# Contemporary Computer Architecture TDSN13

LECTURE 7 – CUDA MEMORY ACCESS PATTERNS
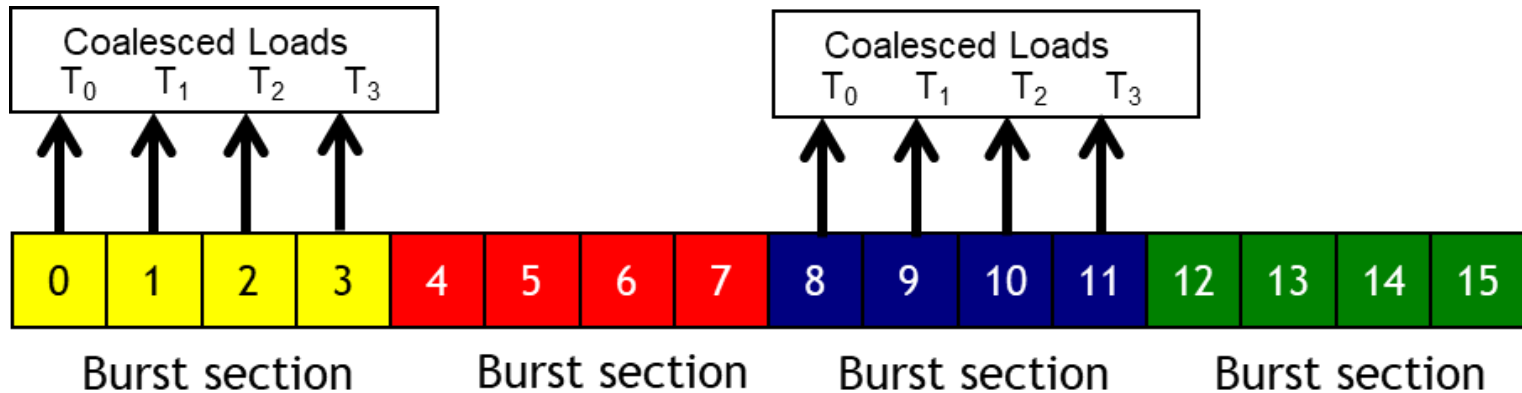
ANDREAS AXELSSON (ANDREAS.AXELSSON@JU.SE)

# DRAM Burst – A System View

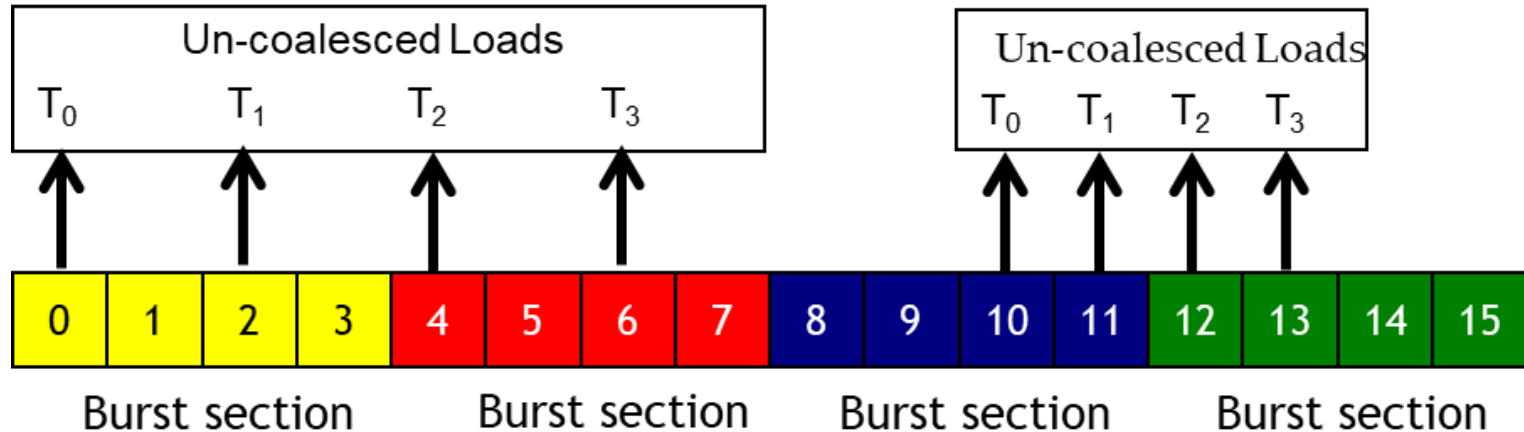| Burst section | | | | Burst section | | | | Burst section | | | | Burst section | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |

– 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.
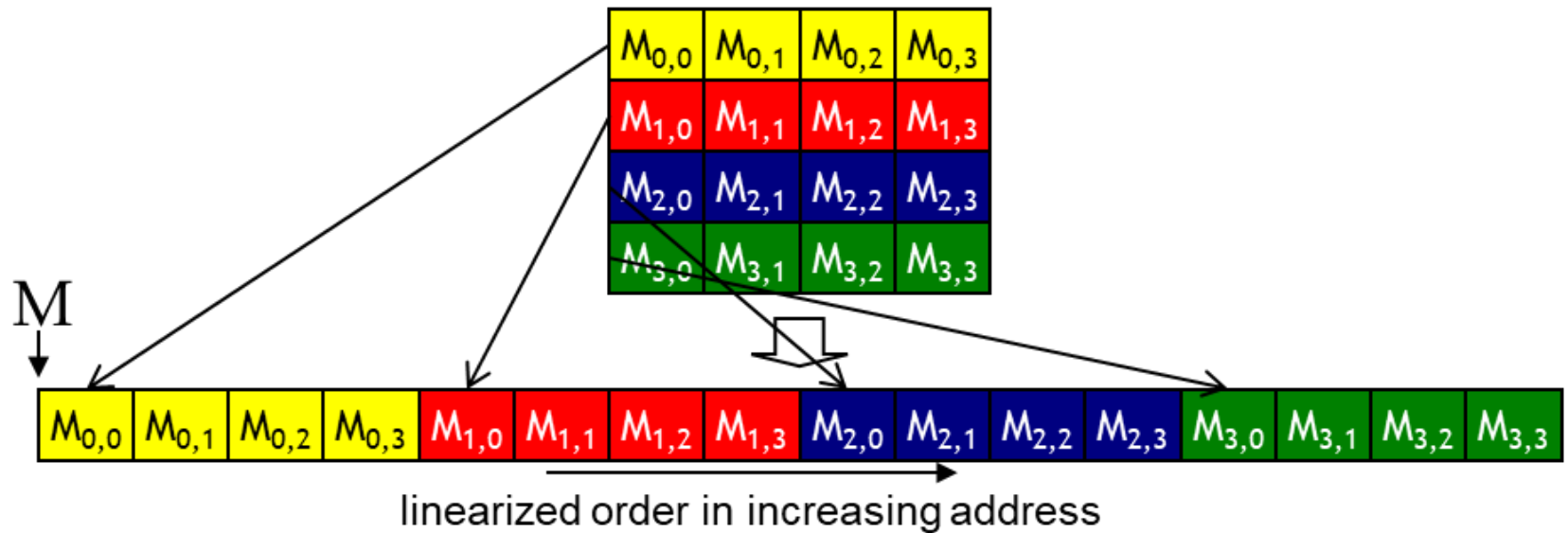
# Memory Coalescing



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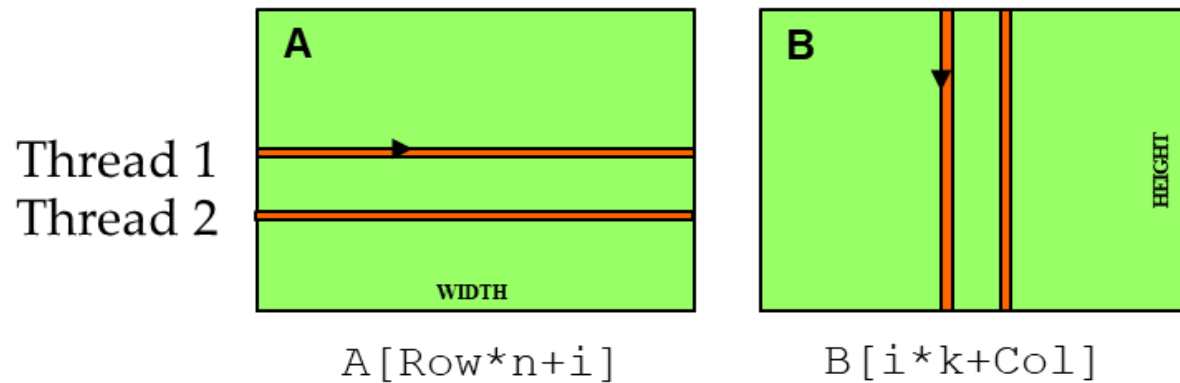- – Some of the bytes accessed and transferred are not used by the threads

# Is access 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];
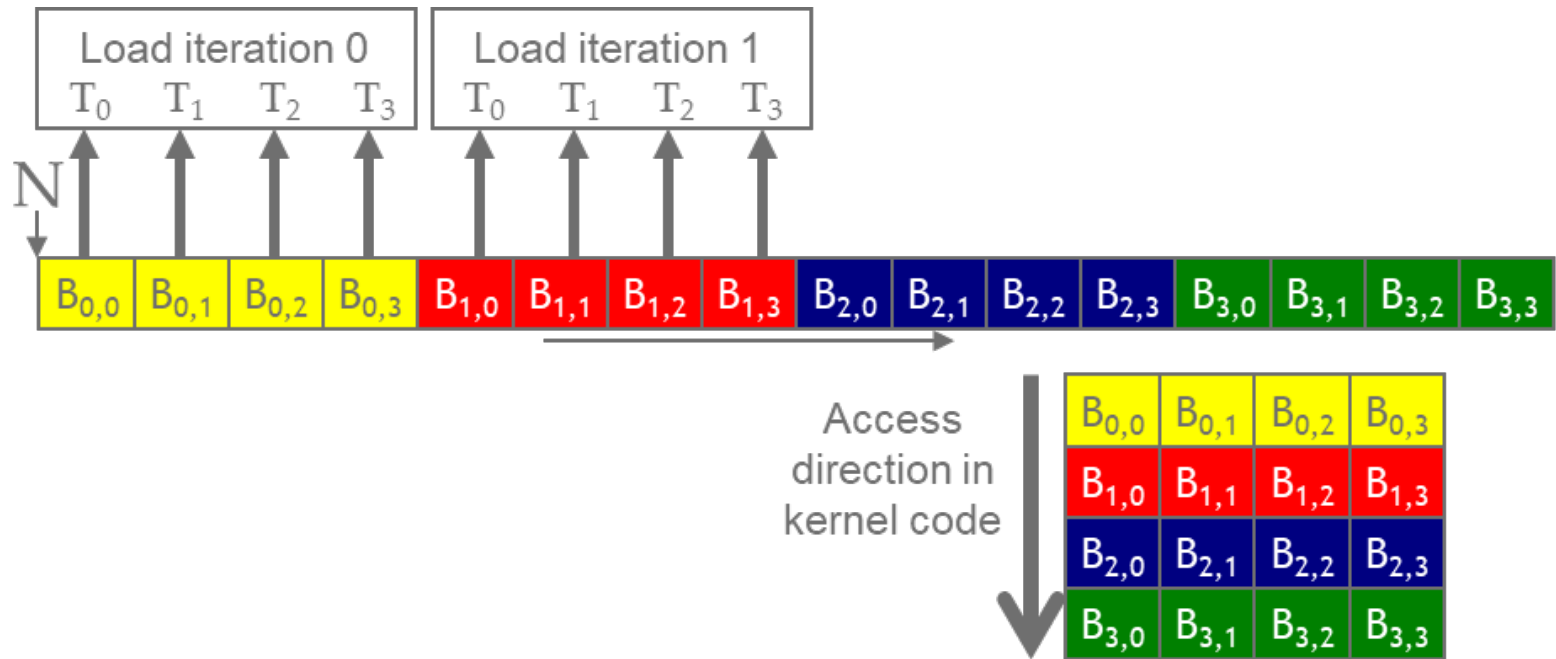
# 2D Array in Linear Memory Space



linearized order in increasing address

# Access Pattern Multiplication



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



- Coalesced accesses in a warp are to consecutive locations if the index in an array access is in the form of
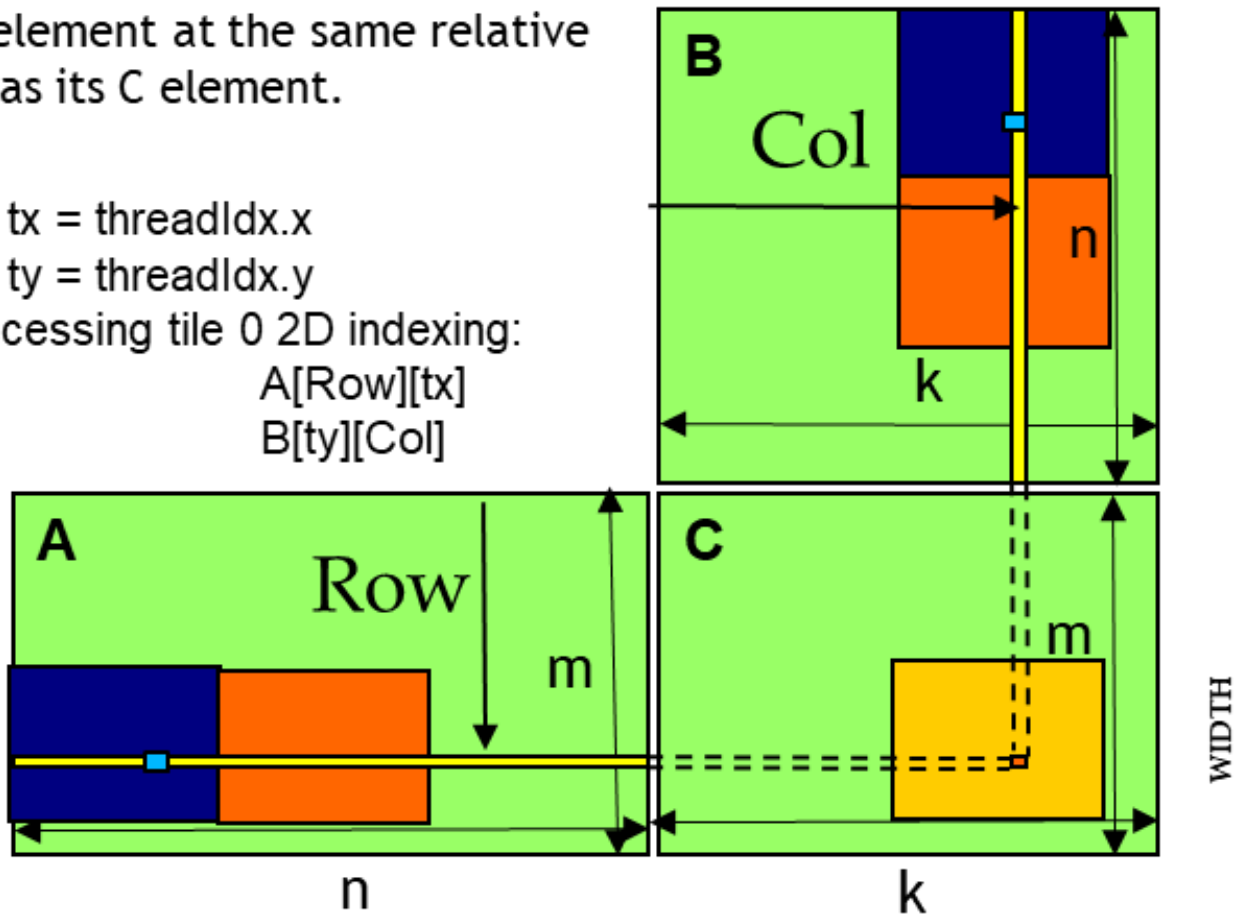  - A[(expression with terms independent of threadIdx.x) + threadIdx.x];

# Corner Turning

– What if we must iterate through data along the row direction?
  – We can use the shared memory to enable memory coalescing — this is the technique called Corner Turning

– how?
  – Use a tiled algorithm
  – Threads of a block can first cooperatively load the tiles into the shared memory
  – Care must be taken to ensure that these tiles are loaded in a coalesced pattern
  – Once the data is in shared memory, they can be accessed either on a row basis or a column basis with much less performance variation because the shared memories are implemented as intrinsically high-speed on-chip memory that does not require coalescing to achieve high data access rate.
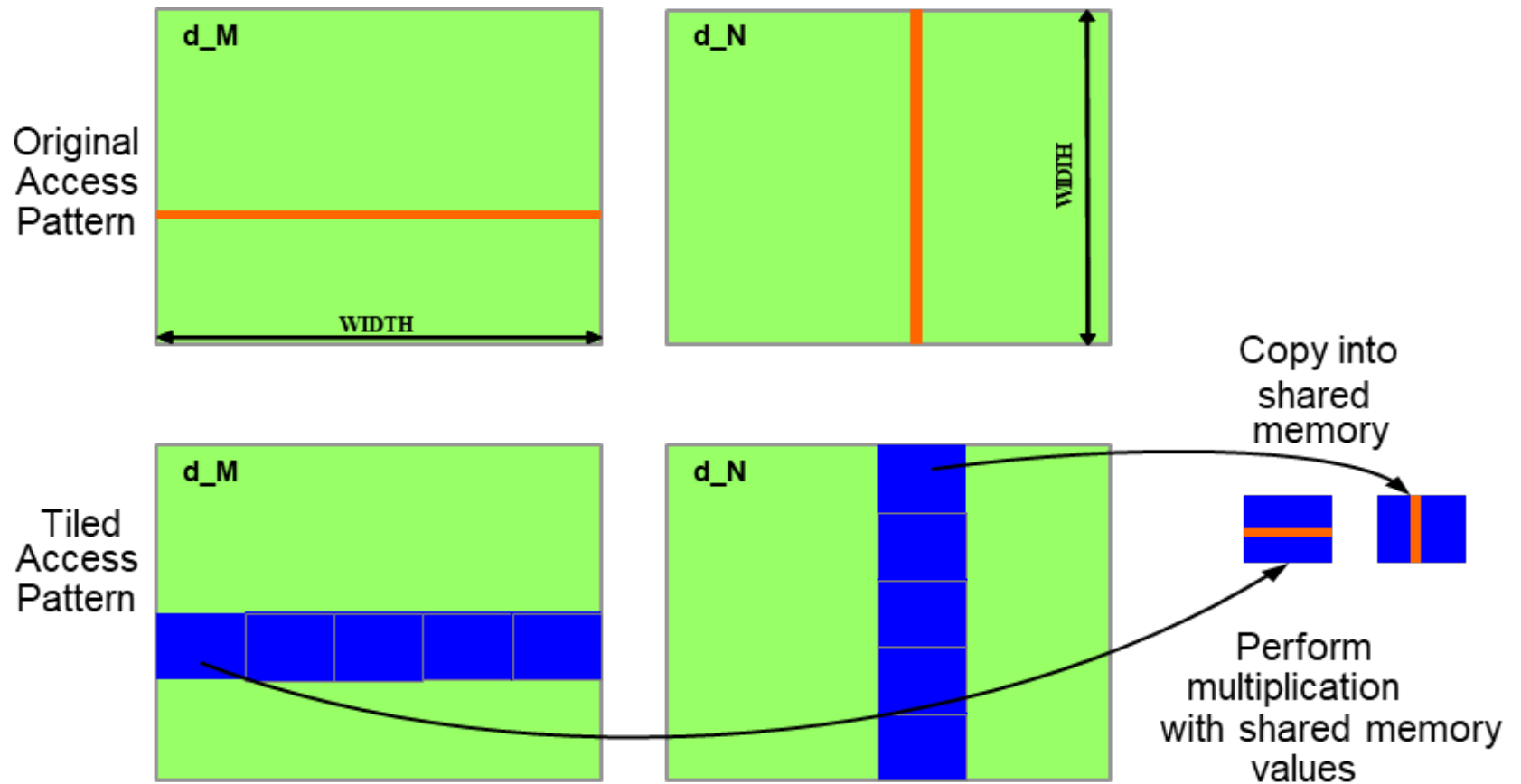
# Corner Turning – Loading a tile



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

# Tiled Matrix Multiplication Kernel

```
__global__ void MatrixMulKernel(float* M, float* N, float* P, Int Width)
{
    __shared__ float ds_M[TILE_WIDTH][TILE_WIDTH];
    __shared__ float ds_N[TILE_WIDTH][TILE_WIDTH];

    int bx = blockIdx.x;  int by = blockIdx.y;
    int tx = threadIdx.x; int ty = threadIdx.y;

    int Row = by * blockDim.y + ty;
    int Col = bx * blockDim.x + tx;
    float Pvalue = 0;

    // Loop over the M and N tiles required to compute the P element
    for (int p = 0; p < n/TILE_WIDTH; ++p) {
        // Collaborative loading of M and N tiles into shared memory
        ds_M[ty][tx] = M[Row*Width + p*TILE_WIDTH+tx];
        ds_N[ty][tx] = N[(t*TILE_WIDTH+ty)*Width + Col];
        __syncthreads();

        for (int i = 0; i < TILE_WIDTH; ++i)Pvalue += ds_M[ty][i] *
        ds_N[i][tx];
        __syncthreads();
    }
    P[Row*Width+Col] = Pvalue;
}
```

Remember that this Kernel runs TILE_WIDTH x TILE_WIDTH number of threads in parallel in each block (SM)

# Tiled Matrix Multiplication Kernel

```
__global__ void MatrixMulKernel(float* M, float* N, float* P, Int Width)
{
    __shared__ float ds_M[TILE_WIDTH][TILE_WIDTH];
    __shared__ float ds_N[TILE_WIDTH][TILE_WIDTH];

    int bx = blockIdx.x;  int by = blockIdx.y;
    int tx = threadIdx.x; int ty = threadIdx.y;

    int Row = by * blockDim.y + ty;
    int Col = bx * blockDim.x + tx;
    float Pvalue = 0;

    // Loop over the M and N tiles required to compute the P element
    for (int p = 0; p < n/TILE_WIDTH; ++p) {
        // Collaborative loading of M and N tiles into shared memory
        ds_M[ty][tx] = M[Row*Width + p*TILE_WIDTH+tx];
        ds_N[ty][tx] = N[(t*TILE_WIDTH+ty)*Width + Col];
        __syncthreads();

        for (int i = 0; i < TILE_WIDTH; ++i)Pvalue += ds_M[ty][i] *
        ds_N[i][tx];
        __syncthreads();
    }
    P[Row*Width+Col] = Pvalue;
}
```

# Tiled Matrix Multiplication Kernel

```
__global__ void MatrixMulKernel(float* M, float* N, float* P, Int Width)
{
    __shared__ float ds_M[TILE_WIDTH][TILE_WIDTH];
    __shared__ float ds_N[TILE_WIDTH][TILE_WIDTH];

    int bx = blockIdx.x;   int by = blockIdx.y;
    int tx = threadIdx.x; int ty = threadIdx.y;

    int Row = by * blockDim.y + ty;
    int Col = bx * blockDim.x + tx;
    float Pvalue = 0;

// Loop over the M and N tiles required to compute the P element
 for (int p = 0; p < n/TILE_WIDTH; ++p) {
    // Collaborative loading of M and N tiles into shared memory
    ds_M[ty][tx] = M[Row*Width + p*TILE_WIDTH+tx];
    ds_N[ty][tx] = N[(t*TILE_WIDTH+ty)*Width + Col];
    __syncthreads();

    for (int i = 0; i < TILE_WIDTH; ++i)Pvalue += ds_M[ty][i] *
    ds_N[i][tx];
    __syncthreads();
  }
  P[Row*Width+Col] = Pvalue;
}
```
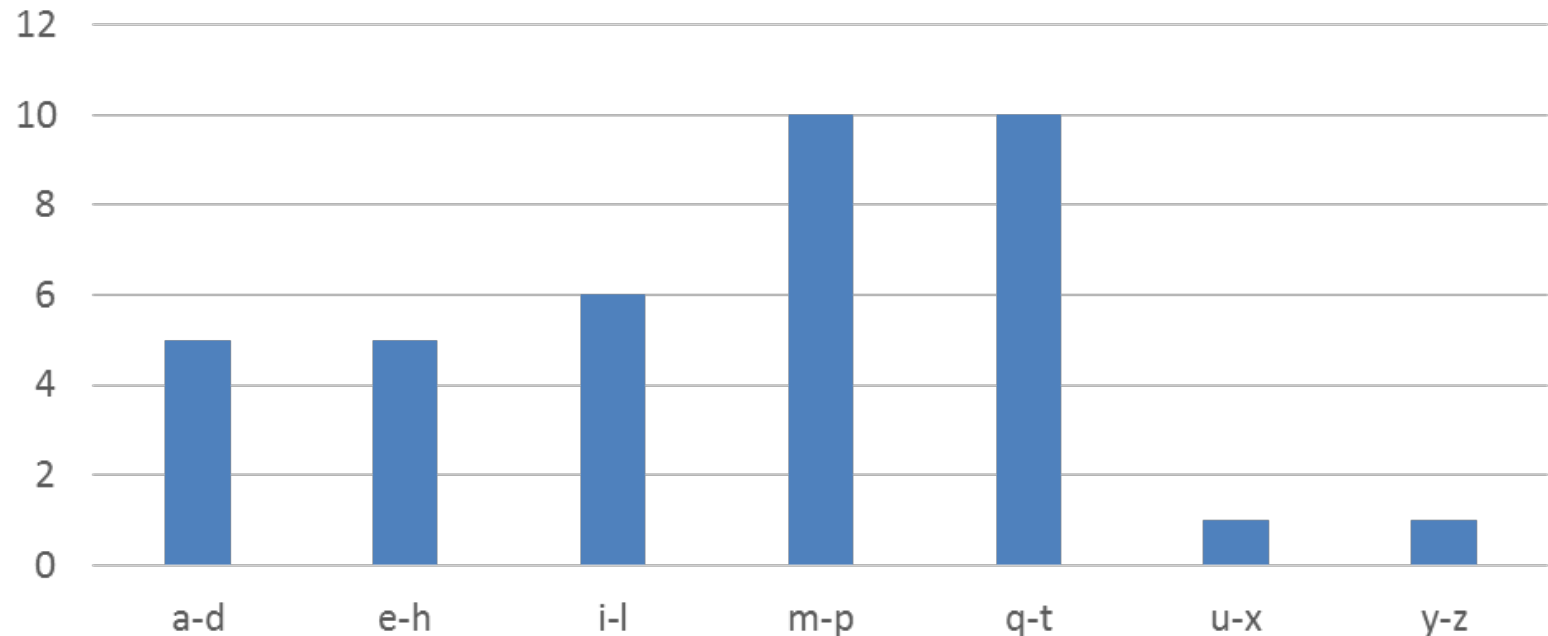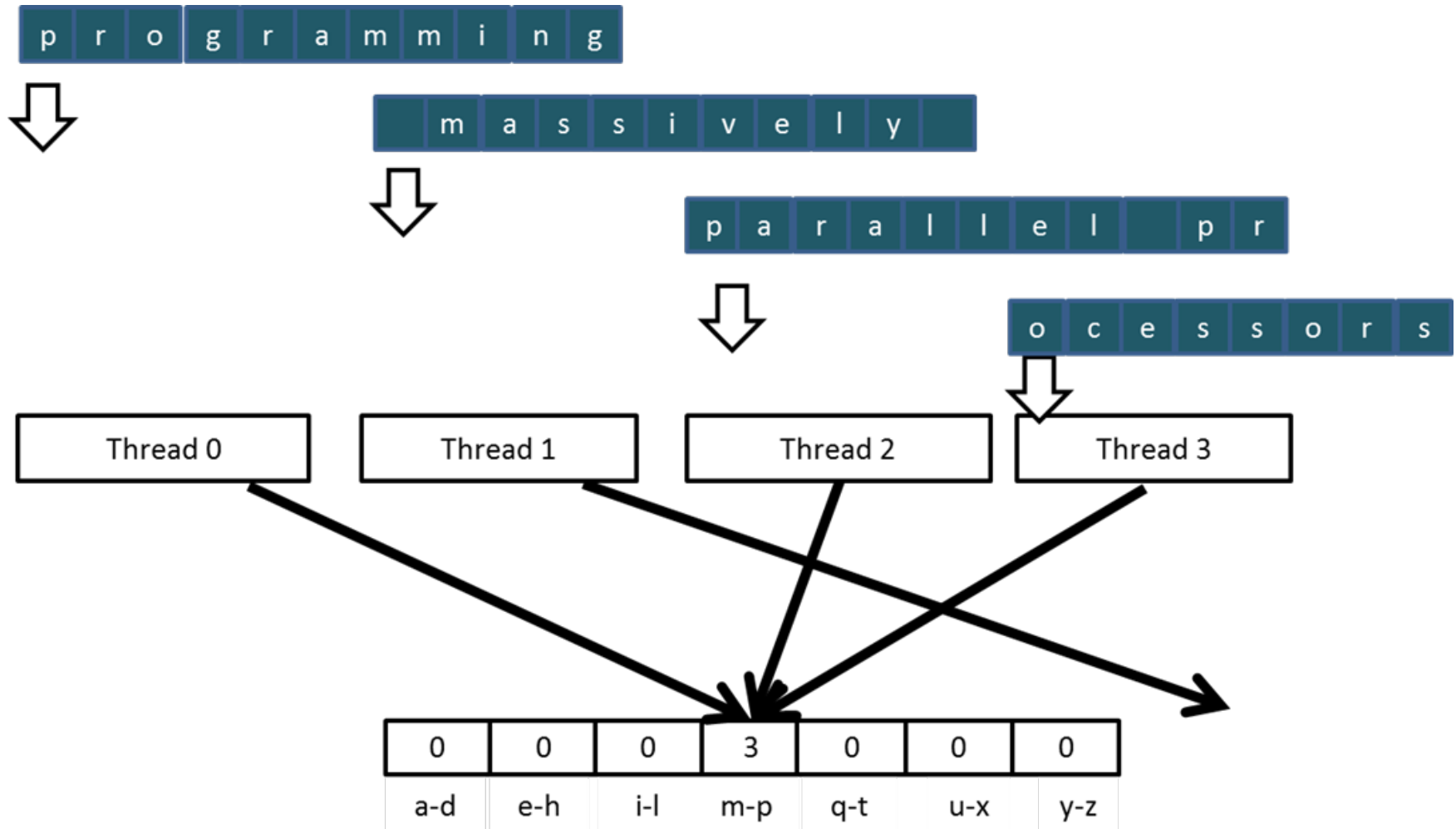
# Histogram

- A method for extracting notable features and patterns from large data sets
  - Feature extraction for object recognition in images
  - Fraud detection in credit card transactions
  - Correlating heavenly object movements in astrophysics
  - …

- Basic histograms - for each element in the data set, use the value to identify a "bin counter" to increment
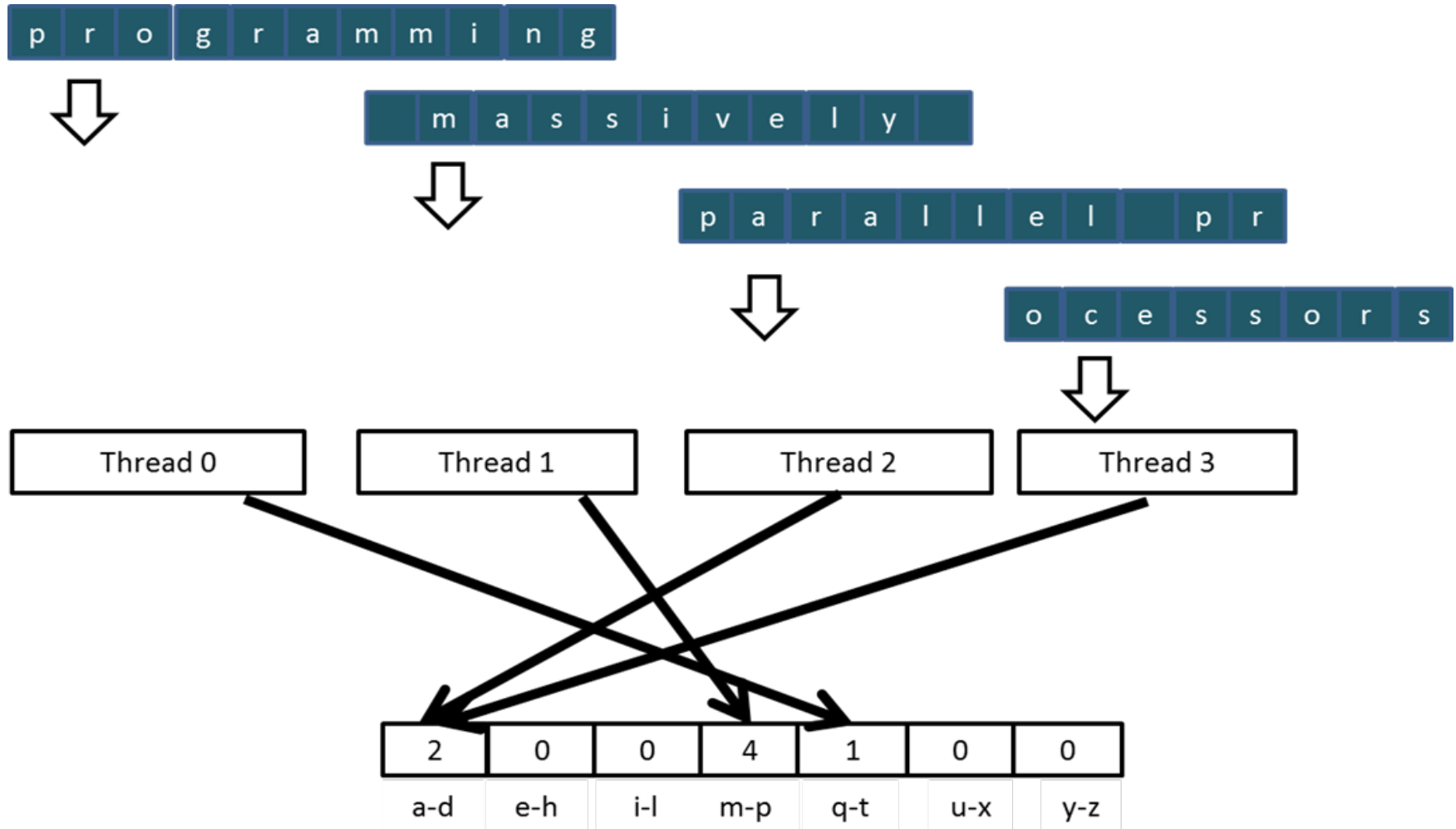
# A Text Histogram Example

– Define the bins as four-letter sections of the alphabet: a-d, e-h, i-l, n-p, …
– For each character in an input string, increment the appropriate bin counter.
– In the phrase "Programming Massively Parallel Processors" the output histogram is shown below:

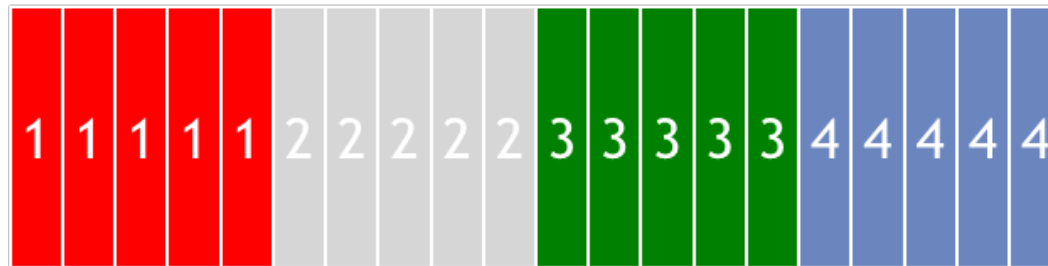# Sectioned Partition (Iteration #1)

# Sectioned Partition (Iteration #2)

# Partitioning Access Efficiency

- Sectioned partitioning results in poor memory access efficiency
  - Adjacent threads do not access adjacent memory locations
  - Accesses are not coalesced
  - DRAM bandwidth is poorly utilized

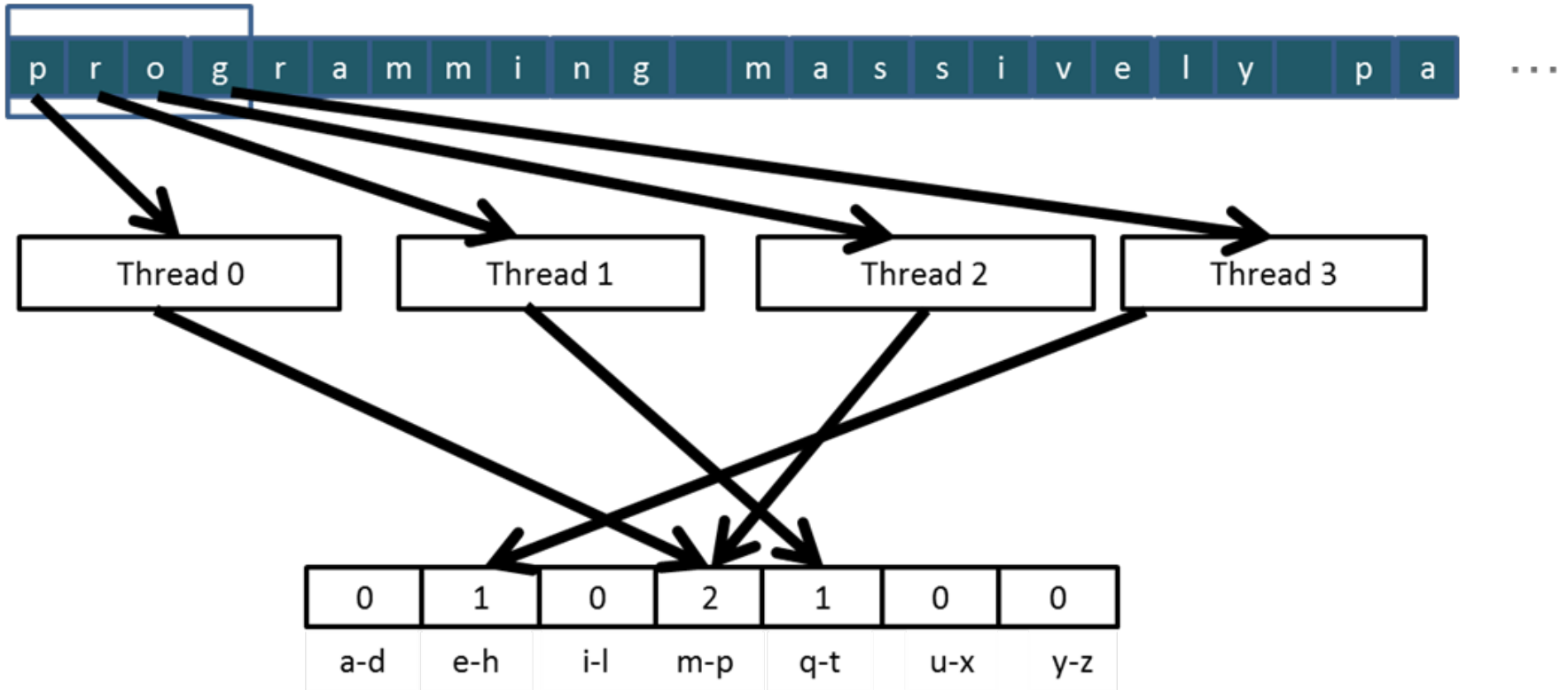| 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 4 | 4 |

- Change to interleaved partitioning
  - All threads process a contiguous section of elements
  - They all move to the next section and repeat
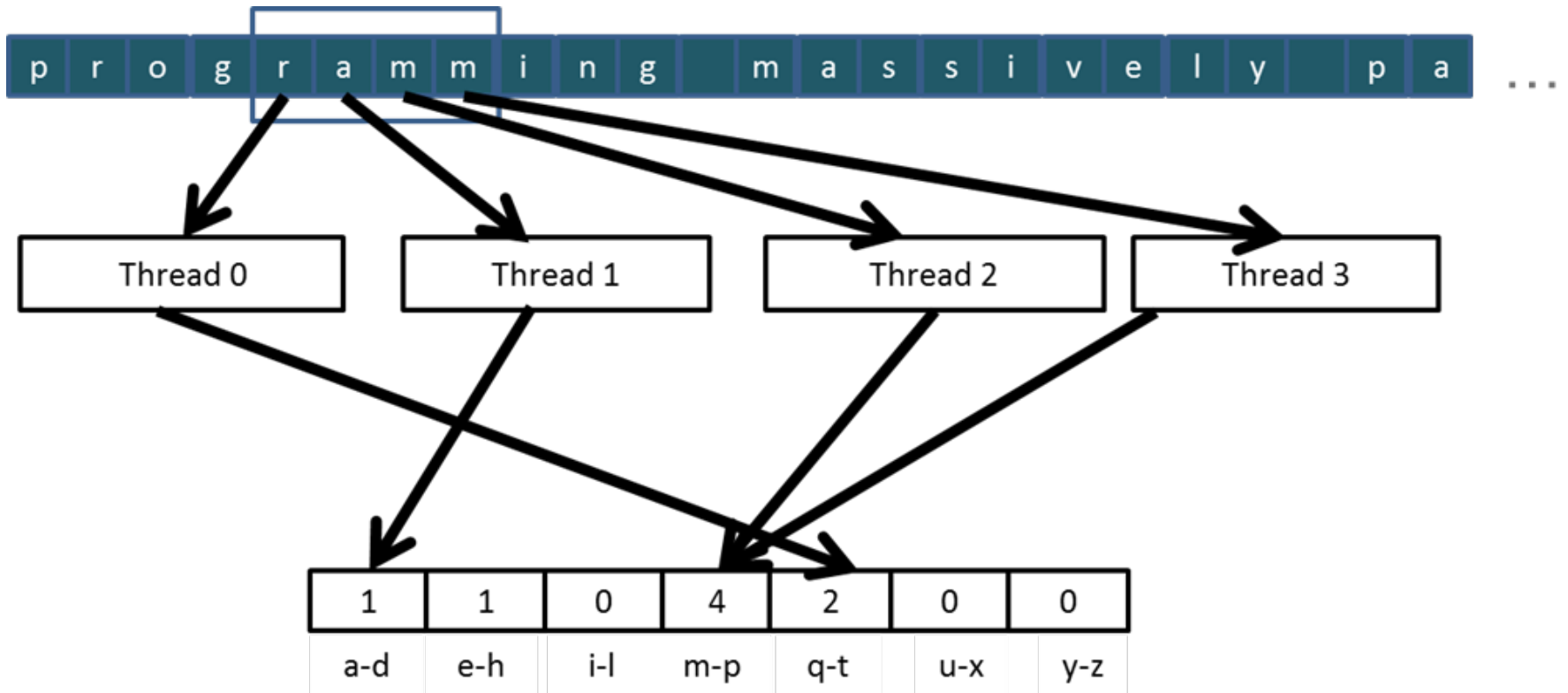  - The memory accesses are coalesced

| 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |

# Interleaved Partitioning of Input

– For coalescing and better memory access performance

# Interleaved Partitioning (Iter #2)

# Data Race in Parallel Threads

thread1: Old ← Mem[x]
          New ← Old + 1
          Mem[x] ← New

thread2: Old ← Mem[x]
          New ← Old + 1
          Mem[x] ← New

Old and New are per-thread register variables.

Question 1: If Mem[x] was initially 0, what would the value of Mem[x] be after threads 1 and 2 have completed?

Question 2: What does each thread get in their Old variable?

Unfortunately, the answers may vary according to the relative execution timing between the two threads, which is referred to as a **data race**.

# Timing Scenario #1

| Time | Thread 1 | Thread 2 |
|------|----------|----------|
| 1 | (0) Old ← Mem[x] | |
| 2 | (1) New ← Old + 1 | |
| 3 | (1) Mem[x] ← New | |
| 4 | | (1) Old ← Mem[x] |
| 5 | | (2) New ← Old + 1 |
| 6 | | (2) Mem[x] ← New |

- Thread 1 Old = 0
- Thread 2 Old = 1
- Mem[x] = 2 after the sequence

# Timing Scenario #2

| Time | Thread 1 | Thread 2 |
|------|----------|----------|
| 1 | | (0) Old ← Mem[x] |
| 2 | | (1) New ← Old + 1 |
| 3 | | (1) Mem[x] ← New |
| 4 | (1) Old ← Mem[x] | |
| 5 | (2) New ← Old + 1 | |
| 6 | (2) Mem[x] ← New | |

– Thread 1 Old = 1
– Thread 2 Old = 0
– Mem[x] = 2 after the sequence

# Timing Scenario #3

| Time | Thread 1 | Thread 2 |
|:---:|:---:|:---:|
| 1 | (0) Old ← Mem[x] | |
| 2 | (1) New ← Old + 1 | |
| 3 | | (0) Old ← Mem[x] |
| 4 | (1) Mem[x] ← New | |
| 5 | | (1) New ← Old + 1 |
| 6 | | (1) Mem[x] ← New |

- Thread 1 Old = 0
- Thread 2 Old = 0
- Mem[x] = 1 after the sequence

# Timing Scenario #4

| Time | Thread 1 | Thread 2 |
|------|----------|----------|
| 1 | | (0) Old ← Mem[x] |
| 2 | | (1) New ← Old + 1 |
| 3 | (0) Old ← Mem[x] | |
| 4 | | (1) Mem[x] ← New |
| 5 | (1) New ← Old + 1 | |
| 6 | (1) Mem[x] ← New | |

- Thread 1 Old = 0
- Thread 2 Old = 0
- Mem[x] = 1 after the sequence

# Key Concept of Atomic Operations

- A read-modify-write operation performed by a single hardware instruction on a memory location *address*
  - Read the old value, calculate a new value, and write the new value to the location
- The hardware ensures that no other threads can perform another read-modify-write operation on the same location until the current atomic operation is complete
  - Any other threads that attempt to perform an atomic operation on the same location will typically be held in a queue
  - All threads perform their atomic operations **serially** on the same location

# Atomic Operations in CUDA

– Performed by calling functions that are translated into single instructions (a.k.a. *intrinsic functions* or *intrinsics*)

  – Atomic add, sub, inc, dec, min, max, exch (exchange), CAS (compare and swap)

  – Read CUDA C programming Guide for details

– Atomic Add

  ```
  int atomicAdd(int* address, int val);
  ```

  – reads the 32-bit word **old** from the location pointed to by **address** in global or shared memory, computes (**old + val**), and stores the result back to memory at the same address. These three operations are performed in one atomic transaction. The function returns **old**.

https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#atomic-functions

# More Atomic Adds in CUDA

– Unsigned 32-bit integer atomic add

```
unsigned int atomicAdd(unsigned int* address,
    unsigned int val);
```

– Unsigned 64-bit integer atomic add

```
unsigned long long int atomicAdd(unsigned long long
    int* address, unsigned long long int val);
```

– Single-precision floating-point atomic add (Compute capability 2.x+)

```
float atomicAdd(float* address, float val);
```

– Double-precision floating-point atomic add (Compute capability 6.x+)

```
double atomicAdd(double* address, double val);
```

– 16-bit floating-point atomic add (Compute capability 7.x+)

```
__half atomicAdd(__half* address, __half val);
```

# A Generic Text Histogram Kernel

- The kernel receives a pointer to the input buffer of byte values
- Each thread process the input in a strided pattern

```
__global__ void histo_kernel(unsigned char *buffer,
        long size, unsigned int *histo)
{
    int i = threadIdx.x + blockIdx.x * blockDim.x;

// stride is total number of threads
    int stride = blockDim.x * gridDim.x;

 // All threads handle blockDim.x * gridDim.x
    // consecutive elements
    while (i < size) {
        atomicAdd( &(histo[buffer[i]]), 1);
        i += stride;
    }
}
```

# Our Text Histogram Kernel (cont.)

- The kernel receives a pointer to the input buffer of byte values
- Each thread process the input in a strided pattern

```
__global__ void histo_kernel(unsigned char *buffer,
        long size, unsigned int *histo)
{
    int i = threadIdx.x + blockIdx.x * blockDim.x;

// stride is total number of threads
    int stride = blockDim.x * gridDim.x;

 // All threads handle blockDim.x * gridDim.x
   // consecutive elements
   while (i < size) {
       int alphabet_position = buffer[i] – "a";
       if (alphabet_position >= 0 && alpha_position < 26)
        atomicAdd(&(histo[alphabet position/4]), 1);
        i += stride;
   }
}
```

# Introduction to CUDA

## Questions?

Contact information

Andreas Axelsson

Email: andreas.axelsson@ju.se

Mobile: 0709-467760