# PES University Electronic City Campus

Hosur Rd, Konappana Agrahara, Electronic City, Bengaluru, Karnataka 560100

A Project Report on

## *"CUDA PROGRAM AND CACHE CONFIGURATION"*

Submitted by:
- **MADHU L          (PES2UG20EC802)**
- **ADITHYA NAIR (PES2UG20EC008)**

Under the guidance of:
**PROF. MAHESH AWATI**
**Associate Professor,**
**Department of Electronics and Communication Engineering**
**PES UNIVERSITY**
Electronic City, Bangalore-560100
2022 - 2023

# ASSIGNMENT 1: CUDA PROGRAM

## N X N MATRIX MULTIPLICATION

## CODE:

```c
#include "cuda_runtime.h"
#include "device_launch_parameters.h"
#include <stdio.h>

#define N 2

__global__ void matrixMult(int* A, int* B, int* C, int n)
{
    int i = threadIdx.x;
    int j = threadIdx.y;
    int Pvalue = 0;
    for (int k = 0; k < n; k++)
    {
        Pvalue += A[i * n + k] * B[k * n + j];
    }
    C[i * n + j] = Pvalue;
}

int main()
{
    int A[N * N] = { 6, 7,8, 2 };
    int B[N * N] = { 5, 6, 7, 8 };
    int C[N * N] = { 0, 0, 0, 0 };
    int* dev_A = NULL, * dev_B = NULL, * dev_C = NULL;

    cudaError_t err = cudaMalloc((void**)&dev_A, N * N *
sizeof(int));
    if (err != cudaSuccess)
    {
        printf("Error allocating memory for dev_A: %s\n",
cudaGetErrorString(err));
        return -1;
    }

    err = cudaMalloc((void**)&dev_B, N * N * sizeof(int));
    if (err != cudaSuccess)
    {
        printf("Error allocating memory for dev_B: %s\n",
cudaGetErrorString(err));
```

```
        cudaFree(dev_A);
        return -1;
    }

    err = cudaMalloc((void**)&dev_C, N * N * sizeof(int));
    if (err != cudaSuccess)
    {
        printf("Error allocating memory for dev_C: %s\n",
cudaGetErrorString(err));
        cudaFree(dev_A);
        cudaFree(dev_B);
        return -1;
    }

    err = cudaMemcpy(dev_A, A, N * N * sizeof(int),
cudaMemcpyHostToDevice);
    if (err != cudaSuccess)
    {
        printf("Error copying A to device: %s\n",
cudaGetErrorString(err));
        cudaFree(dev_A);
        cudaFree(dev_B);
        cudaFree(dev_C);
        return -1;
    }

    err = cudaMemcpy(dev_B, B, N * N * sizeof(int),
cudaMemcpyHostToDevice);
    if (err != cudaSuccess)
    {
        printf("Error copying B to device: %s\n",
cudaGetErrorString(err));
        cudaFree(dev_A);
        cudaFree(dev_B);
        cudaFree(dev_C);
        return -1;
    }

    dim3 threadsPerBlock(N, N);
    matrixMult << <1, threadsPerBlock >> > (dev_A, dev_B, dev_C,
N);

    err = cudaMemcpy(C, dev_C, N * N * sizeof(int),
cudaMemcpyDeviceToHost);
    if (err != cudaSuccess)
    {
        printf("Error copying C from device: %s\n",
cudaGetErrorString(err));
        cudaFree(dev_A);
        cudaFree(dev_B);
```

```
        cudaFree(dev_C);
        return -1;
    }

    cudaFree(dev_A);
    cudaFree(dev_B);
    cudaFree(dev_C);

    printf("Matrix A:\n");
    for (int i = 0; i < N; i++)
    {
        for (int j = 0; j < N; j++)
        {
            printf("%d ", A[i * N + j]);
        }
        printf("\n");
    }

    printf("Matrix B:\n");
    for (int i = 0; i < N; i++)
    {
        for (int j = 0; j < N; j++)
        {
            printf("%d ", B[i * N + j]);
        }
        printf("\n");
    }
    printf("Resultant Matrix C:\n");
    for (int i = 0; i < N; i++)
    {
        for (int j = 0; j < N; j++)
        {
            printf("%d ", C[i * N + j]);
        }
        printf("\n");
    }

    return 0;
}
```

# NOTE:

```
cudaMalloc((void **)&dev_A, N*N*sizeof(int));
cudaMalloc((void **)&dev_B, N*N*sizeof(int));
cudaMalloc((void **)&dev_C, N*N*sizeof(int));
```

**ALLOCATE MEMORY ON GPU FOR dev_A,dev_B,dev_C using cudaMalloc**

```
cudaMemcpy(dev_A, A, N*N*sizeof(int), cudaMemcpyHostToDevice);
cudaMemcpy(dev_B, B, N*N*sizeof(int), cudaMemcpyHostToDevice);
```

**COPIES THE MATRICES FROM THE HOST TO GPU USING cudaMemcpy( ).**

```
dim3 threadsPerBlock(N,N);
matrixMult<<<1,threadsPerBlock>>>(dev_A, dev_B, dev_C);
```
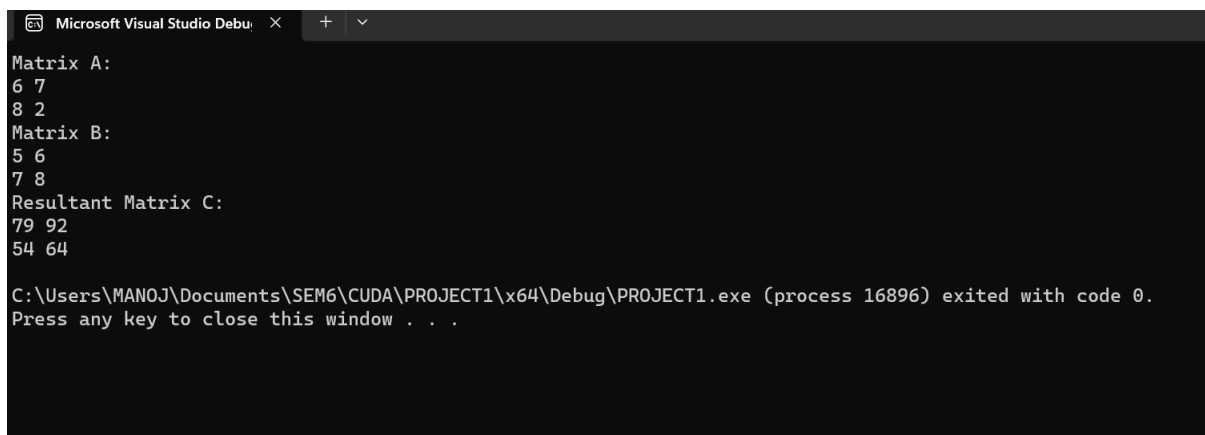
**DEFINES THE NUMBER OF THREADS PER BLOCK AS 2D BLOCK WITH DIM NxN**

```
cudaMemcpy(C, dev_C, N*N*sizeof(int), cudaMemcpyDeviceToHost);
```
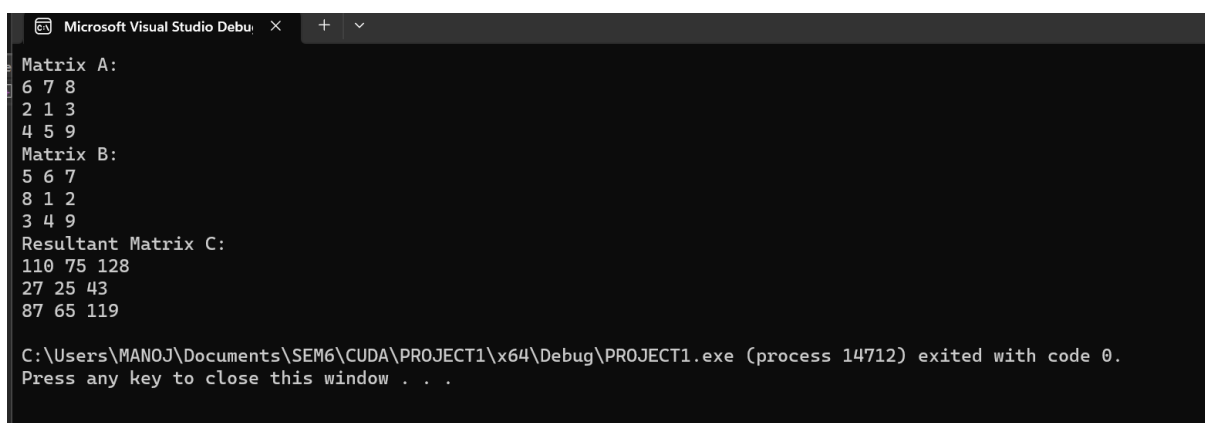
**COPIES MATRIX C FROM GPU TO HOST CPU USING cudaMemcpy( ).**

# OUTPUT:

1.



2.

# ASSIGNMENT 2: CACHE CONFIGURATION

1. Configure a 64-entry 8-word direct mapped cache.

2. Configure a 16-entry 4-word 4-way set-associative cache.
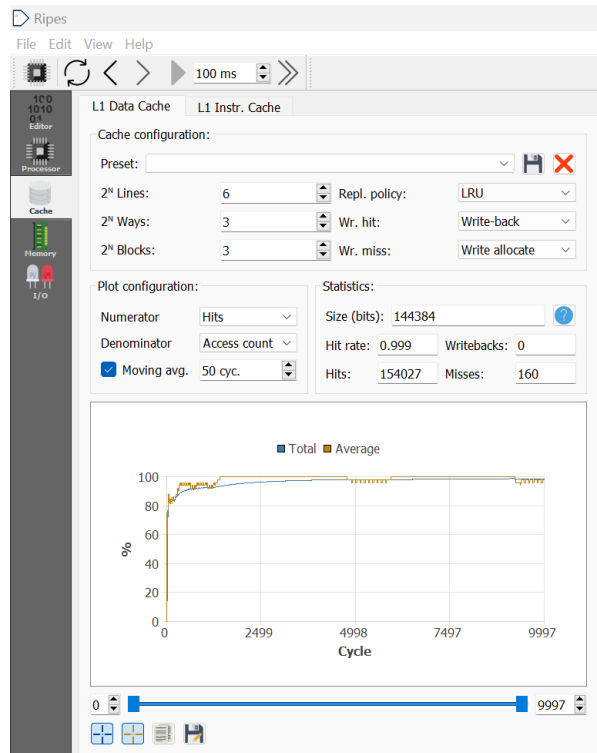
3. Configure a 16-entry 2-word 4-way set-associative cache with write-through.
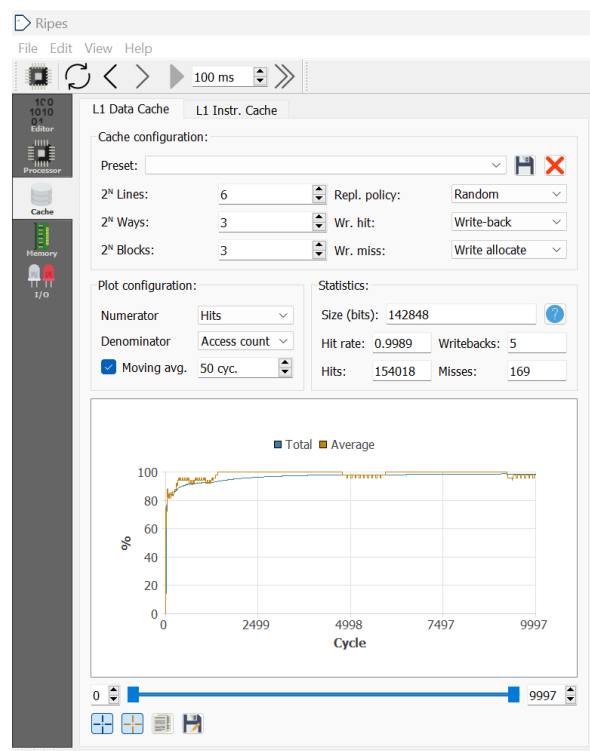
4. Configure a 64-entry 8-word fully associative cache with Least Recently Used replacement policy and report the numbers. Change the replacement policy to Random and report the numbers for the same cache.
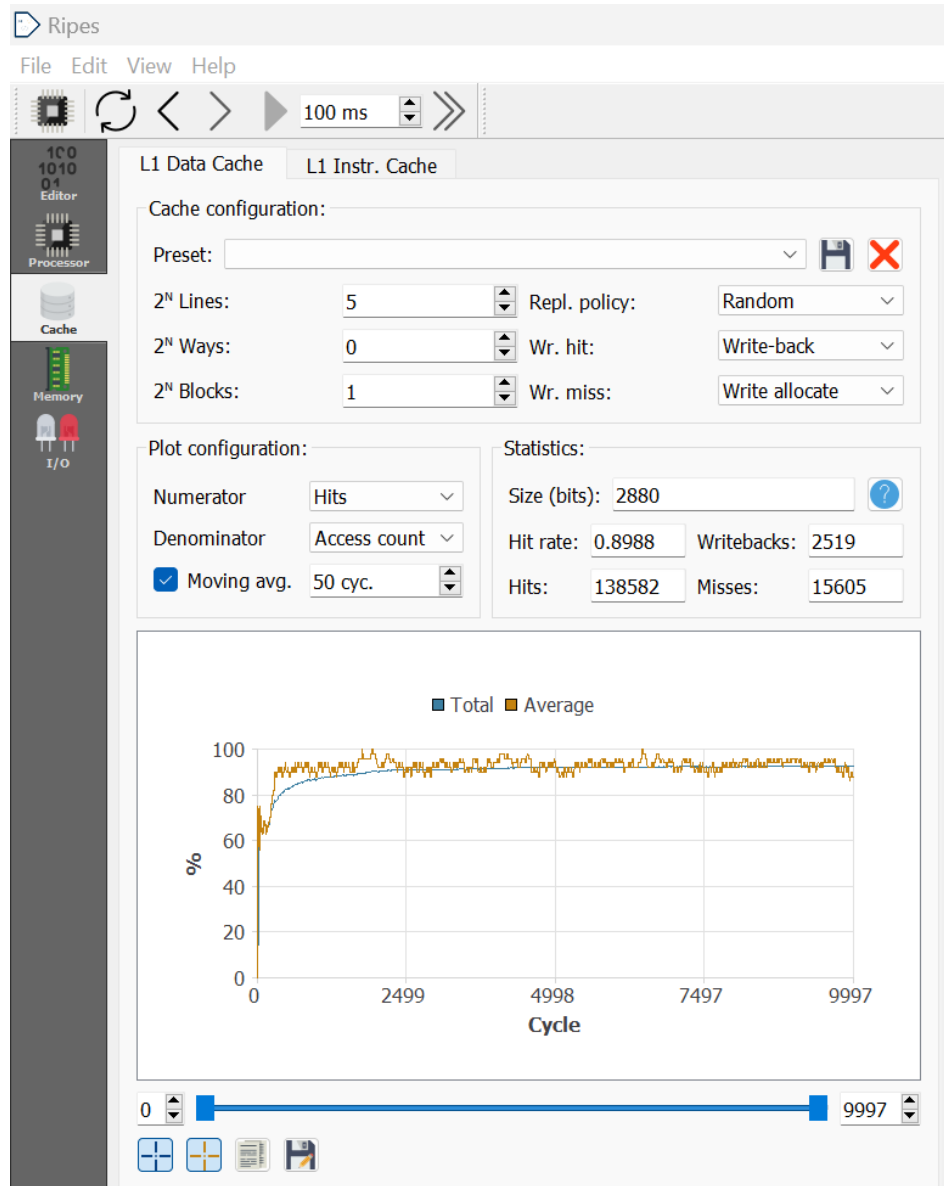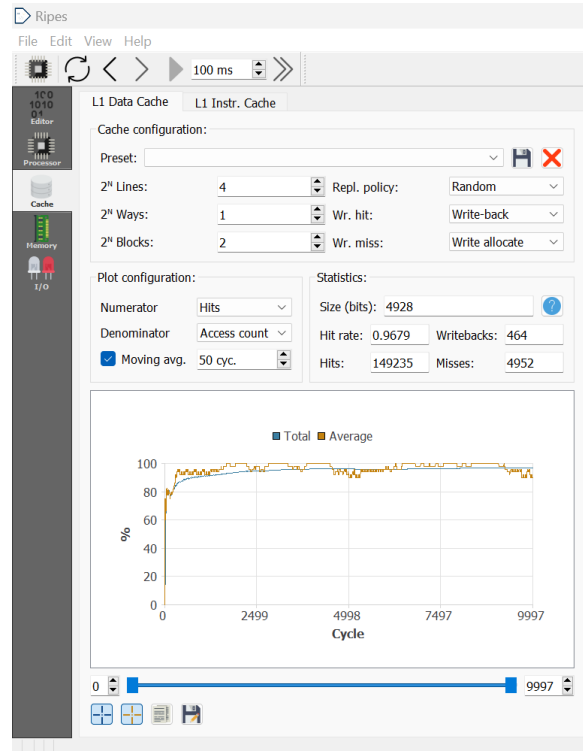
1. LRU



2.RANDOM

5. Configure a 32-entry 2-word direct-mapped cache and plot a graph using the plot configuration with numerator as Hits and denominator as Access count. Explain why the number of hits increases and decreases back down before increasing again.
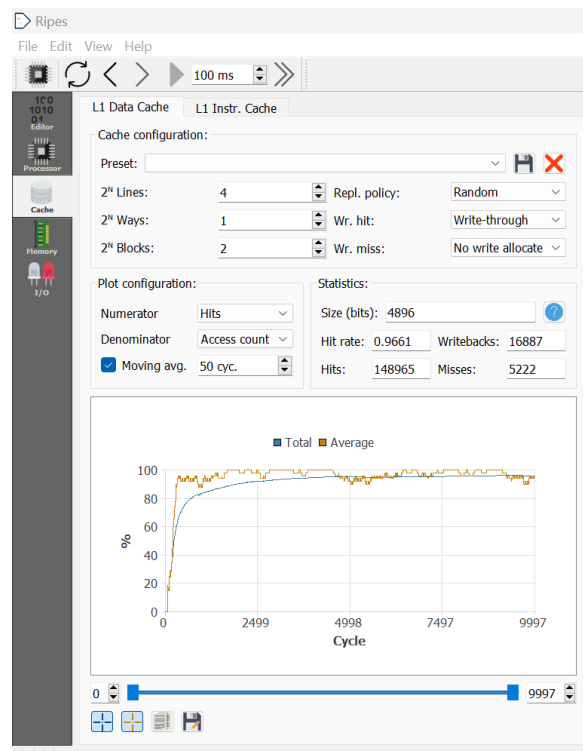


In a 32-entry 2-word direct-mapped cache, the number of hits increases and decreases back down before increasing again because of cache conflicts. When a program accesses memory blocks that map to the same cache line, the cache can only hold one of the blocks at a time, so it will evict one and fetch the other, resulting in a cache miss. As the program continues to access these memory blocks, the cache will alternate between holding one block and the other, resulting in alternating hits and misses on the cache line. Therefore, the number of hits may increase up to a maximum of 16 per cache line (since there are only 16 cache lines available for each set of two memory blocks that map to the same cache line), but will eventually decrease due to evictions and cache misses, before increasing again when the program accesses the memory blocks that were evicted earlier.

6. Configure a 16-entry 4-word 2-way set-associative cache with write-back and write allocate and report the numbers. For the same configuration of cache, use write-through and no write allocate and report the numbers. Explain the differences between the two cache configurations.

## 1. WITH WRITE BACK AND WRITE ALLOCATE



## 2.WITH WRITE THROUGH AND NO WRITE ALLOCATE

Cache with write-back and write-allocate policies, write operations are first written to the cache and marked as "dirty", but not immediately written to main memory. When a cache line containing dirty data is evicted from the cache, the dirty data is then written back to main memory. Write-allocate means that when a write operation misses in the cache, the entire cache block is brought into the cache before the write is performed.

In a Cache with write-through and no-write-allocate policies, write operations are immediately written both to the cache and main memory. No-write-allocate means that when a write operation misses in the cache, the write is directly performed in main memory without bringing the block into the cache.