a Matrix in C

Access direction in Kernel code

| | | | |
|---|---|---|---|
| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ |
| $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ |
| $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ |
| $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

Time Period 2

$T_1$     $T_2$     $T_3$     $T_4$

...

Time Period 1

$T_1$     $T_2$     $T_3$     $T_4$

M

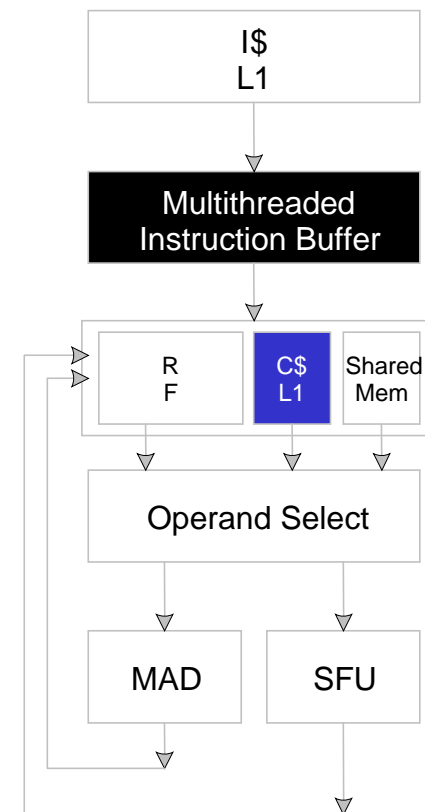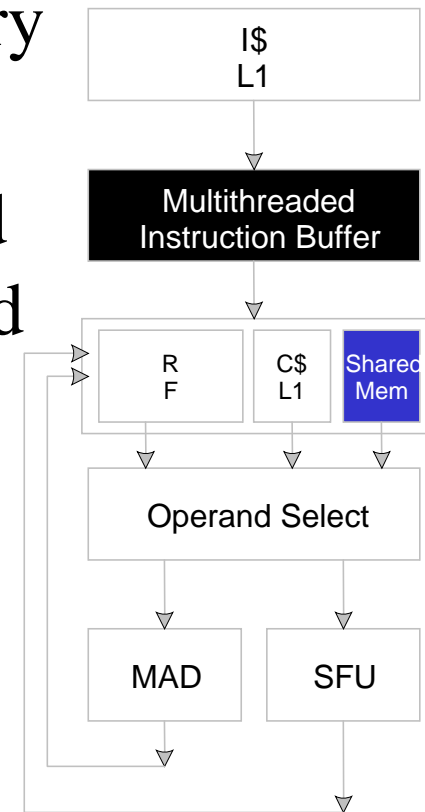| $M_{0,0}$ | $M_{1,0}$ | $M_{2,0}$ | $M_{3,0}$ | $M_{0,1}$ | $M_{1,1}$ | $M_{2,1}$ | $M_{3,1}$ | $M_{0,2}$ | $M_{1,2}$ | $M_{2,2}$ | $M_{3,2}$ | $M_{0,3}$ | $M_{1,3}$ | $M_{2,3}$ | $M_{3,3}$ |

# Constants

- Immediate address constants

- Indexed address constants

- Constants stored in DRAM, and cached on chip
    - L1 per SM

- A constant value can be broadcast to all threads in a Warp
    - Extremely efficient way of accessing a value that is common for all threads in a block!
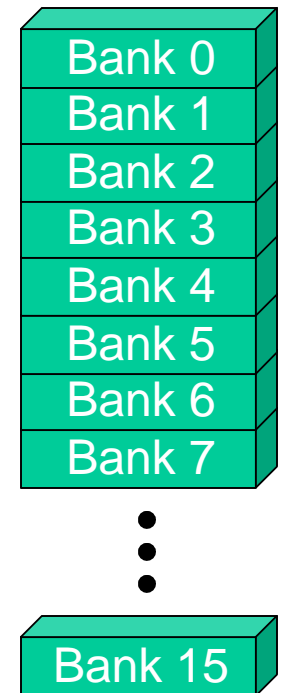
| I$ L1 |
| --- |

| Multithreaded Instruction Buffer |
| --- |

| R F | C$ L1 | Shared Mem |
| --- | --- | --- |

| Operand Select |
| --- |

| MAD | SFU |
| --- | --- |

# Shared Memory

- ### Each SM has 16 KB of Shared Memory
  - 16 banks of 32bit words
- ### CUDA uses Shared Memory as shared storage visible to all threads in a thread block
  - read and write access
- ### Not used explicitly for pixel shader programs
  - we dislike pixels talking to each other ☺

| I$ L1 |
|---|
| **Multithreaded Instruction Buffer** |

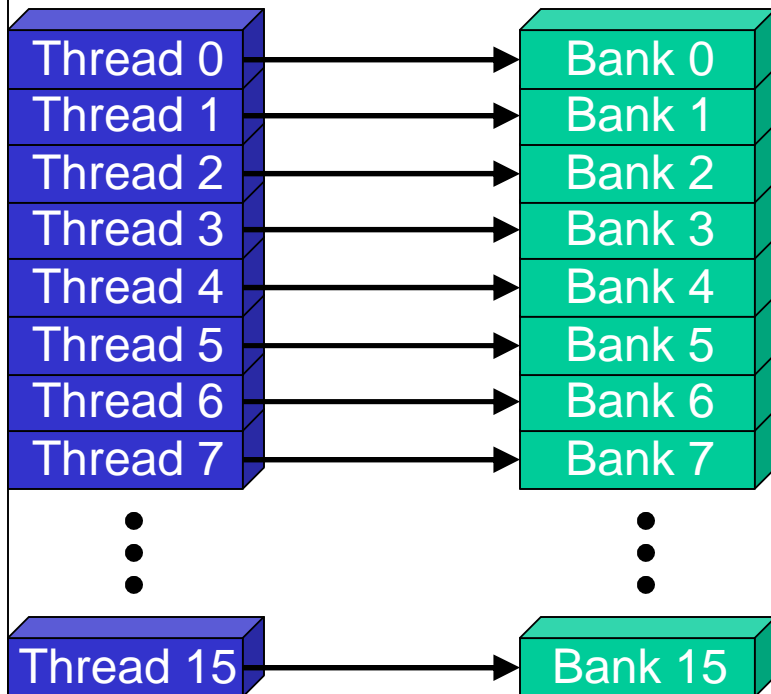| R F | C$ L1 | Shared Mem |
|---|---|---|

| Operand Select |
|---|

| MAD | SFU |
|---|---|

# Parallel Memory Architecture

- In a parallel machine, many threads access memory
    - Therefore, memory is divided into banks
    - Essential to achieve high bandwidth

- Each bank can service one address per cycle
    - A memory can service as many simultaneous accesses as it has banks

- Multiple simultaneous accesses to a bank result in a bank conflict
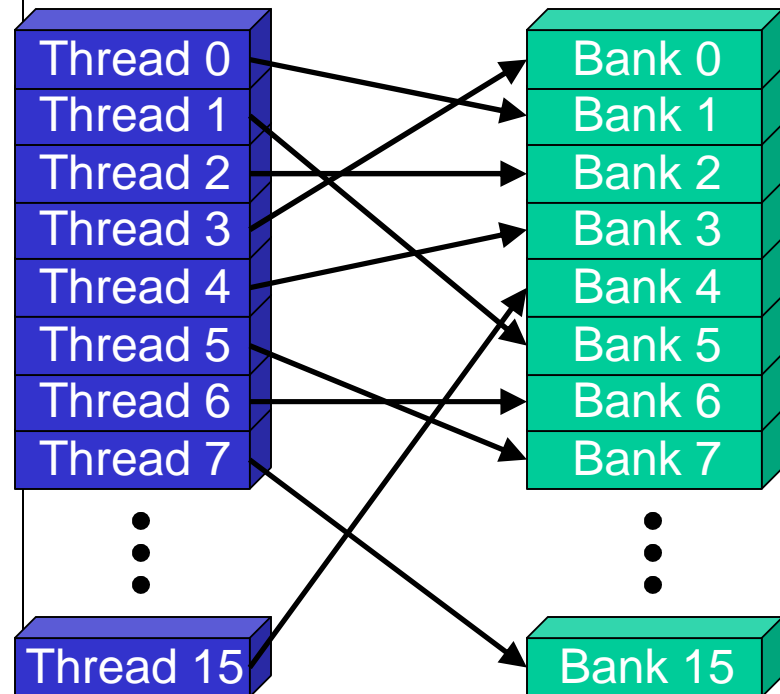    - Conflicting accesses are serialized

| Bank 0 |
| Bank 1 |
| Bank 2 |
| Bank 3 |
| Bank 4 |
| Bank 5 |
| Bank 6 |
| Bank 7 |

...

| Bank 15 |

# Bank Addressing Examples

- **No Bank Conflicts**
  - Linear addressing stride == 1
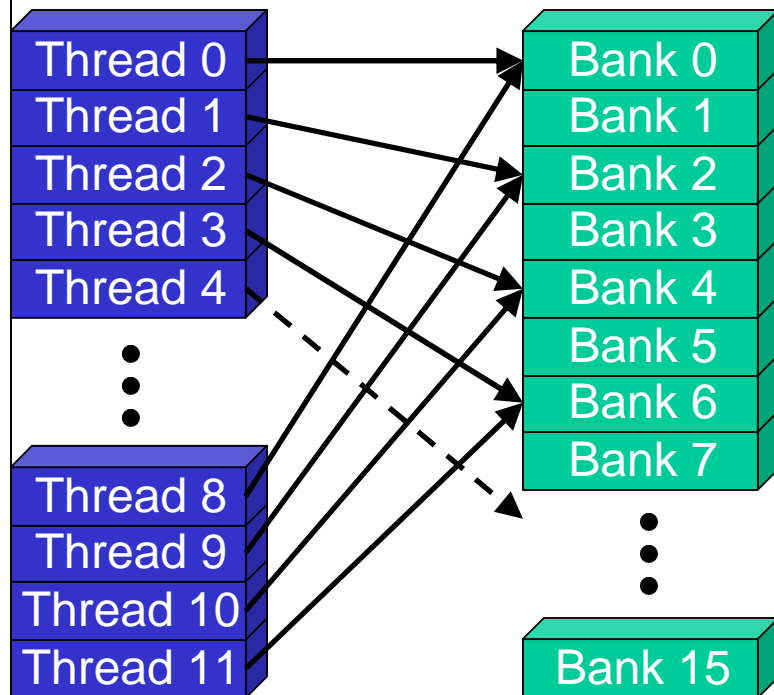
- **No Bank Conflicts**
  - Random 1:1 Permutation
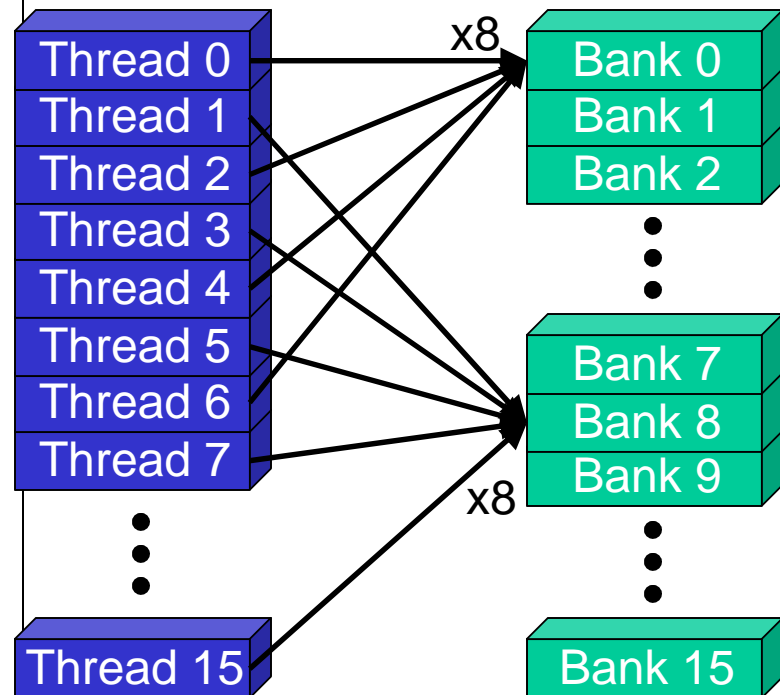
# Bank Addressing Examples

- 2-way Bank Conflicts
  - Linear addressing
    stride == 2

- 8-way Bank Conflicts
  - Linear addressing
    stride == 8

# How addresses map to banks on G80

- Each bank has a bandwidth of 32 bits per clock cycle
- Successive 32-bit words are assigned to successive banks
- G80 has 16 banks
  - So bank = address % 16
  - Same as the size of a half-warp
    - No bank conflicts between different half-warps, only within a single half-warp

# Shared memory bank conflicts

- Shared memory is as fast as registers if there are no bank conflicts

- The fast case:
  - If all threads of a half-warp access different banks, there is no bank conflict
  - If all threads of a half-warp access the identical address, there is no bank conflict (broadcast)

- The slow case:
  - Bank Conflict: multiple threads in the same half-warp access the same bank
  - Must serialize the accesses
  - Cost = max # of simultaneous accesses to a single bank

# Linear Addressing

- Given:

```
__shared__ float shared[256];
float foo =
   shared[baseIndex + s *
   threadIdx.x];
```

- This is only bank-conflict-free if s shares no common factors with the number of banks
  - 16 on G80, so s must be odd

s=1

| Thread 0 | → | Bank 0 |
| Thread 1 | → | Bank 1 |
| Thread 2 | → | Bank 2 |
| Thread 3 | → | Bank 3 |
| Thread 4 | → | Bank 4 |
| Thread 5 | → | Bank 5 |
| Thread 6 | → | Bank 6 |
| Thread 7 | → | Bank 7 |
| Thread 15 | → | Bank 15 |

s=3

| Thread 0 | Bank 0 |
| Thread 1 | Bank 1 |
| Thread 2 | Bank 2 |
| Thread 3 | Bank 3 |
| Thread 4 | Bank 4 |
| Thread 5 | Bank 5 |
| Thread 6 | Bank 6 |
| Thread 7 | Bank 7 |
| Thread 15 | Bank 15 |