



**COURSE CODE: DJ19DSL802**

**YEAR: 2024-2025**

**COURSE NAME: High Performance Computing Laboratory**

**CLASS: B.Tech (D2)**

**Experiment 6**

**Aim: Accelerating a For Loop with a Mismatched Execution Configuration**

**Theory:**

It may be the case that an execution configuration cannot be expressed that will create the exact number of threads needed for parallelizing a loop.

A common example has to do with the desire to choose optimal block sizes. For example, due to GPU hardware traits, blocks that contain a number of threads that are a multiple of 32 are often desirable for performance benefits. Assuming that we wanted to launch blocks each containing 256 threads (a multiple of 32), and needed to run 1000 parallel tasks (a trivially small number for ease of explanation), then there is no number of blocks that would produce an exact total of 1000 threads in the grid, since there is no integer value 32 can be multiplied by to equal exactly 1000.

This scenario can be easily addressed in the following way:

Write an execution configuration that creates more threads than necessary to perform the allotted work.

Pass a value as an argument into the kernel (N) that represents to the total size of the data set to be processed, or the total threads that are needed to complete the work.

After calculating the thread's index within the grid (using  $tid + bid * blockDim$ ), check that this index does not exceed N, and only perform the pertinent work of the kernel if it does not.

Here is an example of an idiomatic way to write an execution configuration when both N and the number of threads in a block are known, and an exact match between the number of threads in the grid and N cannot be guaranteed. It ensures that there are always at least as many threads as needed for N, and only 1 additional block's worth of threads extra, at most:

```
// Assume `N` is known
```

```
int N = 100000;
```

```
// Assume we have a desire to set `threads_per_block` exactly to `256`
```

```
size_t threads_per_block = 256;
```

```
// Ensure there are at least `N` threads in the grid, but only 1 block's worth extra
```

```
size_t number_of_blocks = (N + threads_per_block - 1) / threads_per_block;
```

```
some_kernel<<<number_of_blocks, threads_per_block>>>(N);
```



Shri Vile Parle Kelavani Mandal's

**DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING**

(Autonomous College Affiliated to the University of Mumbai)

NAAC Accredited with "A" Grade (CGPA : 3.18)



**Department of Computer Science and Engineering (Data Science)**

Because the execution configuration above results in more threads in the grid than N, care will need to be taken inside of the some\_kernel definition so that some\_kernel does not attempt to access out of range data elements, when being executed by one of the "extra" threads:

```
_global__some_kernel(int N)
{
    int idx = threadIdx.x + blockIdx.x * blockDim.x;

    if (idx < N) // Check to make sure `idx` maps to some value within `N`
    {
        // Only do work if it does
    }
}
```

Lab Experiment to be performed:

The program in 02-mismatched-config-loop.cu allocates memory, using cudaMallocManaged for a 1000 element array of integers, and then seeks to initialize all the values of the array in parallel using a CUDA kernel. This program assumes that both N and the number of threads\_per\_block are known. Your task is to complete the following two objectives, refer to the solution if you get stuck:

1. Assign a value to number\_of\_blocks that will make sure there are at least as many threads as there are elements in a to work on.
2. Update the initializeElementsTo kernel to make sure that it does not attempt to work on data elements that are out of range.



Shri Vile Parle Kelavani Mandal's

**DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING**

(Autonomous College Affiliated to the University of Mumbai)

NAAC Accredited with "A" Grade (CGPA : 3.18)



**Department of Computer Science and Engineering (Data Science)**

✓  
0s

[1] `!nvcc --version`

nvcc: NVIDIA (R) Cuda compiler driver  
Copyright (c) 2005-2023 NVIDIA Corporation  
Built on Tue\_Aug\_15\_22:02:13\_PDT\_2023  
Cuda compilation tools, release 12.2, V12.2.140  
Build cuda\_12.2.r12.2/compiler.33191640\_0

✓  
5s

[2] `!pip install nvcc4jupyter`

Collecting nvcc4jupyter  
  Downloading nvcc4jupyter-1.2.1-py3-none-any.whl.metadata (5.1 kB)  
  Downloading nvcc4jupyter-1.2.1-py3-none-any.whl (10 kB)  
Installing collected packages: nvcc4jupyter  
Successfully installed nvcc4jupyter-1.2.1

✓  
0s

[3] `%load_ext nvcc4jupyter`

Detected platform "Colab". Running its setup...  
Source files will be saved in "/tmp/tmpypb91x7c".



## Lab 6

```
1s 1s %%cuda

#include <stdio.h>

__global__ void initializeElementsTo(int initialValue, int *a, int N)
{
    int i = threadIdx.x + blockIdx.x * blockDim.x;
    if (i < N) a[i] = initialValue;
}

int main()
{
    int N = 1000;

    int *a;
    size_t size = N * sizeof(int);

    cudaMallocManaged(&a, size);

    size_t threads_per_block = 256;

    size_t number_of_blocks = 4;

    int initialValue = 6;

    initializeElementsTo<<<number_of_blocks, threads_per_block>>>(initialValue, a, N);
    cudaDeviceSynchronize();

    for (int i = 0; i < N; ++i)
    {
        if(a[i] != initialValue)
        {
            printf("FAILURE: target value: %d\t a[%d]: %d\n", initialValue, i, a[i]);
            cudaFree(a);
            exit(1);
        }
    }
    printf("SUCCESS!\n");

    cudaFree(a);
}
```

➡ SUCCESS!

### Conclusion:

In this experiment, we explored how to accelerate a for-loop using a mismatched execution configuration in CUDA. Since the number of required threads (N) may not always be an exact multiple of the chosen block size, we ensured that the execution configuration launched enough



Shri Vile Parle Kelavani Mandal's

**DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING**

(Autonomous College Affiliated to the University of Mumbai)

NAAC Accredited with "A" Grade (CGPA : 3.18)



**Department of Computer Science and Engineering (Data Science)**

threads while preventing out-of-bounds memory access. By calculating the required number of blocks using the formula  $(N + \text{threads\_per\_block} - 1) / \text{threads\_per\_block}$ , we guaranteed that all elements in the array were processed. Additionally, we implemented a boundary check within the kernel to ensure that excess threads did not attempt to access invalid memory locations. This approach optimizes parallel execution while maintaining correctness, demonstrating an efficient method for handling mismatched execution configurations in GPU programming.