# CS 380 - GPU and GPGPU Programming
# Lecture 23: CUDA Memory, Pt. 2

Markus Hadwiger, KAUST

# Reading Assignment #13 (until Nov 27)

Read (required):

- Optimizing Parallel Reduction in CUDA, Mark Harris,

  `https://developer.download.nvidia.com/assets/cuda/files/reduction.pdf`

- Programming Massively Parallel Processors book, 4th edition
  **Chapter 8**: Prefix Sum (Scan) – work efficiency in parallel algorithms


Read (optional):

- Programming Massively Parallel Processors book, 4th edition
  **Chapter 7**: Reduction

- GPU Gems 3 book, Chapter 39: Parallel Prefix Sum (Scan) with CUDA

  `https://developer.nvidia.com/gpugems/GPUGems3/gpugems3_ch39.html`

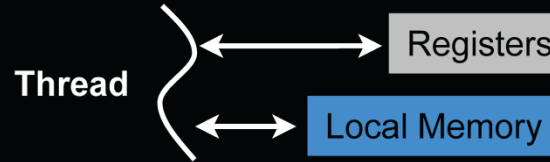- Faster Parallel Reductions on Kepler, Justin Luitjens

  `https://devblogs.nvidia.com/parallelforall/faster-parallel-reductions-kepler/`

# Next Lectures

Lecture 24: Tue, Nov 22 (make-up lecture; 16:00 – 17:15)

Lecture 25: Wed, Nov 23
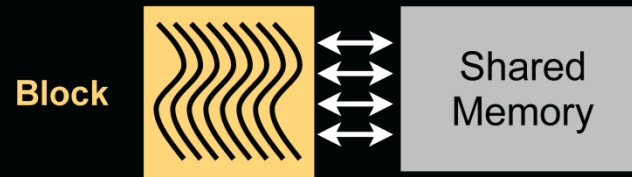
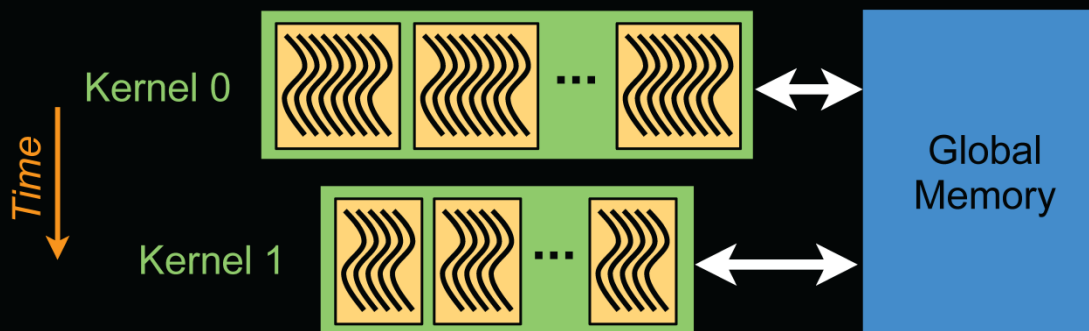# CUDA Memory

# Kernel Memory Access

- **Per-thread**

  Thread
  - Registers — On-chip
  - Local Memory — Off-chip, ~~uncached~~

    cached on Fermi or newer!

- **Per-block**

  Block — Shared Memory
  - On-chip, small
  - Fast

- **Per-device**

  Kernel 0 ... / Kernel 1 ... — Global Memory

  Time
  - Off-chip, large
  - ~~Uncached~~
  - Persistent across kernel launches
  - Kernel I/O

  cached on Fermi or newer!

NVIDIA.

# CUDA Memory:
# Shared Memory

# Memory and Cache Types

Global memory

- [Device]      **L2 cache**

- [SM]            **L1 cache** (shared mem carved out;  *or*  L1 shared with tex cache)

- [SM/TPC]    **Texture cache** (separate, or shared with L1 cache)

- [SM]            **Read-only data cache** (storage might be same as tex cache)

## Shared memory

- [SM]  Shareable only between threads in same thread block
        (Hopper/CC 9.x: also thread block clusters)

Constant memory: Constant (uniform) cache

Unified memory programming: Device/host memory sharing

# L1 Cache vs. Shared Memory

Different configs on Fermi and Kepler; carveout on Maxwell and newer

- More shared memory on newer GPUs (64KB, 96KB, 100KB, 164KB, ...)

  Carveout from unified L1/read-only data cache

  (See CUDA C Programming Guide!)

```
// Device code
__global__ void MyKernel(...)
{
    __shared__ float buffer[BLOCK_DIM];
    ...
}

// Host code
int carveout = 50; // prefer shared memory capacity 50% of maximum
// Named Carveout Values:
// carveout = cudaSharedmemCarveoutDefault;   //  (-1)
// carveout = cudaSharedmemCarveoutMaxL1;     //   (0)
// carveout = cudaSharedmemCarveoutMaxShared; // (100)
cudaFuncSetAttribute(MyKernel, cudaFuncAttributePreferredSharedMemoryCarveout,
 carveout);
MyKernel <<<gridDim, BLOCK_DIM>>>(...);
```

# NVIDIA GH100 SM

Multiprocessor: SM (CC 9.0)

- 128 FP32 + 64 INT32 cores
- 64 FP64 cores
- 4x 4th gen tensor cores
- ++ thread block clusters, DPX insts., FP8, TMA

4 partitions inside SM

- 32 FP32 + 16 INT32 cores
- 16 FP64 cores
- 8x LD/ST units; 4 SFUs each
- 1x 4th gen tensor core each
- Each has: warp scheduler, dispatch unit, 16K register file

# Compute Capab. 9.x (Hopper, Part 2)

## K.8.3. Shared Memory

Similar to the NVIDIA Ampere GPU architecture, the amount of the unified data cache reserved for shared memory is configurable on a per kernel basis. For the *NVIDIA H100 Tensor Core GPU architecture*, the unified data cache has a size of 256 KB for devices of compute capability 9.0. The shared memory capacity can be set to 0, 8, 16, 32, 64, 100, 132, 164, 196 or 228 KB.

As with the NVIDIA Ampere GPU architecture, an application can configure its preferred shared memory capacity, i.e., the `carveout`. Devices of compute capability 9.0 allow a single thread block to address up to 227 KB of shared memory. Kernels relying on shared memory allocations over 48 KB per block are architecture-specific, and must use dynamic shared memory rather than statically sized shared memory arrays. These kernels require an explicit opt-in by using `cudaFuncSetAttribute()` to set the `cudaFuncAttributeMaxDynamicSharedMemorySize`; see Shared Memory for the Volta architecture.

Note that the maximum amount of shared memory per thread block is smaller than the maximum shared memory partition available per SM. The 1 KB of shared memory not made available to a thread block is reserved for system use.

# Shared Memory Allocation

- **2 modes**
- **Static size within kernel**

```
__shared__ float vec[256];
```

- **Dynamic size when calling the kernel**

```
// in main
int VecSize = MAX_THREADS * sizeof(float4);
vecMat<<< blockGrid, threadBlock, VecSize >>>( p1, p2, …);

// declare as extern within kernel
extern __shared__ float vec[];
```
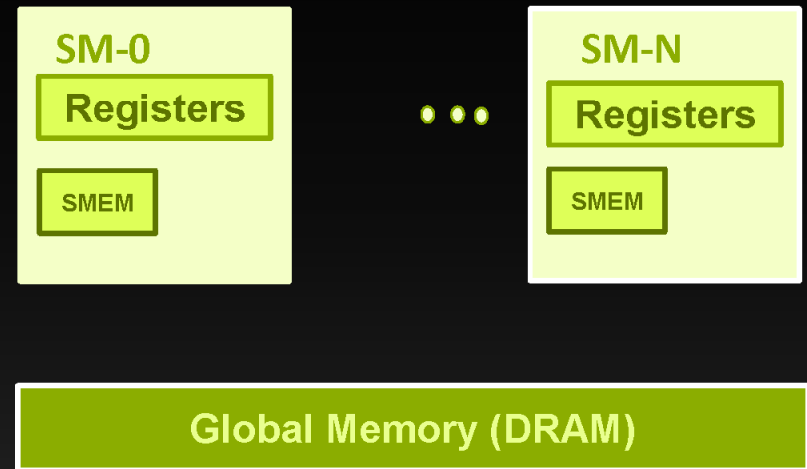
# Shared Memory

- **Accessible by all threads in a block**

- **Fast compared to global memory**
  - **Low access latency**
  - **High bandwidth**

- **Common uses:**
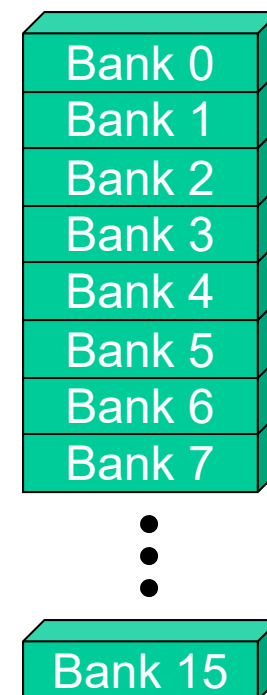  - **Software managed cache**
  - **Data layout conversion**

**SM-0**
Registers
SMEM

• • •

**SM-N**
Registers
SMEM

**Global Memory (DRAM)**

# Shared Memory/L1 Sizing

- **Shared memory and L1 use the same 64KB**
  - **Program-configurable split:**
    - **Fermi: 48:16, 16:48**
    - **Kepler: 48:16, 16:48, 32:32**          use cudaFuncSetAttribute()
  - **CUDA API:** ~~cudaDeviceSetCacheConfig(), cudaFuncSetCacheConfig()~~
- **Large L1 can improve performance when:**
  - **Spilling registers (more lines in the cache -> fewer evictions)**
- **Large SMEM can improve performance when:**
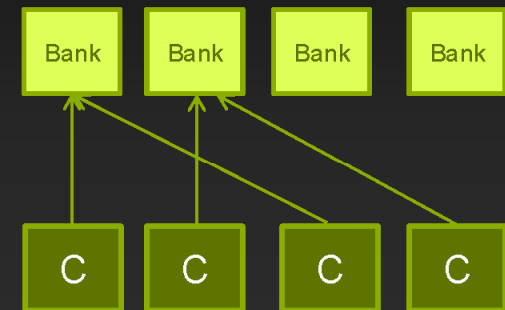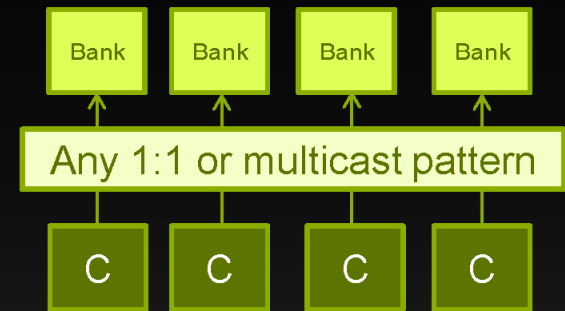  - **Occupancy is limited by SMEM**

# Shared Memory

- ## Uses:
  - Inter-thread communication within a block
  - Cache data to reduce redundant global memory accesses
  - Use it to improve global memory access patterns

- ## Organization:
  - 32 banks, 4-byte (or 8-byte) banks
  - Successive words accessed through different banks

# 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

# Memory Banks

Fermi/Kepler/Maxwell and newer:

32 banks

default:
4B / bank

Kepler or newer:
configurable
to 8B / bank

# Shared Memory

- **Uses:**
  - Inter-thread communication within a block
  - Cache data to reduce redundant global memory accesses
  - Use it to improve global memory access patterns

- **Performance:**
  - smem accesses are issued per warp
  - Throughput is 4 (or 8) bytes per bank per clock per multiprocessor
  - **serialization:** if *N* threads of 32 access different words in the same bank, *N* accesses are executed serially
  - **multicast:** *N* threads access the same word in one fetch
    - Could be different bytes within the same word

# Shared Memory Organization

- **Organized in 32 independent banks**

- **Optimal access: no two words from same bank**
  - **Separate banks per thread**
  - **Banks can multicast**

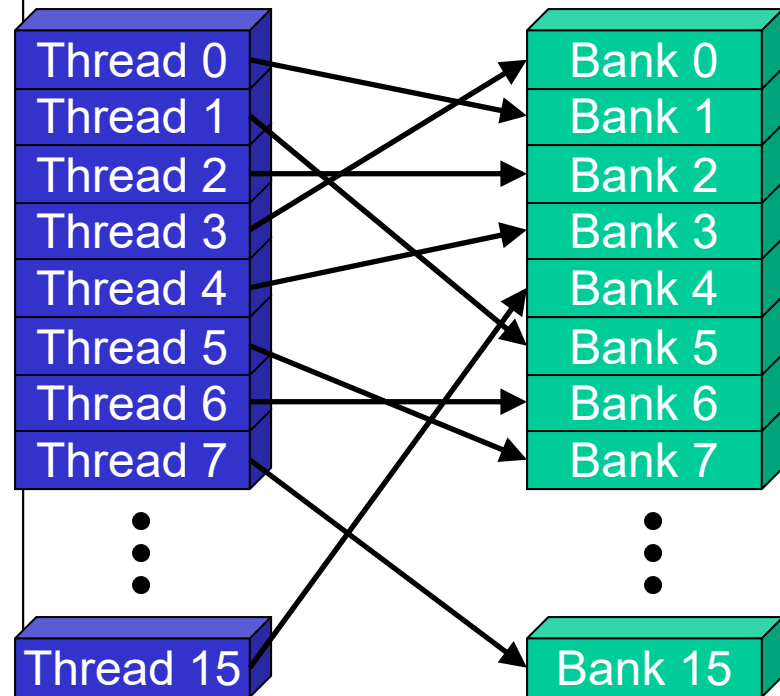- **Multiple words from same bank serialize**

| Bank | Bank | Bank | Bank |

Any 1:1 or multicast pattern

| C | C | C | C |

| Bank | Bank | Bank | Bank |

| C | C | C | C |

# Bank Addressing Examples

- No Bank Conflicts
  - Linear addressing stride == 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 |

- No Bank Conflicts
  - Random 1:1 Permutation

| 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 |

# 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

Fermi and newer have 32 banks,
considers full warps instead of half warps!

# 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

full warps instead of half warps on Fermi and newer!

# 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

# Data Types and Bank Conflicts

- **This has no conflicts if type of `shared` is 32-bits:**

  ```
  foo = shared[baseIndex + threadIdx.x]
  ```

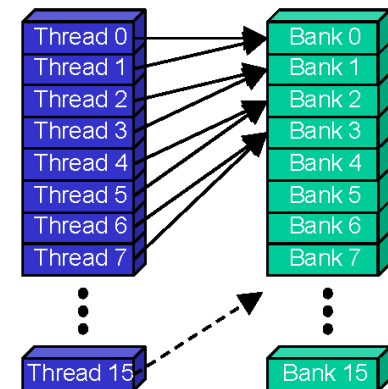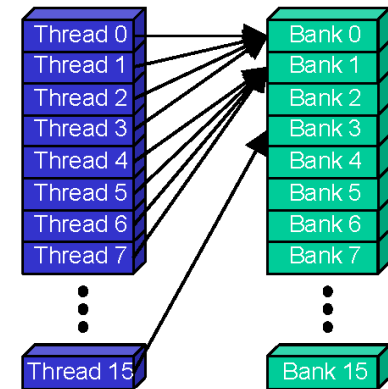- **But not if the data type is smaller**
  - 4-way bank conflicts:
  ```
  __shared__ char shared[];
  foo = shared[baseIndex + threadIdx.x];
  ```

not true on Fermi, because of multi-cast!

  - 2-way bank conflicts:
  ```
  __shared__ short shared[];
  foo = shared[baseIndex + threadIdx.x];
  ```
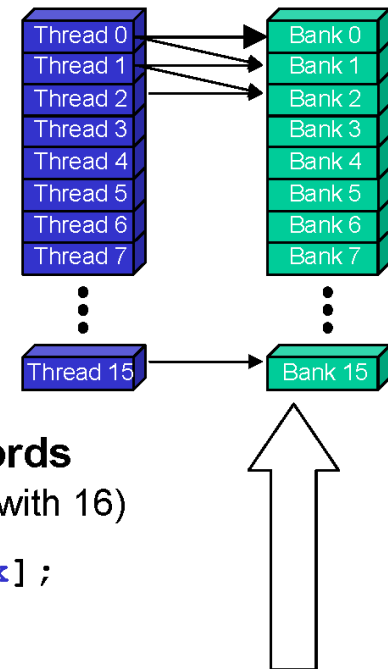
not true on Fermi, because of multi-cast!

# Structs and Bank Conflicts

- **Struct assignments compile into as many memory accesses as there are struct members:**

```
struct vector { float x, y, z; };
struct myType {
    float f;
    int c;
};
__shared__ struct vector vectors[64];
__shared__ struct myType myTypes[64];
```



- **This has no bank conflicts for vector; struct size is 3 words**
  - 3 accesses per thread, contiguous banks (no common factor with 16)

```
struct vector v = vectors[baseIndex + threadIdx.x];
```

- **This has 2-way bank conflicts for myType; (each bank will be accessed by 2 threads simultaneously)**

```
struct myType m = myTypes[baseIndex + threadIdx.x];
```
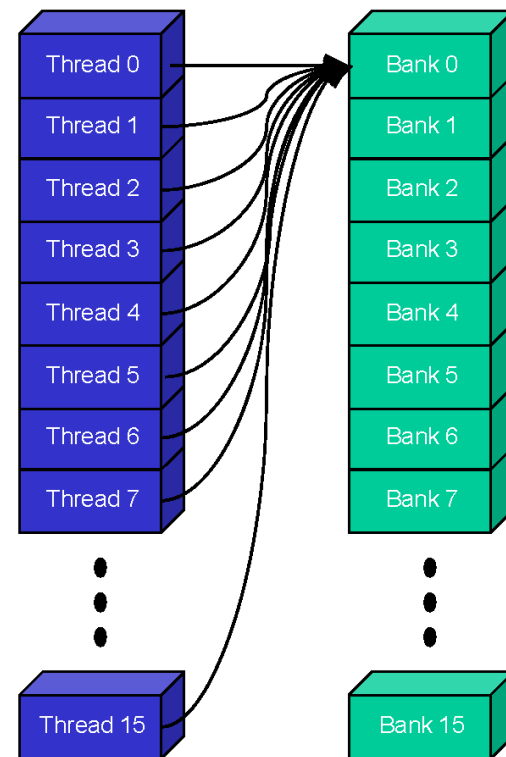
# Broadcast on Shared Memory

- **Each thread loads the same element – no bank conlict**

    `x = shared[0];`

- **Will be resolved implicitly**

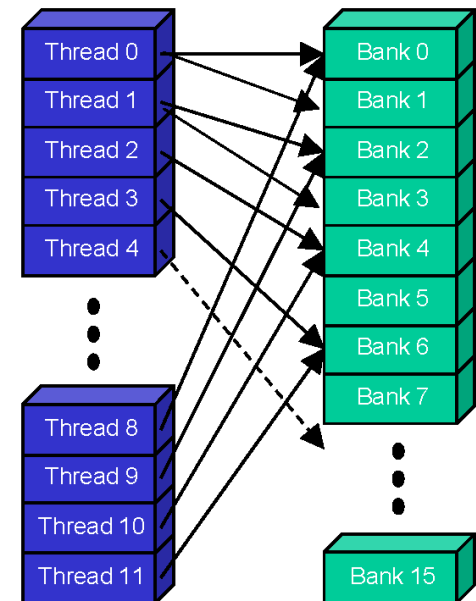multi-cast on Fermi and newer!

# Common Array Bank Conflict Patterns 1D

- **Each thread loads 2 elements into shared mem:**
  - 2-way-interleaved loads result in 2-way bank conflicts:

```
int tid = threadIdx.x;
shared[2*tid] = global[2*tid];
shared[2*tid+1] = global[2*tid+1];
```
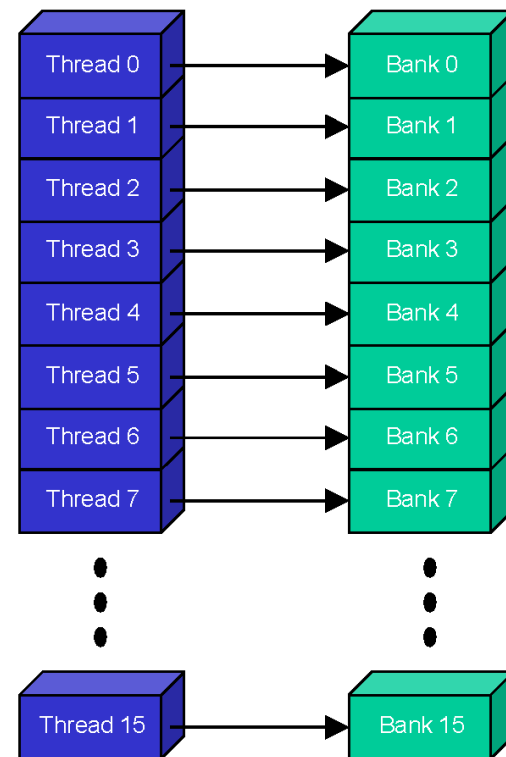
- **This makes sense for traditional CPU threads, locality in cache line usage and reduced sharing traffic.**
  - Not in shared memory usage where there is no cache line effects but banking effects

# A Better Array Access Pattern

- **Each thread loads one element in every consecutive group of `blockDim` elements.**

```
shared[tid] = global[tid];
shared[tid + blockDim.x] =
    global[tid + blockDim.x];
```
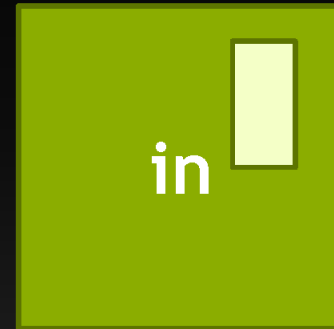
Hendrik Lensch and Robert Strzodka

# OPTIMIZE

**Kernel Optimizations:** *Shared Memory Accesses*

# Case Study: Matrix Transpose

- Coalesced read
- Scattered write (stride N)

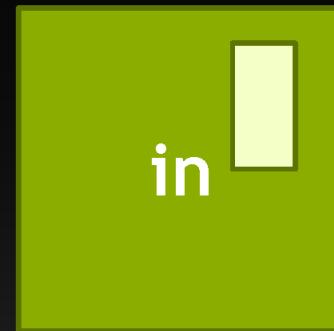⇒ **Process matrix tile, not single row/column, per block**

⇒ **Transpose matrix tile within block**

# Case Study: Matrix Transpose

- Coalesced read
- Scattered write (stride N)

- Transpose matrix tile within block

$\Rightarrow$ **Need threads in a block to cooperate: use shared memory**

in

out

# Transpose with coalesced read/write

```
__global__ transpose(float in[], float out[])
{

    __shared__ float tile[TILE][TILE];

    int glob_in = xIndex + (yIndex)*N;
    int glob_out = xIndex + (yIndex)*N;

    tile[threadIdx.y][threadIdx.x] = in[glob_in];

    __syncthreads();

    out[glob_out] = tile[threadIdx.x][threadIdx.y];

}
```
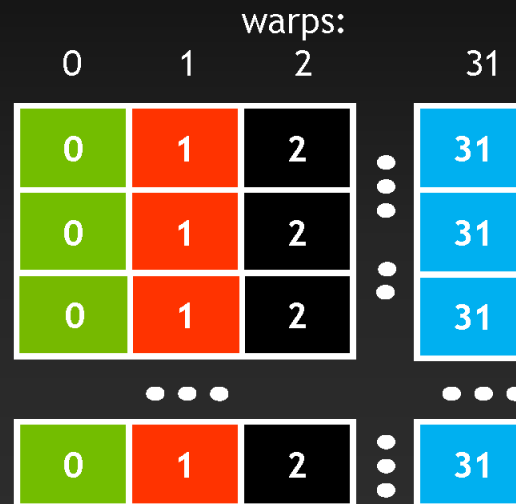
**Fixed GMEM coalescing, but introduced SMEM bank conflicts**

```
transpose<<<grid, threads>>>(in, out);
```

# Shared Memory: Avoiding Bank Conflicts

- **Example: 32x32 SMEM array**
- **Warp accesses a column:**
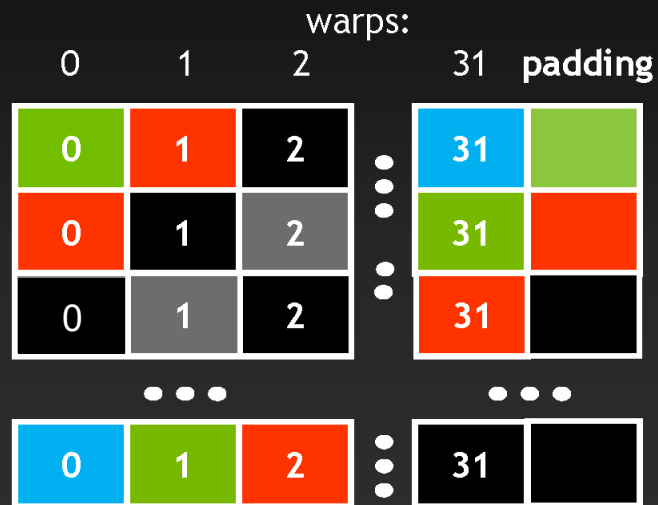  - **32-way bank conflicts (threads in a warp access the same bank)**

# Shared Memory: Avoiding Bank Conflicts

- **Add a column for padding:**
  - **32x33 SMEM array**
- **Warp accesses a column:**
  - **32 different banks, no bank conflicts**



Bank 0
Bank 1
...
Bank 31

Thank you.