# GP-GPU and High Performances Computing
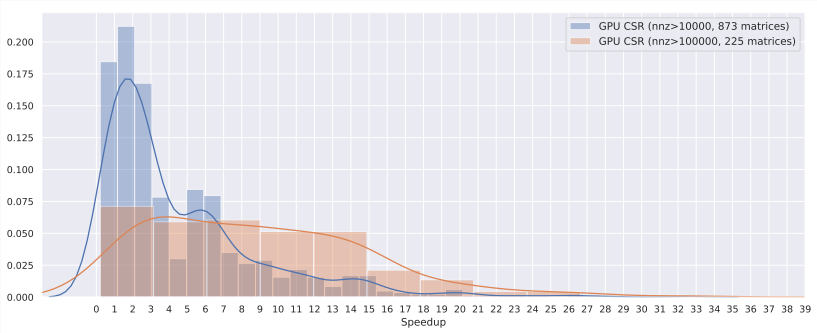
Lecture 10 – Histogram
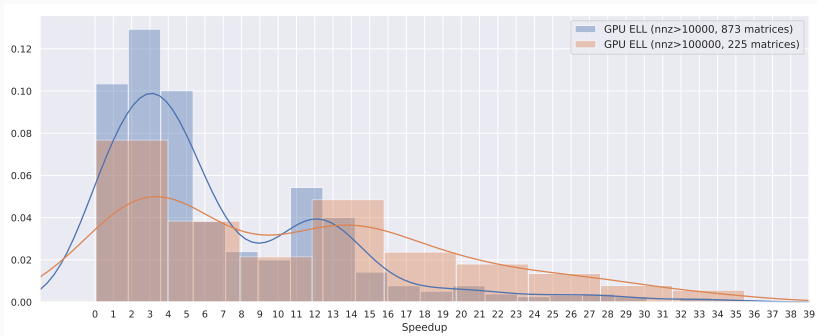
December 1, 2023
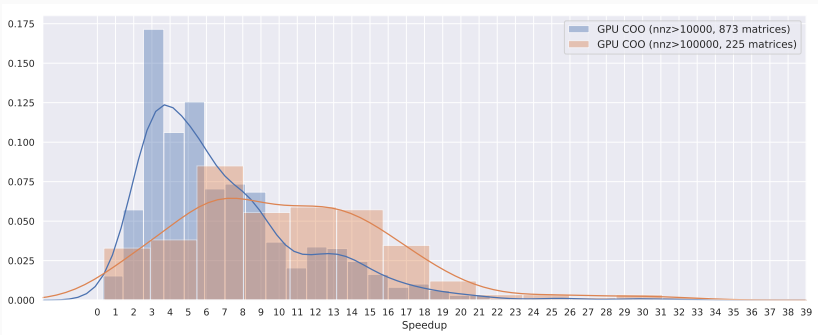
Organize storage of sparse matrices in order to

- ➤ minimize memory occupancy
- ➤ increase throughput
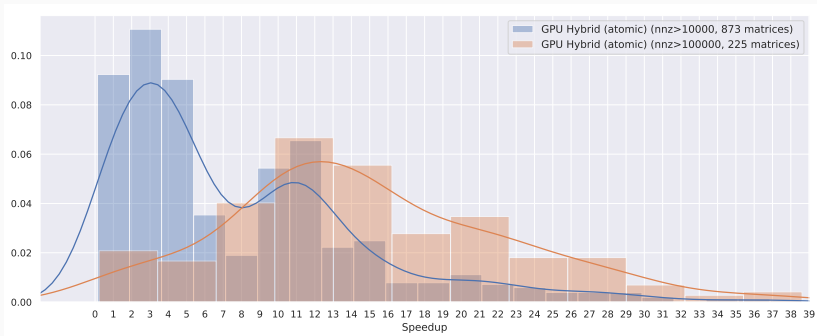- ➤ limit data duplication
- ➤ limit tasks duplication

# Performances comparison

```
1    reduce(int *g_idata, int *g_odata) {
2      extern __shared__ int sdata[];
3      // load shared mem
4      unsigned int tid = threadIdx.x;
5      unsigned int i = blockIdx.x * blockDim.x + threadIdx.x;
6      sdata[tid] = g_idata[i];
7      // do reduction in shared mem
8      for (unsigned int s = 1; s < blockDim.x; s *= 2) {
9        __syncthreads();
10       int index = 2 * s * tid;
11       if (index < blockDim.x) {
12         sdata[tid] = sdata[tid] + sdata[tid + s];
13       }
14       // Thread 0 writes result for this block to global mem
15       if (tid == 0) g_odata[blockIdx.x] = sdata[0];
16     }
17   }
```

```
1    reduce(int *g_idata, int *g_odata) {
2      extern __shared__ int sdata[];
3      // load shared mem
4      unsigned int tid = threadIdx.x;
5      unsigned int i = blockIdx.x * blockDim.x + threadIdx.x;
6      sdata[tid] = g_idata[i];
7      // do reduction in shared mem
8      for (int s = 1; s < blockDim.x; s *= 2) {
9        __syncthreads();
10       if (threadIdx.x % (2 * s) == 0)
11         sdata[threadID] += sdata[threadIdx.x + s];
12     }
13     // Thread 0 writes result for this block to global mem
14     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
15   }
```

Expected gain : $\times 2.5$

```
1    reduce(int *g_idata, int *g_odata) {
2      extern __shared__ int sdata[];
3      // load shared mem
4      unsigned int tid = threadIdx.x;
5      unsigned int i = blockIdx.x*blockDim.x + threadIdx.x;
6      sdata[tid] = g_idata[i];
7      __syncthreads();
8      // do reduction in shared mem
9      for (unsigned int s = blockDim.x/2; s > 0; s >>= 1) {
10       if (tid < s) {
11         sdata[tid] += sdata[tid + s];
12       }
13       __syncthreads();
14     }
15     // write result for this block to global mem
16     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
17   }
```

Expected gain : $\times 2$.

```
1    reduce(int *g_idata, int *g_odata) {
2      extern __shared__ int sdata[];
3      // load shared mem
4      unsigned int tid = threadIdx.x;
5      unsigned int i = blockIdx.x*(blockDim.x*2) + threadIdx.x;
6      sdata[tid] = g_idata[i] + g_idata[i+blockDim.x];
7      __syncthreads();
8      // do reduction in shared mem
9      for (unsigned int s = blockDim.x/2; s > 0; s >>= 1) {
10       if (tid < s) {
11         sdata[tid] += sdata[tid + s];
12       }
13       __syncthreads();
14     }
15     // write result for this block to global mem
16     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
17   }
```

Expected gain : $\times 1.8$

```
1    reduce(int *g_idata, int *g_odata) {
2      extern __shared__ int sdata[];
3      // load shared mem
4      unsigned int tid = threadIdx.x;
5      unsigned int i = blockIdx.x*(blockDim.x*2) + threadIdx.x;
6      sdata[tid] = g_idata[i] + g_idata[i+blockDim.x];
7      __syncthreads();
8      // do reduction in shared mem
9      for (unsigned int s = blockDim.x/2; s > 32; s >>= 1) {
10       if (tid < s) {
11         sdata[tid] += sdata[tid + s];
12       }
13       __syncthreads();
14     }
15     if (tid < 32) warpReduce(sdata, tid);
16     // write result for this block to global mem
17     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
18   }
```

Expected gain : $\times 1.8$

```
1    __device__ void warpReduce(volatile int* sdata, int tid) {
2      sdata[tid] += sdata[tid + 32];
3      sdata[tid] += sdata[tid + 16];
4      sdata[tid] += sdata[tid + 8];
5      sdata[tid] += sdata[tid + 4];
6      sdata[tid] += sdata[tid + 2];
7      sdata[tid] += sdata[tid + 1];
8    }
```

```
1    Template <unsigned int blockSize>
2    __device__ void warpReduce(volatile int* sdata, int tid) {
3      if (blockSize >= 64) sdata[tid] += sdata[tid + 32];
4      if (blockSize >= 32) sdata[tid] += sdata[tid + 16];
5      if (blockSize >= 16) sdata[tid] += sdata[tid + 8];
6      if (blockSize >= 8) sdata[tid] += sdata[tid + 4];
7      if (blockSize >= 4) sdata[tid] += sdata[tid + 2];
8      if (blockSize >= 2) sdata[tid] += sdata[tid + 1];
9    }
```

Expected gain : $\times 1.4$

# Reduction ↻: unroll warp

```
1    reduce(int *g_idata, int *g_odata) {
2      extern __shared__ int sdata[];
3      // load shared mem
4      unsigned int tid = threadIdx.x;
5      unsigned int i = blockIdx.x*(blockDim.x*2) + threadIdx.x;
6      unsigned int gridSize = blockSize*2*gridDim.x;
7      sdata[tid] = 0;
8      while (i < n) {
9        sdata[tid] += g_idata[i] + g_idata[i+blockSize];
10       i += gridSize;
11     }
12     __syncthreads();
13     // do reduction in shared mem
14     for (unsigned int s = blockDim.x/2; s > 32; s >>= 1) {
15       if (tid < s) {
16         sdata[tid] += sdata[tid + s];
17       }
18       __syncthreads();
19     }
20     if (tid < 32) warpReduce(sdata, tid);
21     // write result for this block to global mem
22     if (tid == 0) g_odata[blockIdx.x] = sdata[0];
23   }
```
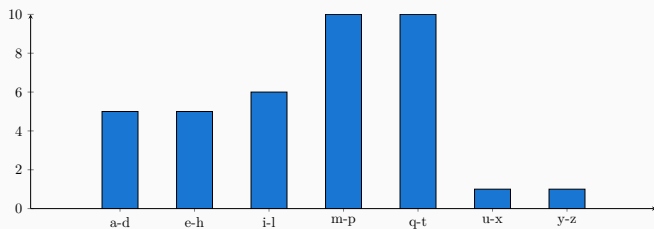
Expected gain : ×1.4

Histogram computation

The key techniques for compacting input data in parallel sparse methods for reduced consumption of memory bandwidth

➤ better utilization of on-chip memory

➤ fewer bytes transferred to on-chip memory

➤ retaining regularity

- 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

- ▸ 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.
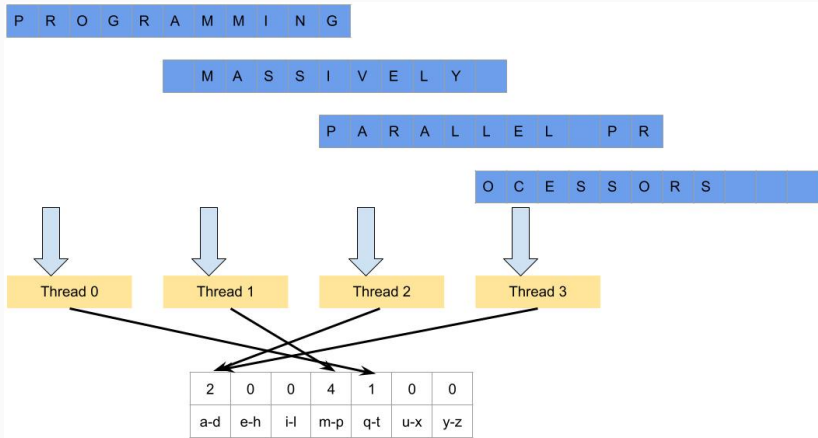- ▸ In the phrase "Programming Massively Parallel Processors" the output histogram is:

- Partition the input into sections
    - Have each thread to take a section of the input
    - Each thread iterates through its section.
    - For each letter, increment the appropriate bin counter

➤ 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 memory accesses

A reduction point of view

|  | a-d | e-h | i-l | m-p | q-t | u-x | y-z |
|---|---|---|---|---|---|---|---|
| P | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| O | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| G | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| A | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| M | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| M | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| I | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| N | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| G | 0 | 1 | 0 | 0 | 0 | 0 | 0 |

|   | a-d | e-h | i-l | m-p | q-t | u-x | y-z |
|---|-----|-----|-----|-----|-----|-----|-----|
| P | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| O | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| G | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| R | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| A | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| M | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| M | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| I | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| N | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| G | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
|   | 1 | 2 | 1 | 5 | 2 | 0 | 0 |

Data races

For example, multiple bank tellers count the total amount of cash in the safe

- ➤ Each grab a pile and count
- ➤ Have a central display of the running total
- ➤ Whenever someone finishes counting a pile, read the current running total (read) and add the subtotal of the pile to the running total (modify-write)
- ➤ A bad outcome

  Some of the piles were not accounted for in the final total

## A Common Parallel Service Pattern

For example, multiple customer service agents serving waiting customers

- ➤ The system maintains two numbers,
- ➤ the number to be given to the next incoming customer (I)
- ➤ the number for the customer to be served next (S)
- ➤ The system gives each incoming customer a number (read I) and increments the number to be given to the next customer by 1 modify (write I)
- ➤ A central display shows the number for the customer to be served next
- ➤ When an agent becomes available, he/she calls the number (read S) and increments the display number by 1 (modify-write S)
- ➤ Bad outcomes

    Multiple customers receive the same number, only one of them receives service

    Multiple agents serve the same number

For example, multiple customers booking airline tickets in parallel Each

- Brings up a flight seat map (read)
- Decides on a seat
- Updates the seat map and marks the selected seat as taken (modify-write)
- A bad outcome
    Multiple passengers ended up booking the same seat

Thread 0                          Thread 1

  Old = Mem[x]                      Old = Mem[x]
  New = Old + 1                     New = Old + 1
  Mem[x] = New                      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

Thread 0                                    Thread 1

```
1   Old = Mem[x] // 0              1
2   New = Old + 1 // 1             2
3   Mem[x] = New // 1              3
4                                  4   Old = Mem[x] // 1
5                                  5   New = Old + 1 // 2
6                                  6   Mem[x] = New // 2
```

- ▸ Thread 1 : `Old = 0`
- ▸ Thread 2 : `Old = 1`
- ▸ After the sequence : `Mem[x] = 2`

Thread 0

1
2
3
4     Old = Mem[x] // 1
5     New = Old + 1 // 2
6     Mem[x] = New // 2

Thread 1

1     Old = Mem[x] // 0
2     New = Old + 1 // 1
3     Mem[x] = New // 1
4
5
6

➤ Thread 1 : Old = 1

➤ Thread 2 : Old = 0

➤ After the sequence : Mem[x] = 2

Thread 0

```
1    Old = Mem[x] // 0
2    New = Old + 1 // 1
3
4    Mem[x] = New // 1
5
6
```

Thread 1

```
1
2
3    Old = Mem[x] // 1
4
5    New = Old + 1 // 1
6    Mem[x] = New // 1
```

- ➤ Thread 1 : $Old = 0$
- ➤ Thread 2 : $Old = 0$
- ➤ After the sequence : $Mem[x] = 1$

Thread 0

```
1
2
3    Old = Mem[x] // 1
4
5    New = Old + 1 // 1
6    Mem[x] = New // 1
```

Thread 1

```
1    Old = Mem[x] // 0
2    New = Old + 1 // 1
3
4    Mem[x] = New // 1
5
6
```

► Thread 1 : Old = 0
► Thread 2 : Old = 0
► After the sequence : Mem[x] = 1

## Atomic operation

The goal of atomic operation is to ensure that

Thread 0

1    Old = Mem[x] // 0
2    New = Old + 1 // 1
3    Mem[x] = New // 1
4
5
6

Thread 1

1
2
3
4    Old = Mem[x] // 1
5    New = Old + 1 // 2
6    Mem[x] = New // 2

or

Thread 0

1
2
3
4    Old = Mem[x] // 1
5    New = Old + 1 // 2
6    Mem[x] = New // 2

Thread 1

1    Old = Mem[x] // 0
2    New = Old + 1 // 1
3    Mem[x] = New // 1
4
5
6

# Atomic operation

```
Mem[x]= 0
Thread 0                              Thread 1
1   Old = Mem[x] // 0               1
2   New = Old + 1 // 1              2
3                                   3   Old = Mem[x] // 1
4   Mem[x] = New // 1               4
5                                   5   New = Old + 1 // 1
6                                   6   Mem[x] = New // 1
```

➤ Both threads receive 0 in Old

➤ Mem[x] becomes 1

➤ 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

➤ Performed by calling functions that are translated into single instructions:
add, sub, inc, dec, min, max, exch (exchange), CAS (compare and swap)

➤ 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. The function returns old.

- Unsigned 32-bit integer atomic add

  ```
  unsigned int atomicAdd(unsigned int* address, unsigne
  ```

- Unsigned 64-bit integer atomic add

  ```
  unsigned long long int atomicAdd(unsigned long long i
  ```

- Single-precision floating-point atomic add (capability > 2.0)

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

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

```
1  __global__ void histo_kernel(unsigned char *buffer, long size,
2                                unsigned int *histo)
3  {
4    int i = threadIdx.x + blockIdx.x * blockDim.x;
5    // stride is total number of threads
6    int stride = blockDim.x * gridDim.x;
7    // All threads handle blockDim.x * gridDim.x
8    // consecutive elements
9    while (i < size) {
10     int alphabet_position = buffer[i] - "a";
11     if (alphabet_position >= 0 && alpha_position < 26)
12       atomicAdd(&(histo[alphabet_position/4]), 1);
13     i += stride;
14   }
15 }
```

# Performances

- An atomic operation on a DRAM location starts with a read, which has a latency of a few hundred cycles
- The atomic operation ends with a write to the same location, with a latency of a few hundred cycles
- During this whole time, no one else can access the location

Each Read-Modify-Write has two full memory access delays

All atomic operations on the same variable (DRAM location) are serialized.

- Throughput of atomic operations on the same DRAM location is the rate at which the application can execute an atomic operation.
- The rate for atomic operation on a particular location is limited by the total latency of the read-modify-write sequence, typically more than 1000 cycles for global memory (DRAM) locations.
- This means that if many threads attempt to do atomic operation on the same location (contention), the memory throughput is reduced to < 1/1000 of the peak bandwidth of one memory channel!

1. Some customers realize that they missed an item after they started to check out
2. They run to the isle and get the item while the line waits:
   The rate of checkout is drastically reduced due to the long latency of running to the isle and back.
3. Imagine a store where every customer starts the check out before they even fetch any of the items:
   The rate of the checkout will be 1 / (entire shopping time of each customer)

Atomic operations on Shared Memory

➤ Very short latency
➤ Private to each thread block
➤ Need algorithm work by programmers

Avoid atomic operations as much as possible.

# Privatization

- ➤ Cost
  - ➤ Overhead for creating and initializing private copies
  - ➤ Overhead for accumulating the contents of private copies into the final copy
- ➤ Benefit
  - ➤ Much less contention and serialization in accessing both the private copies and the final copy
  - ➤ The overall performance can often be improved more than 10

➤ Each subset of threads are in the same block

➤ Much higher throughput than DRAM (100x) or L2 (10x) atomics

➤ Less contention – only threads in the same block can access a shared memory variable

➤ This is a very important use case for shared memory!

# Kernel with privatization

```
1   __global__ void histo_kernel(unsigned char *buffer, long size,
2             unsigned int *histo)
3   {
4     /// Create private copies of the histo[] array for each thread block
5     __shared__ unsigned int histo_private[7];
6     // Initialize the bin counters in the private copies of histo[]
7     if (threadIdx.x < 7)
8       histo_private[threadidx.x] = 0;
9     __syncthreads();
10
11    /// Build Private Histogram
12    int i = threadIdx.x + blockIdx.x * blockDim.x;
13    // stride is total number of threads
14    int stride = blockDim.x * gridDim.x;
15    while (i < size) {
16      atomicAdd( &(private_histo[buffer[i]/4), 1);
17      i += stride;
18    }
19    /// Build Final Histogram
20    // wait for all other threads in the block to finish
21    __syncthreads();
22    if (threadIdx.x < 7) {
23      atomicAdd(&(histo[threadIdx.x]), private_histo[threadIdx.x] );
24    }
25  }
```

- ➤ Privatization is a powerful and frequently used technique for parallelizing applications
- ➤ The operation needs to be associative and commutative
- ➤ Histogram add operation is associative and commutative
- ➤ No privatization if the operation does not fit the requirement
- ➤ The private histogram size needs to be small
- ➤ Fits into shared memory
- ➤ If the histogram is too large to privatize: partially privatize an output histogram and use range testing to go to either global memory or shared memory.

# Conclusion

- Themes of this class
    - Memory access
    - Race condition
    - Atomic operation
    - Use private memory