# LECTURE4 PERFORMANCE CONSIDERATIONS

Warps and SIMD

performance impact of control divergence

Parallel reduction

Memory parallelism

# Objective

– To understand how CUDA threads execute on SIMD Hardware
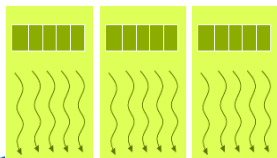  – Warp partitioning
  – SIMD Hardware
  – Control divergence

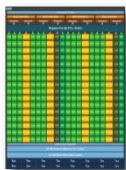电子科技大学
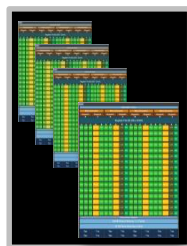University of Electronic Science and Technology of China

# 执行过程

软件

硬件

线程

CUDA Core

**CUDA Core=ALU=SP**
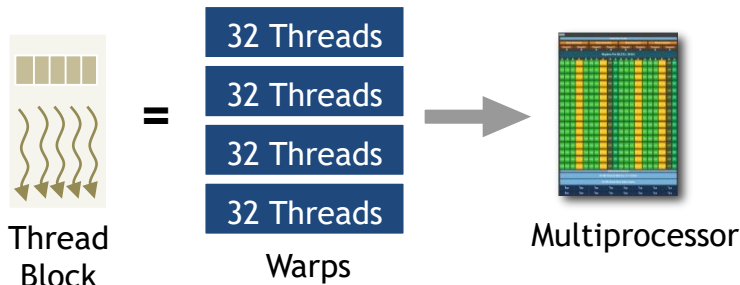
Thread Block

Multiprocessor

**SM=内核=逻辑架构里的CORE**

Device

当调用**kernel**函数时，启动起来很多线程，然后分配给硬件去执行，执行过程中要占用硬件资源。

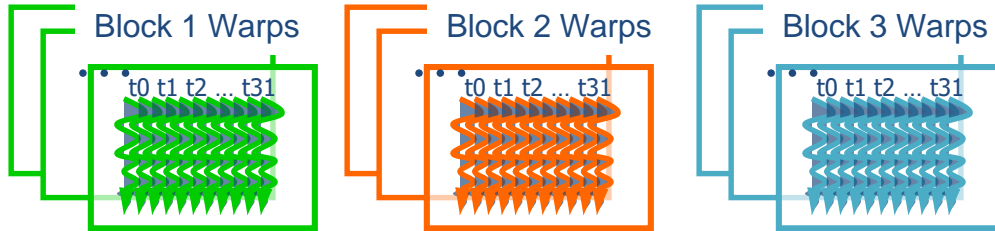# Warps



**A thread block consists of 32-thread *warps***

**A warp is executed physically in parallel (SIMD) on a multiprocessor**

Thread Block  =  Warps  →  Multiprocessor

32 Threads
32 Threads
32 Threads
32 Threads

# Warps as Scheduling Units



Block 1 Warps — t0 t1 t2 ... t31
Block 2 Warps — t0 t1 t2 ... t31
Block 3 Warps — t0 t1 t2 ... t31
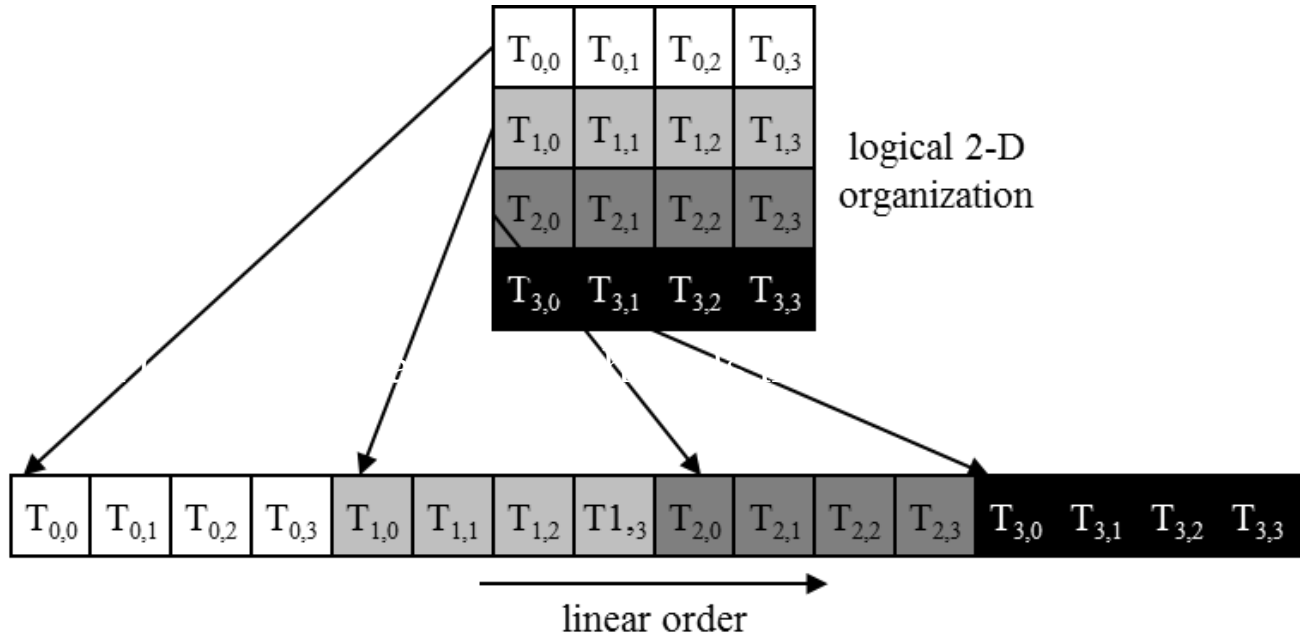
– Each block is divided into 32-thread warps

  – An implementation technique, not part of the CUDA programming model

  – Warps are scheduling units in SM

  – Threads in a warp execute in Single Instruction Multiple Data (SIMD) manner

  – The number of threads in a warp may vary in future generations

# Warps in Multi-dimensional Thread Blocks

- The thread blocks are first linearized into 1D in row major order
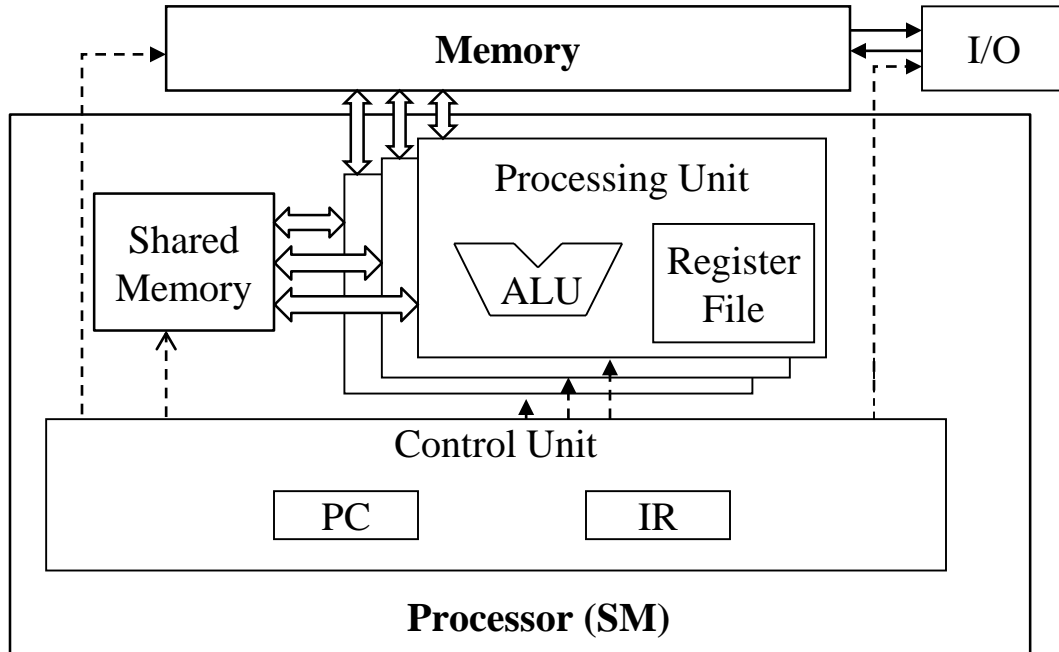  - In x-dimension first, y-dimension next, and z-dimension last



logical 2-D organization

linear order

# Blocks are partitioned after linearization

– Linearized thread blocks are partitioned
  – Thread indices within a warp are consecutive and increasing
  – Warp 0 starts with Thread 0

– Partitioning scheme is consistent across devices
  – Thus you can use this knowledge in control flow
  – However, the exact size of warps may change from generation to generation

– DO NOT rely on any ordering within or between warps
  – If there are any dependencies between threads, you must __syncthreads() to get correct results (more later).
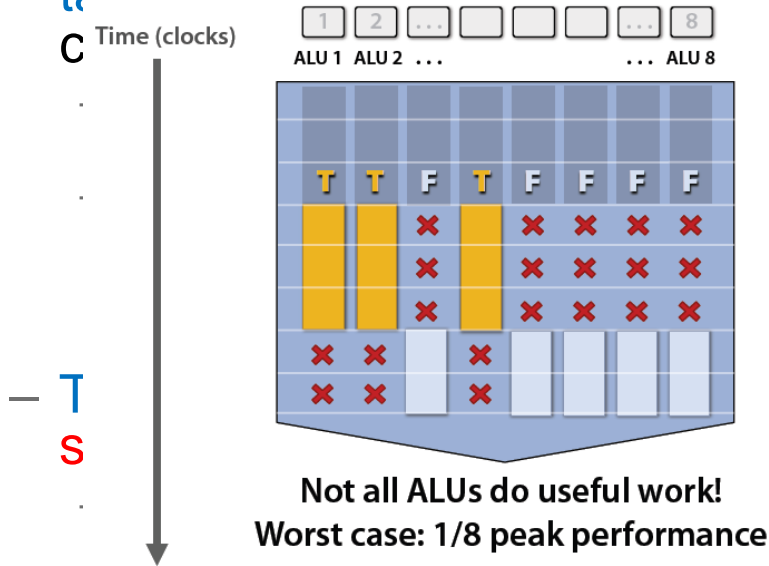
# SMs are SIMD Processors

– Control unit for instruction fetch, decode, and control is shared among multiple processing units

  – Control overhead is minimized (Module 1)

# Control Divergence

– Control divergence occurs when threads in a warp take different control flow paths by making different c...



...ath of an ...

...1s than ...

– T ...
s ...

...hs are ...

...e

**Not all ALUs do useful work!**
**Worst case: 1/8 peak performance**

```
<unconditional
 shader code>

if (x > 0) {
    y = pow(x, exp);
    y *= Ks;
    refl = y + Ka;
} else {
    x = 0;
    refl = Ka;
}
<resume unconditional
 shader code>
```

Time (clocks)
ALU 1  ALU 2 ...        ... ALU 8

– The number of different paths can be large when considering nested control flow statements

# Control Divergence Examples

– Divergence can arise when branch or loop condition is a function of thread indices

– Example kernel statement with divergence:

– if (threadIdx.x > 2) { }

– This creates two different control paths for threads in a block

– Decision granularity < warp size; threads 0, 1 and 2 follow different path than the rest of the threads in the first warp

– Example without divergence:

– If (blockIdx.x > 2) { }

– Decision granularity is a multiple of blocks size; all threads in any given warp follow the same path

# Example: Vector Addition Kernel

## Device Code

```
// Compute vector sum C = A + B
// Each thread performs one pair-wise addition

__global__
void vecAddKernel(float* A, float* B, float* C,
  int n)
{
  int i = threadIdx.x + blockDim.x * blockIdx.x;
   if(i<n) C[i] = A[i] + B[i];
}
```

# Analysis for vector size of 1,000 elements

- Assume that block size is 256 threads
  - 8 warps in each block

- All threads in Blocks 0, 1, and 2 are within valid range
  - i values from 0 to 767
  - There are 24 warps in these three blocks, none will have control divergence

- Most warps in Block 3 will not control divergence
  - Threads in the warps 0-6 are all within valid range, thus no control divergence

- One warp in Block 3 will have control divergence
  - Threads with i values 992-999 will all be within valid range
  - Threads with i values of 1000-1023 will be outside valid range

- Effect of serialization on control divergence will be small
  - 1 out of 32 warps has control divergence
  - The impact on performance will likely be less than 3%

# Warps and SIMD

# performance impact of control divergence

# parallel reduction

# Objective

– To learn to analyze the performance impact of control divergence
   – Boundary condition checking
   – Control divergence is data-dependent

# Performance Impact of Control Divergence

– Boundary condition checks are vital for complete functionality and robustness of parallel code

– The tiled matrix multiplication kernel has many boundary condition checks

– The concern is that these checks may cause significant performance degradation

– For example, see the tile loading code below:

```
if(Row < Width && p * TILE_WIDTH+tx < Width) {
    ds_M[ty][tx] = M[Row * Width + p * TILE_WIDTH + tx];
} else {
  ds_M[ty][tx] = 0.0;
}

if (p*TILE_WIDTH+ty < Width && Col < Width) {
    ds_N[ty][tx] = N[(p*TILE_WIDTH + ty) * Width + Col];
} else {
    ds_N[ty][tx] = 0.0;
}
```
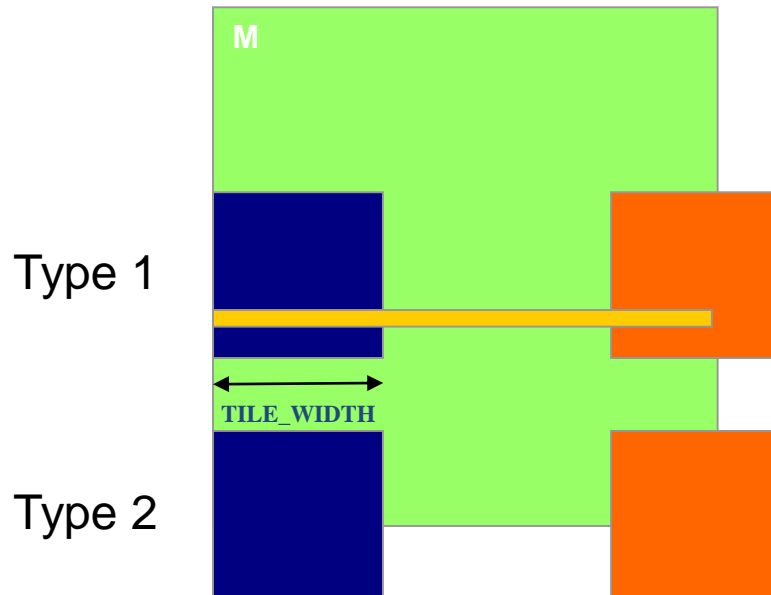
# Two types of blocks in loading M Tiles

- 1. Blocks whose tiles are all within valid range until the last phase.
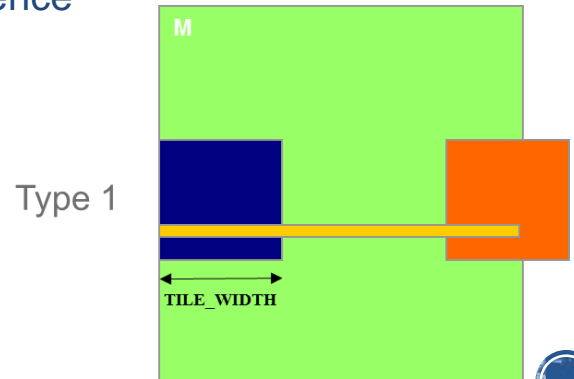- 2. Blocks whose tiles are partially outside the valid range all the way

# Analysis of Control Divergence Impact

- – Assume 16x16 tiles and thread blocks
- – Each thread block has 8 warps (256/32)
- – Assume square matrices of 100x100
- – Each thread will go through 7 phases (ceiling of 100/16)

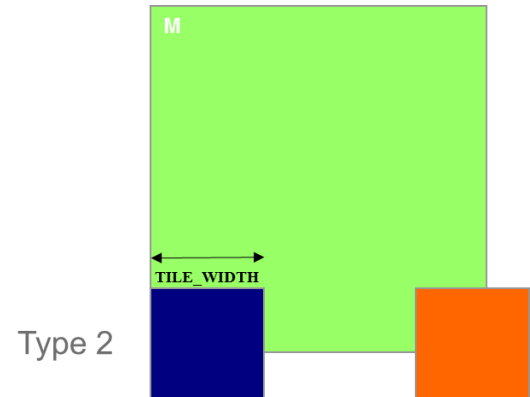- – There are 49 thread blocks (7 in each dimension)

# Control Divergence in Loading M Tiles

- Assume 16x16 tiles and thread blocks
- Each thread block has 8 warps (256/32)
- Assume square matrices of 100x100
- Each warp will go through 7 phases (ceiling of 100/16)

- There are 42 (6*7) Type 1 blocks, with a total of 336 (8*42) warps
- They all have 7 phases, so there are 2,352 (336*7) warp-phases
- The warps have control divergence only in their last phase
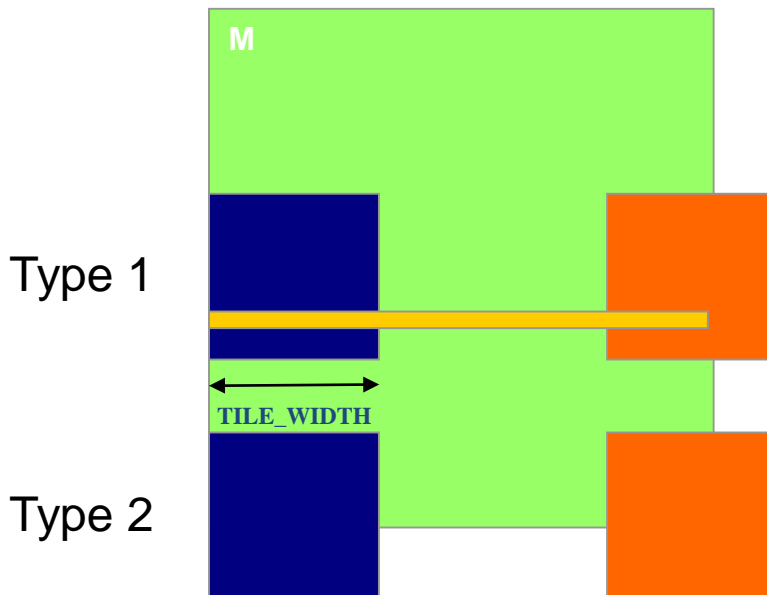- 336 warp-phases have control divergence

Type 1

M

TILE_WIDTH

# Control Divergence in Loading M Tiles (Type 2)

- Type 2: the 7 block assigned to load the bottom tiles, with a total of 56 (8*7) warps
- They all have 7 phases, so there are 392 (56*7) warp-phases
- The first 2 warps in each Type 2 block will stay within the valid range until the last phase
- The 6 remaining warps stay outside the valid range
- So, only 14 (2*7) warp-phases have control divergence

**M**

**TILE_WIDTH**

Type 2

# Overall Impact of Control Divergence

- Type 1 Blocks: 336 out of 2,352 warp-phases have control divergence
- Type 2 Blocks: 14 out of 392 warp-phases have control divergence
- The performance impact is expected to be less than 12% (350/2,944 or (336+14)/(2352+14))

# Additional Comments

– The calculation of impact of control divergence in loading N tiles is somewhat different and is left as an exercise

– The estimated performance impact is data dependent.
  – For larger matrices, the impact will be significantly smaller

– In general, the impact of control divergence for boundary condition checking for large input data sets should be insignificant
  – One should not hesitate to use boundary checks to ensure full functionality

– The fact that a kernel is full of control flow constructs does not mean that there will be heavy occurrence of control divergence

– We will cover some algorithm patterns that naturally incur control divergence (such as parallel reduction) later.

# Objective

– To learn the parallel reduction pattern
  – An important class of parallel computation
  – Work efficiency analysis
  – Resource efficiency analysis

# "Partition and Summarize"

- A commonly used strategy for processing large input data sets
  - There is no required order of processing elements in a data set (associative and commutative)
  - Partition the data set into smaller chunks
  - Have each thread to process a chunk
  - Use a reduction tree to summarize the results from each chunk into the final answer
- E.G., Google and Hadoop MapReduce frameworks support this strategy
- We will focus on the reduction tree step for now

# What is a reduction computation?

– Summarize a set of input values into one value using a "reduction operation"
  – Max
  – Min
  – Sum
  – Product

– Often used with a user defined reduction operation function as long as the operation
  – Is associative(结合律) and commutative(交换律)
  – Has a well-defined identity value (e.g., 0 for sum)
  – For example, the user may supply a custom "max" function for 3D coordinate data sets where the magnitude for the each coordinate data tuple is the distance from the origin.
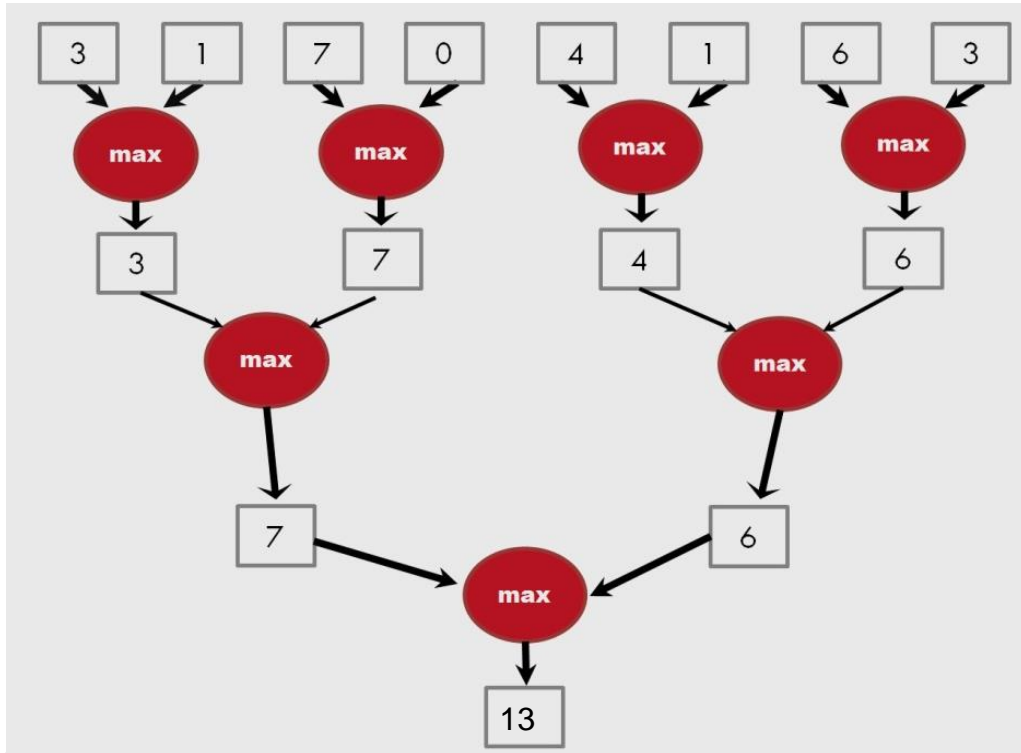
电子科技大学
University of Electronic Science and Technology of China

# An Efficient Sequential Reduction O(N)

– Initialize the result as an identity value for the reduction operation
  – Smallest possible value for max reduction
  – Largest possible value for min reduction
  – 0 for sum reduction
  – 1 for product reduction
– Iterate through the input and perform the reduction operation between the result value and the current input value
  – Each reduction operations performed for N input values
  – an O(N) algorithm
  – This is a computationally efficient algorithm.

# A parallel reduction tree algorithm performs N-1 operations in log(N) steps

# A Quick Analysis

– ## For N input values, the reduction tree performs

  – $(1/2)N + (1/4)N + (1/8)N + \ldots (1/N)N = (1- (1/N))N = N-1$ operations
  – In Log (N) steps – 1,000,000 input values take 20 steps
    – Assuming that we have enough execution resources
  – Average Parallelism (N-1)/Log(N))
    – For N = 1,000,000, average parallelism is 50,000
    – However, peak resource requirement is 500,000
    – This is not resource efficient

– ## This is a work-efficient parallel algorithm

  – The amount of work done is comparable to the an efficient sequential algorithm
  – Many parallel algorithms are not work efficient
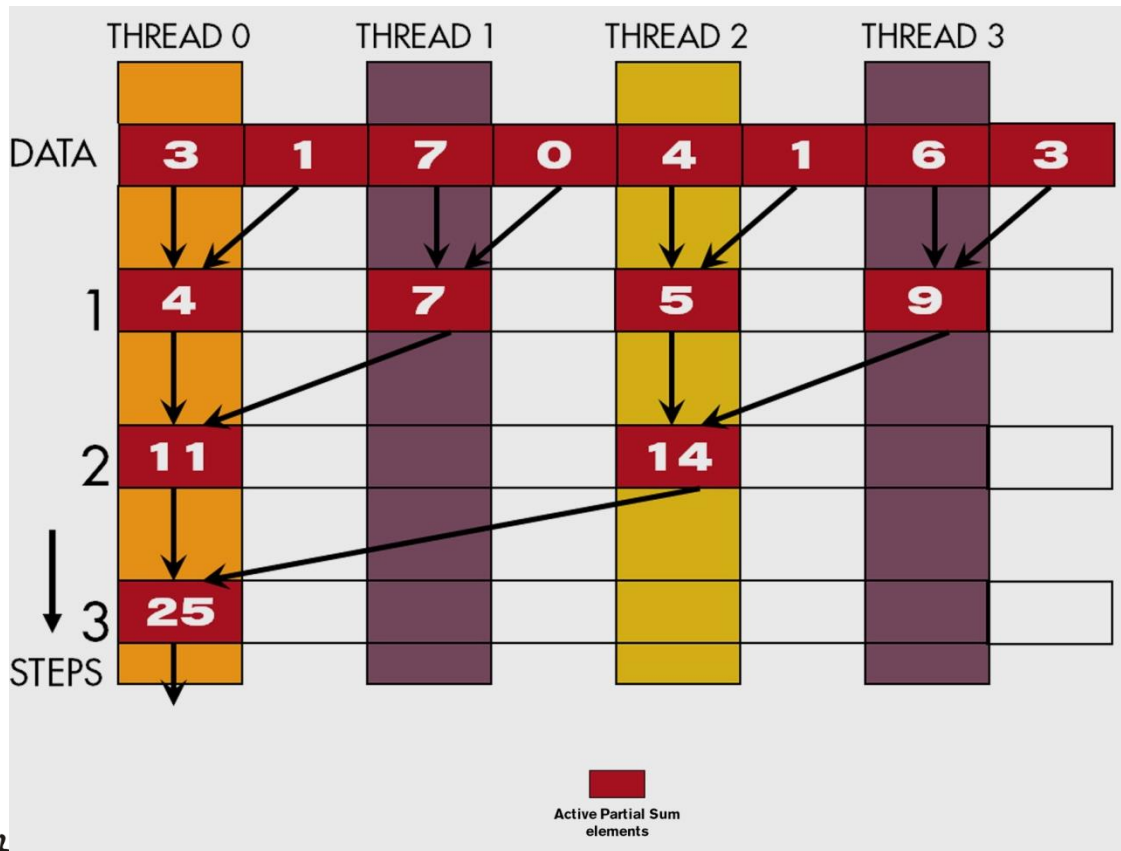
How to implement with CUDA?

# Parallel Sum Reduction

- Parallel implementation
  - Each thread adds two values in each step
  - Recursively halve # of threads
  - Takes log(n) steps for n elements, requires n/2 threads
- Assume an in-place reduction using shared memory
  - The original vector is in device global memory
  - The shared memory is used to hold a partial sum vector
  - Initially, the partial sum vector is simply the original vector
  - Each step brings the partial sum vector closer to the sum
  - The final sum will be in element 0 of the partial sum vector
  - Reduces global memory traffic due to reading and writing partial sum values
  - Thread block size limits n to be less than or equal to 2,048

# A Parallel Sum Reduction Example



Active Partial Sum elements

# A Naive Thread to Data Mapping

- – Each thread is responsible for an even-index location of the partial sum vector (location of responsibility)
- – After each step, half of the threads are no longer needed
- – One of the inputs is always from the location of responsibility
- – In each step, one of the inputs comes from an increasing distance away

# A Simple Thread Block Design

– Each thread block takes 2*BlockDim.x input elements
– Each thread loads 2 elements into shared memory

```
__shared__ float partialSum[2*BLOCK_SIZE];

unsigned int t = threadIdx.x;
unsigned int start = 2*blockIdx.x*blockDim.x;
partialSum[t] = input[start + t];
partialSum[blockDim.x+t] = input[start + blockDim.x+t];
```

# The Reduction Steps

```
for (unsigned int stride = 1;
     stride <= blockDim.x;  stride *= 2)
{
  __syncthreads();
  if (t % stride == 0)
    partialSum[2*t]+= partialSum[2*t+stride];
}
```

Why do we need __syncthreads()?

# Barrier Synchronization

- – __syncthreads() is needed to ensure that all elements of each version of partial sums have been generated before we proceed to the next step

# Back to the Global Picture

– At the end of the kernel, Thread 0 in each block writes the sum of the thread block in partialSum[0] into a vector indexed by the blockIdx.x

– There can be a large number of such sums if the original vector is very large

    – The host code may iterate and launch another kernel

– If there are only a small number of sums, the host can simply transfer the data back and add them together

– Alternatively, Thread 0 of each block could use atomic operations to accumulate into a global sum variable.

# Objective

– To learn to write a better reduction kernel

- – Improved resource efficiency
- – Improved thread to data mapping
- – Reduced control divergence

# Some Observations on the naïve reduction kernel

- In each iteration, two control flow paths will be sequentially traversed for each warp
  - Threads that perform addition and threads that do not
  - Threads that do not perform addition still consume execution resources
- Half or fewer of threads will be executing after the first step
  - All odd-index threads are disabled after first step
  - After the 5th step, entire warps in each block will fail the `if` test, poor resource utilization but no divergence
    - This can go on for a while, up to 5 more steps (stride = 64, 128, 256, 512,1024), where each active warp only has one productive thread until all warps in a block retire

```
if (t % stride == 0)
    partialSum[2*t]+=
    partialSum[2*t+stride];
```
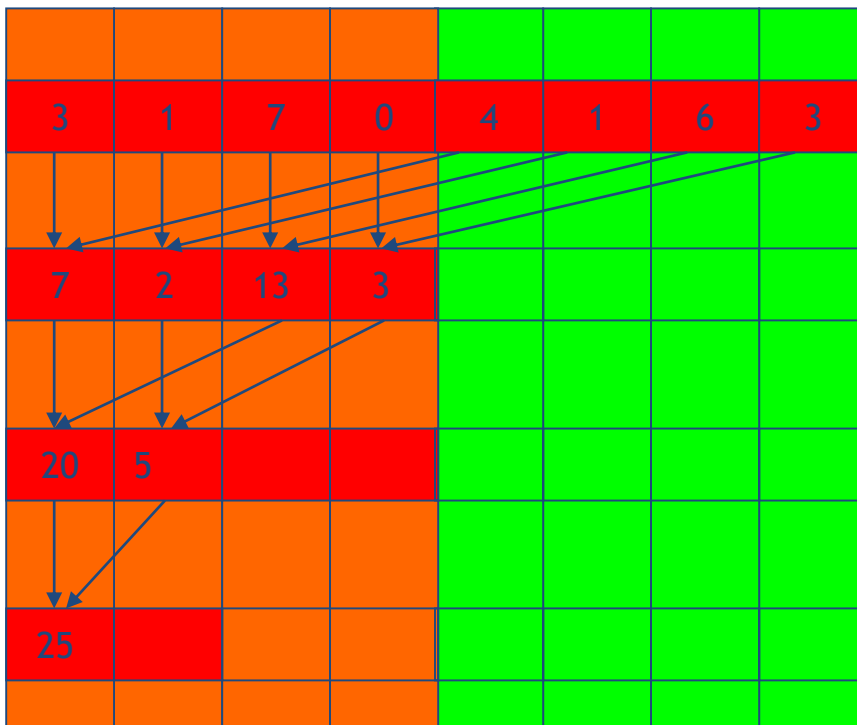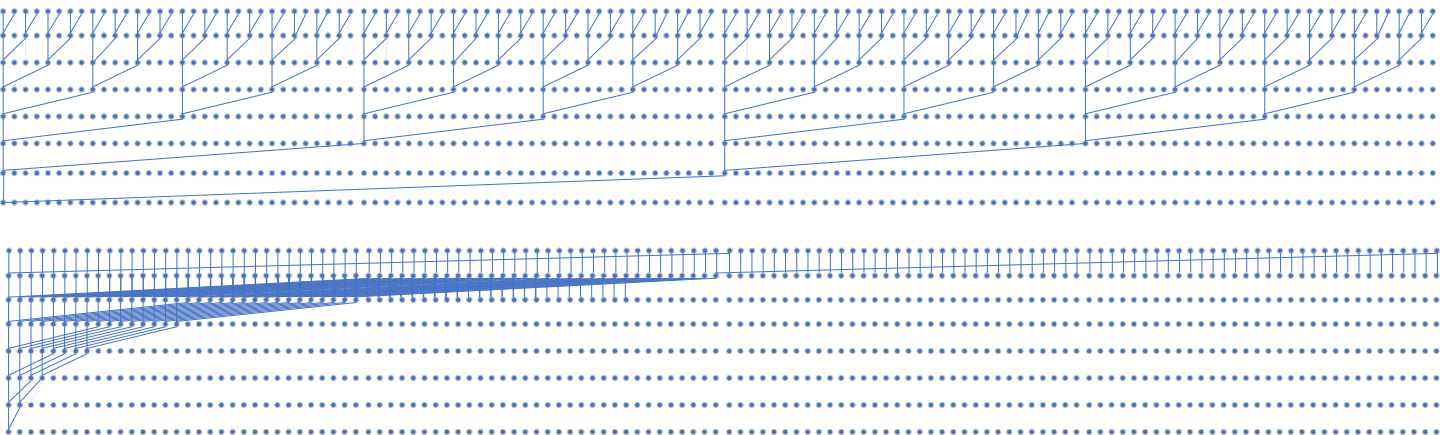
# Thread Index Usage Matters

– In some algorithms, one can shift the index usage to improve the divergence behavior

  – Commutative and associative operators

– Always compact the partial sums into the front locations in the partialSum[ ] array

– Keep the active threads consecutive

# An Example of 4 threads

# A Better Reduction Kernel

```
for (unsigned int stride = blockDim.x;
      stride > 0;  stride /= 2)
{
  __syncthreads();
  if (t < stride)
     partialSum[t] += partialSum[t+stride];
}
```

# A Quick Analysis

– For a 1024 thread block
  – No divergence in the first 6 steps
    – 1024, 512, 256, 128, 64, 32 consecutive threads are active in each step
    – All threads in each warp  either all active or all inactive
  – The final 5 steps will still have divergence

Warps and SIMD

performance impact of control divergence

Parallel reduction

Memory parallelism

# Objective

– To learn that memory bandwidth is a first-order performance factor in a massively parallel processor

   – DRAM bursts, banks, and channels
   – All concepts are also applicable to modern multicore processors

电子科技大学
University of Electronic Science and Technology of China

# Global Memory (DRAM) Bandwidth
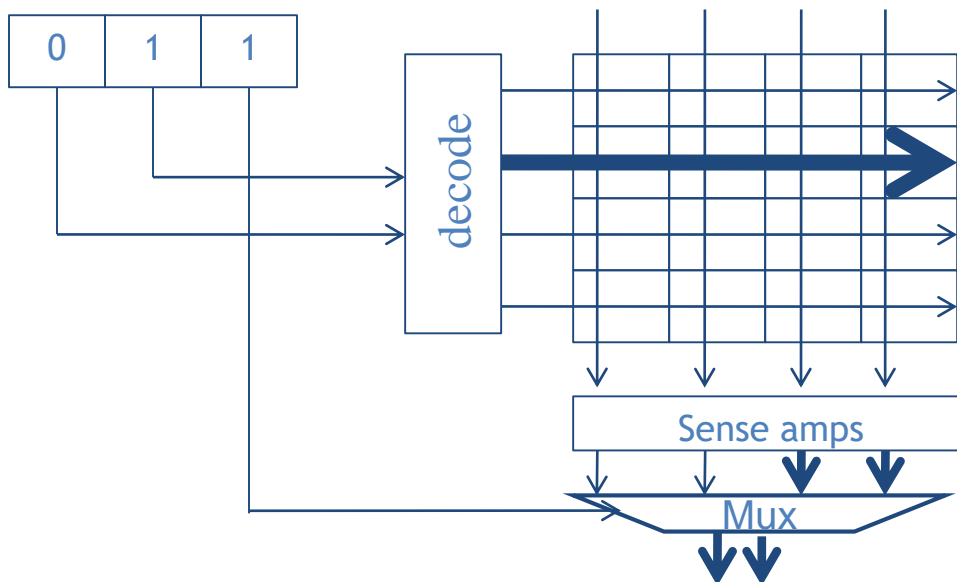
– Ideal



– Reality

# DRAM Core Array Organization

– Each DRAM core array has about 16M bits

– Each bit is stored in a tiny capacitor made of one transistor（1T1C）
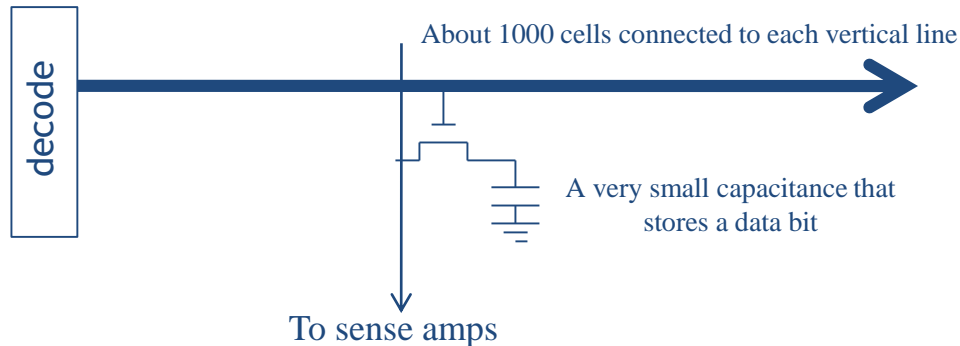
# A very small (8x2-bit) DRAM Core Array

# DRAM Core Arrays are Slow

– Reading from a cell in the core array is a very slow process
  - DDR: Core speed = ½ interface speed
  - DDR2/GDDR3: Core speed = ¼ interface speed
  - DDR3/GDDR4: Core speed = ⅛ interface speed
  - … likely to be worse in the future

About 1000 cells connected to each vertical line

decode

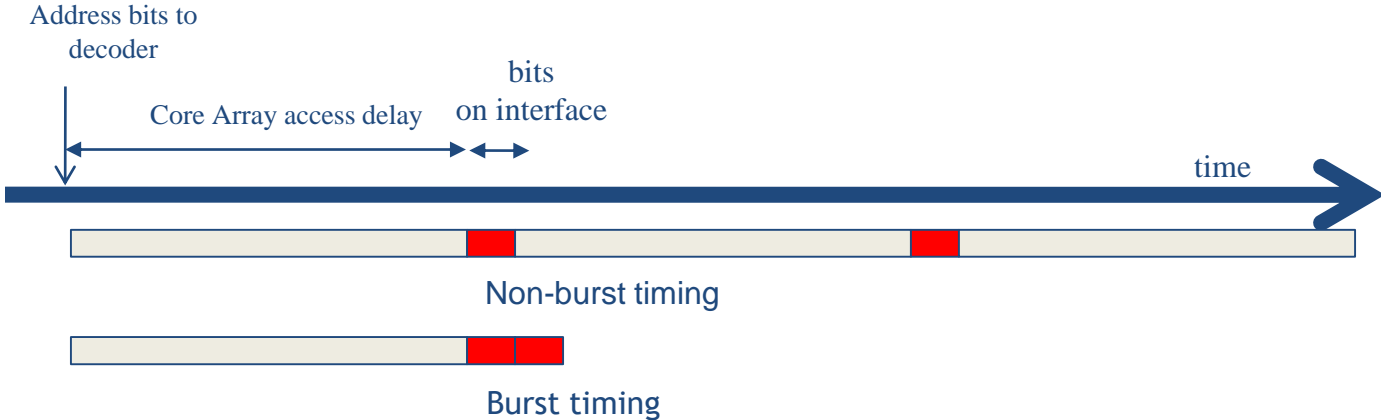A very small capacitance that
stores a data bit

To sense amps

# DRAM Bursting

– For DDR{2,3} SDRAM cores clocked at 1/N speed of the interface:

    – Load (N × interface width) of DRAM bits from the same row at once to an internal buffer, then transfer in N steps at interface speed
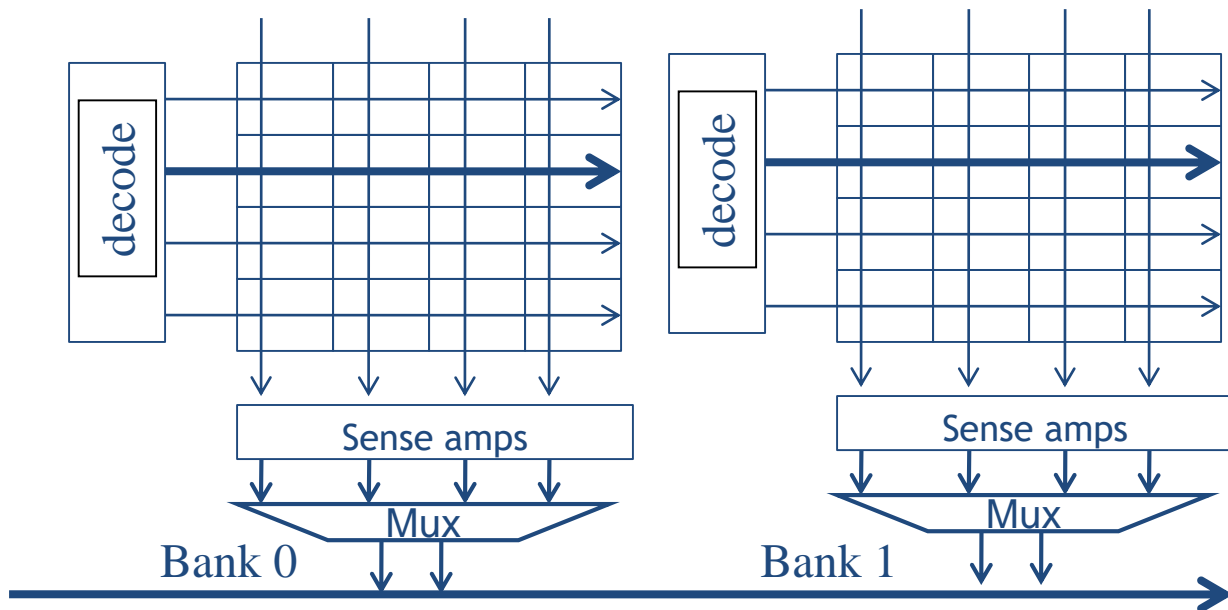
    – DDR3/GDDR4: buffer width = 8X interface width

# DRAM Bursting Timing Example

Address bits to decoder

Core Array access delay

bits on interface

time

Non-burst timing

Burst timing

Modern DRAM systems are designed to always be accessed in burst mode. Burst bytes are transferred to the processor but discarded when accesses are not to sequential locations.

# Multiple DRAM Banks



decode

Sense amps

Mux

Bank 0

decode

Sense amps

Mux

Bank 1

# DRAM Bursting with Banking

Single-Bank burst timing, dead time on interface

Multi-Bank burst timing, reduced dead time

# GPU off-chip memory subsystem

- NVIDIA GTX280 GPU:
  - Peak global memory bandwidth = 141.7GB/s

- Global memory (GDDR3) interface @ 1.1GHz
  - (Core speed @ 276Mhz)
  - For a typical 64-bit interface, we can sustain only about 17.6 GB/s (Recall DDR - 2 transfers per clock)
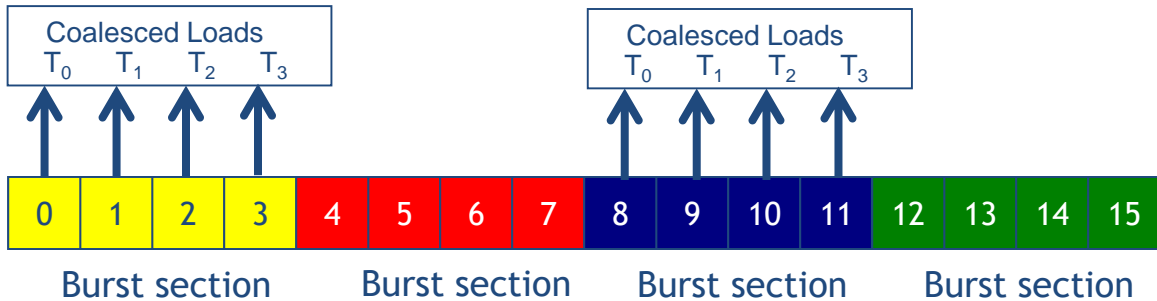  - We need a lot more bandwidth (141.7 GB/s) – thus 8 memory channels

电子科技大学
University of Electronic Science and Technology of China

# DRAM Burst – A System View

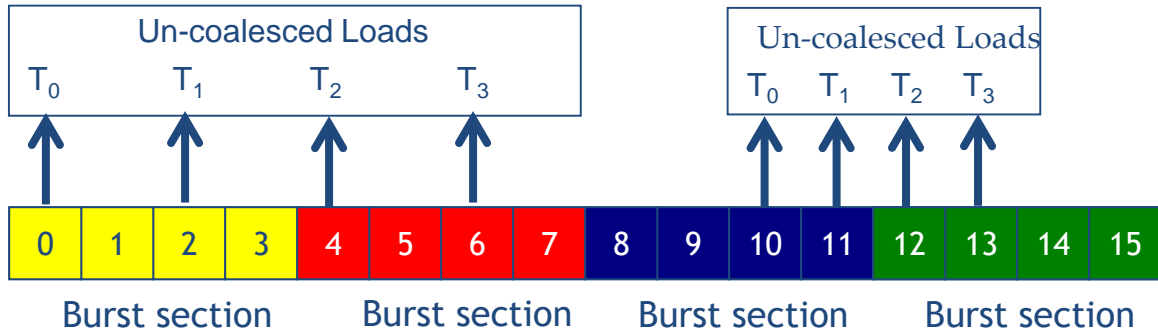| Burst section | | | | Burst section | | | | Burst section | | | | Burst section | | | |
|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|:-:|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

– Each address space is partitioned into burst sections
  – Whenever a location is accessed, all other locations in the same section are also delivered to the processor
– Basic example: a 16-byte address space, 4-byte burst sections
  – In practice, we have at least 4GB address space,  burst section sizes of 128-bytes or more

# Memory Coalescing



- When all threads of a warp execute a load instruction, if all accessed locations fall into the same burst section, only one DRAM request will be made and the access is fully coalesced.

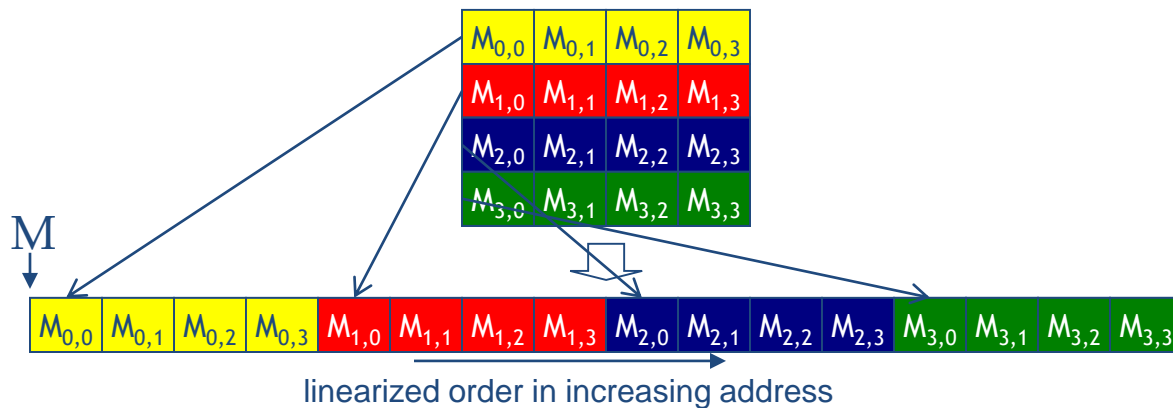# Un-coalesced Accesses



- When the accessed locations spread across burst section boundaries:
  - Coalescing fails
  - Multiple DRAM requests are made
  - The access is not fully coalesced.
- Some of the bytes accessed and transferred are not used by the threads
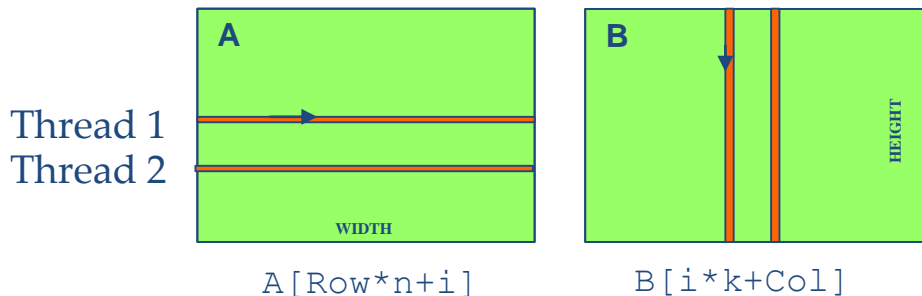
# How to judge if an access is coalesced?

– Accesses in a warp are to consecutive locations if the index in an array access is in the form of

 – A[(expression with terms independent of threadIdx.x) + threadIdx.x];

# A 2D C Array in Linear Memory Space



linearized order in increasing address

# Two Access Patterns of Basic Matrix Multiplication



Thread 1
Thread 2

**A**  WIDTH

**B**  HEIGHT

A[Row*n+i]          B[i*k+Col]

i is the loop counter in the inner product loop of the kernel code

A is m × n, B is n × k
Col = blockIdx.x*blockDim.x + threadIdx.x
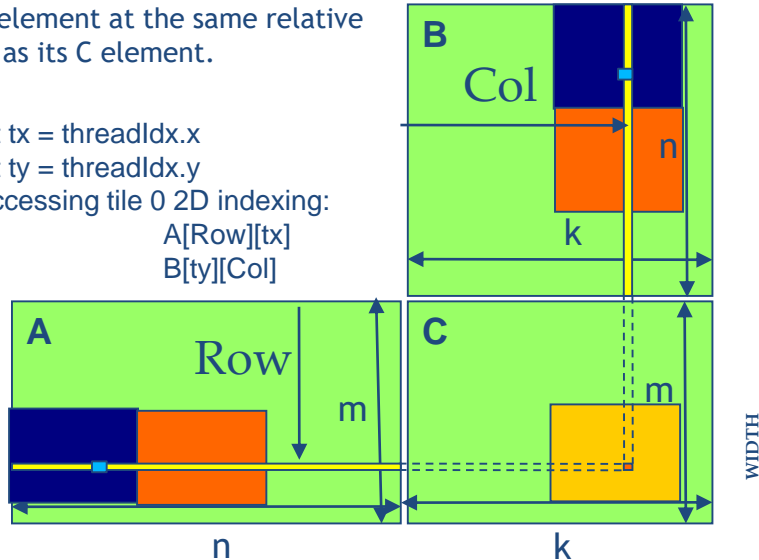
# B accesses are coalesced

# A Accesses are Not Coalesced

# Loading an Input Tile

Have each thread load an A element and a B element at the same relative position as its C element.

int tx = threadIdx.x
int ty = threadIdx.y
Accessing tile 0 2D indexing:
A[Row][tx]
B[ty][Col]

# Corner Turning



Original Access Pattern

d_M — WIDTH

d_N — WIDTH

Tiled Access Pattern

d_M

d_N

Copy into shared memory

Perform multiplication with shared memory values