

CUDA Parallel Programming

Chen Chen

October 16, 2025

Contents

Scan

Reduction

Installing GPGPU-Sim v4: Prerequisites (1/3)

Step 1: Download Sources & CUDA Installer

Download and Prepare Files

Get all necessary files before installation:

```
1 # 1. Download CUDA 11.0.2 installer
2 wget http://developer.download.nvidia.com/compute/cuda/11.0.2/\
3 local_installers/cuda_11.0.2_450.51.05_linux.run
4
5 # 2. Clone GPGPU-Sim v4 Distribution
6 git clone https://github.com/accel-sim/gpgpu-sim_distribution
7
8 # 3. Rename the directory for convenience
9 mv gpgpu-sim_distribution gpgpusim-v4
```

Installing GPGPU-Sim v4: Prerequisites (2/3)

Step 2: Adjust GCC/G++ Version to 9

Configure Compiler Compatibility

CUDA 11.0 requires GCC/G++ 9 for compatibility:

```
1 sudo apt install gcc-9 g++-9
2
3 # Set GCC/G++ version 9 as default
4 sudo update-alternatives --install /usr/bin/gcc gcc /usr/bin/gcc-9 100
5 sudo update-alternatives --install /usr/bin/g++ g++ /usr/bin/g++-9 100
```

Installing GPGPU-Sim v4: CUDA Setup (2/3)

CUDA Installation and Environment Variables

Step 3: Install CUDA 11.0 and Set Paths

Execute the CUDA installer (deselect the driver installation option):

```
1 sudo sh cuda_11.0.2_450.51.05_linux.run
```

Add CUDA paths to your `~/.bashrc`:

```
1 echo 'export PATH=$PATH:/usr/local/cuda/bin' >> ~/.bashrc
2 echo 'export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:/usr/local/cuda/lib:/usr\
3 /local/cuda/lib64' >> ~/.bashrc
4 echo 'export CUDA_INSTALL_PATH=/usr/local/cuda' >> ~/.bashrc
5 source ~/.bashrc
6
7 # Verify installation
8 nvcc -V
9 # Check if output shows: V11.0.194
```

Installing GPGPU-Sim v4: Final Steps (3/3)

Dependencies and Compilation

Step 4: Install OS Dependencies

Install essential build tools and libraries:

```
1 sudo apt-get install -y wget build-essential xutils-dev bison\  
2 zlib1g-dev flex libglu1-mesa-dev git g++ libssl-dev libxml2-dev \  
3 libboost-all-dev vim python3-pip libxi-dev libxmu-dev libglut3-dev
```

Step 5: Compile GPGPU-Sim v4

Load environment variables and compile the simulator:

```
1 cd gpgpusim-v4  
2 source setup_environment  
3 make -j$(nproc)
```

Scan

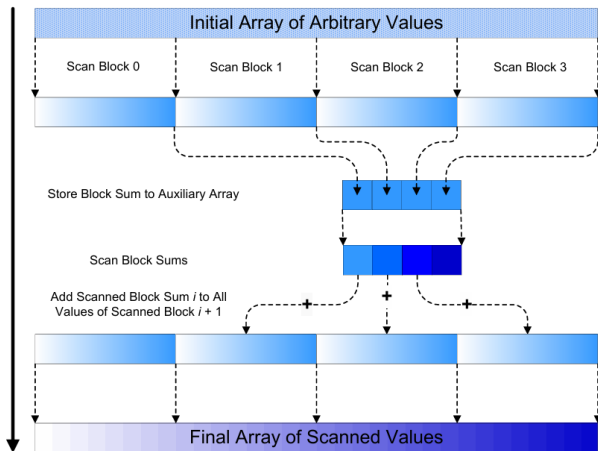


Figure: Algorithm for performing a sum scan on a large array of values.

Header and Definitions

```
1  #include <cuda_runtime.h>
2  #include <iostream>
3  #include <vector>
4  #include <chrono>
5  using namespace std;
6
7  #define MAX_THREADS_PER_BLOCK 1024
8  #define MAX_ELEMENTS_PER_BLOCK (MAX_THREADS_PER_BLOCK * 2)
```


Kernel: Data Load

```
1  __global__ void parallel_large_scan_kernel(int *data, int
   ↪ *prefix_sum, int N, int *sums)
2  {
3      __shared__ int tmp[MAX_ELEMENTS_PER_BLOCK];
4      int tid = threadIdx.x;
5      int bid = blockIdx.x;
6      int block_offset = bid * MAX_ELEMENTS_PER_BLOCK;
7      int leaf_num = MAX_ELEMENTS_PER_BLOCK;
8
9      tmp[tid * 2]      = (tid * 2 + block_offset < N) ? data[tid *
   ↪ 2 + block_offset] : 0;
10     tmp[tid * 2 + 1] = (tid * 2 + 1 + block_offset < N) ?
   ↪ data[tid * 2 + 1 + block_offset] : 0;
11     __syncthreads();
```

Upsweep Stage

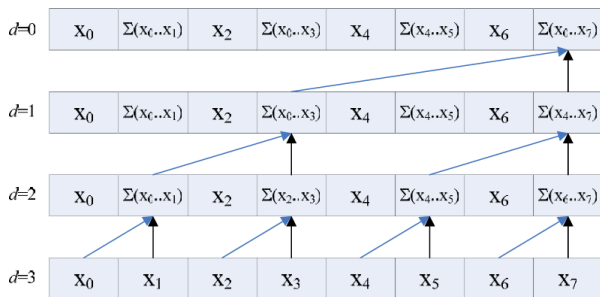


Figure: An illustration of the up-sweep phase of a work-efficient sum scan algorithm.

Kernel: Upsweep Stage

```
1  int offset = 1;
2  for (int d = leaf_num >> 1; d > 0; d >>= 1)
3  {
4      if (tid < d) {
5          int ai = offset * (2 * tid + 1) - 1;
6          int bi = offset * (2 * tid + 2) - 1;
7          tmp[bi] += tmp[ai];
8      }
9      offset *= 2;
10     __syncthreads();
11 }
```

Downsweep Stage

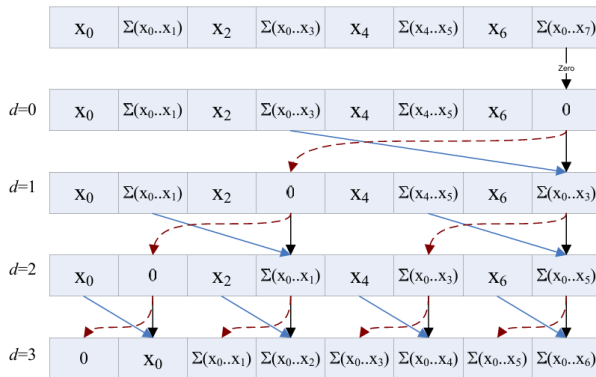


Figure: An illustration of the down-sweep phase of the work efficient parallel sum scan algorithm. Notice that the first step zeros the last element of the array.

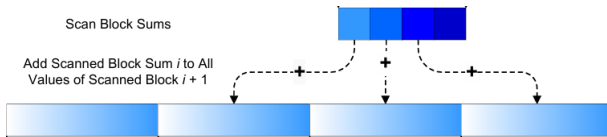
Kernel: Exclusive Downsweep

```
1  if (tid == 0) {
2      sums[bid] = tmp[leaf_num - 1];
3      tmp[leaf_num - 1] = 0;
4  }
5  __syncthreads();
6
7  for (int d = 1; d < leaf_num; d *= 2)
8  {
9      offset >>= 1;
10     if (tid < d) {
11         int ai = offset * (2 * tid + 1) - 1;
12         int bi = offset * (2 * tid + 2) - 1;
13         int t = tmp[ai];
14         tmp[ai] = tmp[bi];
15         tmp[bi] += t;
16     }
17     __syncthreads();
18 }
19 }
```

Write Back

```
1  if (tid * 2 + block_offset < N)
2      prefix_sum[tid * 2 + block_offset] = tmp[tid * 2];
3  if (tid * 2 + 1 + block_offset < N)
4      prefix_sum[tid * 2 + 1 + block_offset] = tmp[tid * 2 +
    ↪ 1];
```

Add Kernel



```
1  __global__ void add_kernel(int *prefix_sum, const int *values,
   ↪  int N)
2  {
3      int tid = threadIdx.x;
4      int bid = blockIdx.x;
5      int block_offset = bid * MAX_ELEMENTS_PER_BLOCK;
6      int ai = tid * 2 + block_offset;
7      int bi = tid * 2 + 1 + block_offset;
8
9      if (ai < N)
10         prefix_sum[ai] += values[bid];
11     if (bi < N)
12         prefix_sum[bi] += values[bid];
13 }
```

Recursive Scan

```
1 void recursive_scan(int *d_data, int *d_prefix_sum, int N)
2 {
3     int block_num = (N + MAX_ELEMENTS_PER_BLOCK - 1) /
4         ↪ MAX_ELEMENTS_PER_BLOCK;
5     int *d_sums = nullptr, *d_sums_prefix = nullptr;
6     cudaMalloc(&d_sums, block_num * sizeof(int));
7     cudaMalloc(&d_sums_prefix, block_num * sizeof(int));
8
9     parallel_large_scan_kernel<<<block_num,
10         ↪ MAX_THREADS_PER_BLOCK>>>(d_data, d_prefix_sum, N,
11         ↪ d_sums);
12     if (block_num > 1) {
13         recursive_scan(d_sums, d_sums_prefix, block_num);
14         add_kernel<<<block_num,
15             ↪ MAX_THREADS_PER_BLOCK>>>(d_prefix_sum, d_sums_prefix,
16             ↪ N);
17     }
18     cudaFree(d_sums);
19     cudaFree(d_sums_prefix);
20 }
```


Performance Comparison

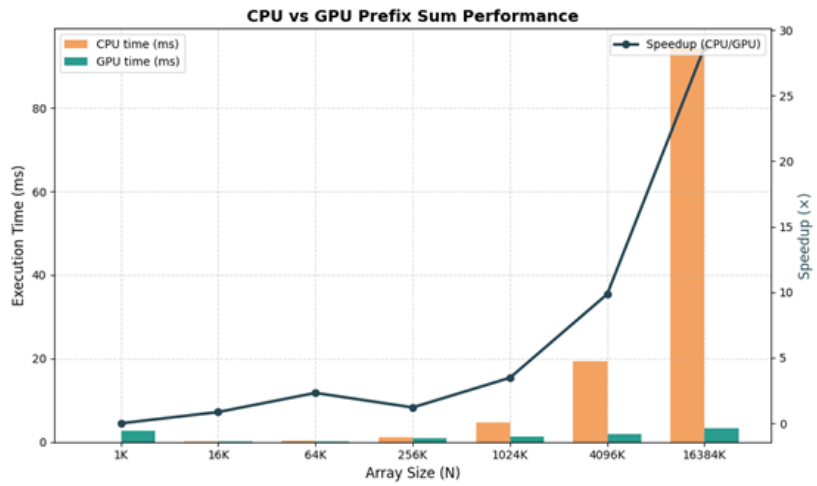


Figure: Performance Comparison of Parallel Scan and Serial Scan at Different Scale.

Parallel Reduction

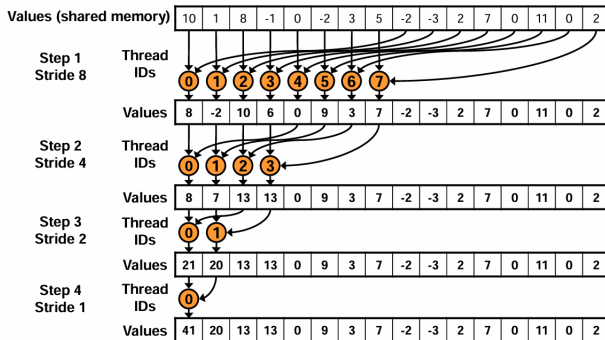


Figure: Sequential addressing is conflict free.

Reduction: Basic Version

Simple In-Block Reduction Kernel

```
1  __global__ void reduce1(int* iData, int* oData, int n) {
2      extern __shared__ int temp[];
3      int thid = threadIdx.x;
4      int gid  = blockIdx.x * blockDim.x + threadIdx.x;
5
6      // Load input data into shared memory
7      temp[thid] = (gid < n) ? iData[gid] : 0;
8      __syncthreads();
9
10     // Perform reduction in shared memory
11     for (int s = blockDim.x >> 1; s > 0; s >>= 1) {
12         if (thid < s) {
13             temp[thid] += temp[thid + s];
14         }
15         __syncthreads();
16     }
17
18     // Write result for this block to global memory
19     if (thid == 0)
20         oData[blockIdx.x] = temp[0];
21 }
```

Warp Reduce

When $s \leq 32$, we have only one warp left. Instructions are SIMD synchronous within a warp. That means when $s \leq 32$: We don't need to `__syncthreads()`

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

Reduction

```
1  template <unsigned int blockSize>
2  __global__ void reduce2(int* iData, int* oData, int n) {
3      extern __shared__ int temp[];
4      unsigned int thid = threadIdx.x;
5      unsigned int i = blockIdx.x * (blockSize * 2) + threadIdx.x;
6      unsigned int gridSize = blockSize * 2 * gridDim.x;
7      temp[thid] = 0;
8
9      while (i < n) {
10         int a = iData[i];
11         int b = (i + blockSize < n) ? iData[i + blockSize] : 0;
12         temp[thid] += a + b;
13         i += gridSize;
14     }
```

Reduction

```
1  __syncthreads();
2
3  if (blockSize >= 512) { if (thid < 256) temp[thid] +=
    ↪ temp[thid + 256]; __syncthreads(); }
4  if (blockSize >= 256) { if (thid < 128) temp[thid] +=
    ↪ temp[thid + 128]; __syncthreads(); }
5  if (blockSize >= 128) { if (thid < 64) temp[thid] +=
    ↪ temp[thid + 64]; __syncthreads(); }
6
7  if (thid < 32) warpReduce<blockSize>(temp, thid);
8  if (thid == 0) oData[blockIdx.x] = temp[0];
9 }
```

Performance Comparison

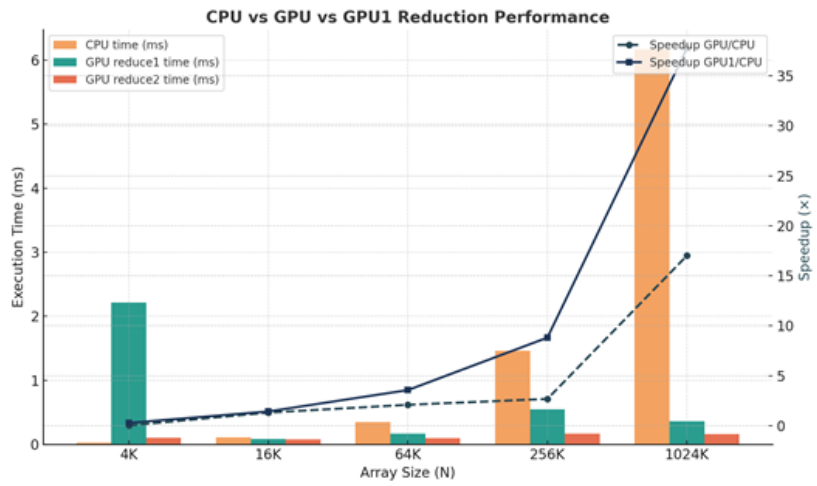


Figure: Performance Comparison of Parallel Reduction and Serial Reduction at Different Scale.