

CUDA Assignment Week 2

Complete Solutions

Krishna Kukreja

GPU Programming using CUDA and Triton (WiDS'25)

December 2025

Contents

1	Introduction	3
1.1	Assignment Overview	3
1.2	Key Concepts	3
1.2.1	Thread Indexing	3
1.2.2	Grid Size Calculation	3
2	Task 1: Basic CUDA Kernel Implementation	3
2.1	Kernel 1: Vector Addition	3
2.1.1	Mathematical Formulation	3
2.1.2	CUDA Implementation	3
2.1.3	Complexity Analysis	4
2.2	Kernel 2: Element-wise Multiply and Scale	4
2.2.1	Mathematical Formulation	4
2.3	Kernel 3: ReLU Activation	4
2.3.1	Mathematical Formulation	4
3	Host Code Implementation	4
3.1	Memory Allocation Pattern	4
3.2	Kernel Launch	5
4	Task 2: Grid and Block Configuration Analysis	5
4.1	Configuration Results	5
4.2	Answer to Q1: Total Threads Exceeding n	6
4.3	Answer to Q2: Why Bounds Checking Required	6
4.4	Answer to Q3: Failed Configurations	6
4.5	Answer to Q4: Block Size Selection Trade-offs	7
4.5.1	Small Block Size (32-64 threads)	7
4.5.2	Large Block Size (256-512 threads)	7
4.5.3	Recommendation	7
5	Task 3: Timing and Performance Analysis	7
5.1	Experimental Setup	7
5.2	Bandwidth Calculation	8
5.3	Observed Results	8
5.4	Analysis	8
6	Compilation and Execution	8
6.1	Compilation	8
6.2	Execution	9
7	Correctness Verification	9
7.1	Verification Strategy	9
7.2	Error Tolerance	9
8	Summary and Key Takeaways	10
8.1	Critical Learning Points	10

1 Introduction

This document provides complete solutions for CUDA Assignment Week 2: Parallel Thinking and CUDA Programming Basics.

1.1 Assignment Overview

The assignment consists of three main tasks:

1. Implement three basic CUDA kernels (vector addition, multiply-scale, ReLU)
2. Analyze grid and block configuration trade-offs
3. Measure and analyze kernel performance timing

1.2 Key Concepts

1.2.1 Thread Indexing

The global thread ID calculation:

$$\text{idx} = \text{blockIdx.x} \times \text{blockDim.x} + \text{threadIdx.x} \quad (1)$$

1.2.2 Grid Size Calculation

For n elements with block size B :

$$\text{gridSize} = \left\lceil \frac{n}{B} \right\rceil = \frac{n + B - 1}{B} \quad (2)$$

2 Task 1: Basic CUDA Kernel Implementation

2.1 Kernel 1: Vector Addition

2.1.1 Mathematical Formulation

$$c[i] = a[i] + b[i], \quad \forall i \in [0, n) \quad (3)$$

2.1.2 CUDA Implementation

Listing 1: Vector Addition Kernel

```

1  __global__ void vectorAdd(const float *a, const float *b,
2                                float *c, int n) {
3      int idx = blockIdx.x * blockDim.x + threadIdx.x;
4
5      if (idx < n) {
6          c[idx] = a[idx] + b[idx];
7      }
8 }
```

2.1.3 Complexity Analysis

- Time Complexity: $O(n)$
- Space Complexity: $O(n)$ device memory
- Memory Bandwidth: 12 bytes per operation
- Arithmetic Intensity: 0.083 FLOP/byte
- Classification: Memory-bound operation

2.2 Kernel 2: Element-wise Multiply and Scale

2.2.1 Mathematical Formulation

$$c[i] = \alpha \cdot a[i] \cdot b[i] \quad (4)$$

Listing 2: Multiply-Scale Kernel

```

1 --global__ void multiplyScale(const float *a, const float *b,
2                                           float *c, float alpha, int n) {
3     int idx = blockIdx.x * blockDim.x + threadIdx.x;
4
5     if (idx < n) {
6         c[idx] = alpha * a[idx] * b[idx];
7     }
8 }
```

2.3 Kernel 3: ReLU Activation

2.3.1 Mathematical Formulation

$$y[i] = \max(x[i], 0) \quad (5)$$

Listing 3: ReLU Activation Kernel

```

1 --global__ void relu(const float *x, float *y, int n) {
2     int idx = blockIdx.x * blockDim.x + threadIdx.x;
3
4     if (idx < n) {
5         y[idx] = fmaxf(x[idx], 0.0f);
6     }
7 }
```

3 Host Code Implementation

3.1 Memory Allocation Pattern

Listing 4: Memory Allocation and Transfer

```

1 // Host memory allocation
2 float *h_a = (float *)malloc(bytes);
3 float *h_b = (float *)malloc(bytes);
4 float *h_c = (float *)malloc(bytes);
5
6 // Device memory allocation
7 float *d_a, *d_b, *d_c;
8 cudaMalloc((void**)&d_a, bytes);
9 cudaMalloc((void**)&d_b, bytes);
10 cudaMalloc((void**)&d_c, bytes);
11
12 // Copy to device
13 cudaMemcpy(d_a, h_a, bytes, cudaMemcpyHostToDevice);
14 cudaMemcpy(d_b, h_b, bytes, cudaMemcpyHostToDevice);

```

3.2 Kernel Launch

Listing 5: Kernel Launch with Configuration

```

1 int blockSize = 256;
2 int gridSize = (n + blockSize - 1) / blockSize;
3
4 vectorAdd<<<gridSize, blockSize>>>(d_a, d_b, d_c, n);
5
6 cudaDeviceSynchronize();
7
8 cudaMemcpy(h_c, d_c, bytes, cudaMemcpyDeviceToHost);

```

4 Task 2: Grid and Block Configuration Analysis

4.1 Configuration Results

Table 1: Grid and Block Configuration Analysis

Vector Size	Block	Grid	Total	Excess	Status
10^3	32	32	1,024	24	Check needed
10^3	128	8	1,024	24	Check needed
10^3	256	4	1,024	24	Check needed
10^3	512	2	1,024	24	Check needed
10^5	32	3,125	100,000	0	Perfect
10^5	128	782	100,096	96	Check needed
10^5	256	391	100,096	96	Check needed
10^5	512	196	100,352	352	Check needed
10^7	32	312,500	10M	0	Perfect
10^7	128	78,125	10M	0	Perfect
10^7	256	39,063	10M	128	Check needed
10^7	512	19,532	10M	384	Check needed

4.2 Answer to Q1: Total Threads Exceeding n

With Bounds Checking: If the total number of threads exceeds array size n , the excess threads execute the bounds checking condition:

```

1 if (idx < n) {
2     // Excess threads skip this block
3 }
```

Result: Excess threads return early without performing computation. This causes warp divergence but maintains correctness.

Without Bounds Checking: Excess threads access out-of-bounds memory locations causing:

1. Segmentation fault (if memory unallocated)
2. Silent data corruption (if memory allocated to other data)
3. Undefined program behavior

4.3 Answer to Q2: Why Bounds Checking Required

Bounds checking is essential because:

1. **No Hardware Protection:** GPUs lack Memory Management Units (MMUs) with virtual memory protection that CPUs provide
2. **Flat Memory Space:** GPU global memory is contiguous; out-of-bounds writes overwrite adjacent data structures
3. **Silent Corruption:** Invalid memory access doesn't trigger exceptions; data gets corrupted silently
4. **Programmer Responsibility:** CUDA places safety burden on programmer, not hardware

4.4 Answer to Q3: Failed Configurations

Correctness Perspective: All configurations produce correct results when bounds checking is properly implemented.

Performance Perspective: Configurations with excess threads suffer warp divergence penalty:

Table 2: Warp Divergence Analysis

Config	Excess	Divergence Rate	Impact
$n = 10^3, B = 256$	24	9.4%	Minor
$n = 10^5, B = 512$	352	68.8%	Significant
$n = 10^7, B = 256$	128	0.12%	Negligible

The $n = 10^5, B = 512$ case shows significant divergence because 352 excess threads represent most of a warp.

4.5 Answer to Q4: Block Size Selection Trade-offs

4.5.1 Small Block Size (32-64 threads)

Advantages:

- Minimal wasted threads for small vectors
- Lower register pressure per block

Disadvantages:

- More blocks to schedule (overhead)
- Lower occupancy (fewer warps per SM)
- Poor cache locality

4.5.2 Large Block Size (256-512 threads)

Advantages:

- Better SM occupancy
- Fewer blocks (less scheduling overhead)
- Better latency hiding
- Improved cache efficiency

Disadvantages:

- More wasted threads for small vectors
- Higher register pressure (may limit occupancy)

4.5.3 Recommendation

Default choice: **128-256 threads per block**

- Balanced for most workloads
- Good occupancy without excessive waste
- Adjustable based on profiling

5 Task 3: Timing and Performance Analysis

5.1 Experimental Setup

The experiments were conducted on Google Colab using an ****NVIDIA T4 GPU****.

- **Input Size (n):** 10^7 elements (40 MB per vector)
- **Block Size:** 256 threads
- **Grid Size:** 39,063 blocks

5.2 Bandwidth Calculation

The theoretical peak bandwidth of the NVIDIA T4 is approximately **320 GB/s**. The effective bandwidth is calculated as:

$$\text{Bandwidth} = \frac{\text{Total Bytes Transferred}}{\text{Execution Time}} \quad (6)$$

For the **Multiply & Scale** kernel:

- Total Data: $3 \times 10^7 \times 4$ bytes = 120 MB
- Execution Time: 0.562 ms
- Measured Bandwidth: 213.64 GB/s

5.3 Observed Results

Table 3: Performance Results on NVIDIA T4 (Google Colab)

Kernel	Time (ms)	Bandwidth (GB/s)	Status
Vector Addition	≈ 0.56	≈ 213.6	Verified
Multiply-Scale	0.562	213.64	Verified
ReLU Activation	0.434	184.42	Verified

5.4 Analysis

The results show that the kernels achieve approximately **66% of the theoretical peak bandwidth** (213 GB/s vs 320 GB/s). This confirms that these element-wise operations are **memory-bound**; the GPU's compute units are waiting for data to arrive from global memory. The ReLU kernel has slightly lower effective bandwidth, likely due to the overhead of launching kernels relative to the smaller total data volume (80 MB vs 120 MB).

6 Compilation and Execution

6.1 Compilation

Listing 6: CUDA Compilation Commands

```

1 # Basic compilation
2 nvcc -o vector_add vector_add.cu -lm
3
4 # With optimizations
5 nvcc -O3 -o vector_add vector_add.cu -lm
6
7 # With debug symbols
8 nvcc -G -g -O0 -o vector_add vector_add.cu -lm
9
10 # Using Makefile
11 make all

```

6.2 Execution

Listing 7: Execution Commands

```

1 # Run with default size (10M elements)
2 ./vector_add
3
4 # Run with custom size
5 ./vector_add 1000000
6
7 # Run all kernels
8 make run
9
10 # Comprehensive tests
11 make test_all

```

7 Correctness Verification

7.1 Verification Strategy

Listing 8: Verification Against CPU Reference

```

1 // Compute CPU reference
2 for (int i = 0; i < n; i++) {
3     h_c_cpu[i] = h_a[i] + h_b[i];
4 }
5
6 // Compare GPU results
7 float maxError = 0.0f;
8 int errorCount = 0;
9 const float tolerance = 1e-5f;
10
11 for (int i = 0; i < n; i++) {
12     float error = fabsf(h_c[i] - h_c_cpu[i]);
13     if (error > maxError) maxError = error;
14     if (error > tolerance) errorCount++;
15 }
16
17 if (errorCount == 0) {
18     printf("VERIFICATION PASSED\n");
19 }

```

7.2 Error Tolerance

For single-precision floating-point (float):

- Machine epsilon: $\epsilon_m = 1.19 \times 10^{-7}$
- Recommended tolerance: $10 \times \epsilon_m = 1.19 \times 10^{-6}$
- Used in code: 1×10^{-5} (conservative)

8 Summary and Key Takeaways

8.1 Critical Learning Points

1. Bounds Checking is Non-Negotiable

- Single line of code: `if (idx < n)`
- Prevents entire class of memory safety bugs
- Negligible performance cost

2. Memory is the Bottleneck

- Vector operations are memory-bound
- Only 15-20% GPU bandwidth utilized
- Optimization opportunities come from memory, not computation

3. Grid Configuration Matters

- Use ceiling division: $(n + B - 1)/B$
- Excess threads cause warp divergence
- 128-256 is good default block size

4. Always Verify Results

- Compare with CPU reference implementation
- Test with multiple input sizes
- Check numerical accuracy with tolerance