

Name: Nirmal Kumar Sedhumadhavan

**ECE 786 Programming Assignment 1**  
**CUDA Programming and GPGPU Simulator**

**Part- A: CUDA Program**

CUDA Program for Single Quantum Bit Gate Simulation with different memory management methods.

**CUDA Kernel Function:**

```
__global__ void matrix_multiply(const float *input, float *output, const float
*Umatrix, int size, int qbit)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int mask = 1 << qbit;
    int index = i ^ mask;
    if (i <= (size - mask))
    {
        if((i/mask) % 2 != 1)
        {
            output[i] = Umatrix[0] * input[i] + Umatrix[1] * input[index];
            output[index] = Umatrix[2] * input[i] + Umatrix[3] * input[index];
        }
    }
}
```

**Explanation:**

The implemented kernel function `matrix_multiply` takes input vector (`float *input`), unitary matrix vector (`float *Umatrix`), qbit position (`int qbit`), size of the input vector (`int size`), output vector (`float *output`) as arguments.

**`int i = blockIdx.x * blockDim.x + threadIdx.x;`**

Inside the function, index for the threads is calculated based on the thread hierarchy. Each thread processes 2 elements in the output vector.

**`mask = 1 << qbit;`** //Calculates the value of 2 to the power of qbit( $2^{\text{qbit}}$ )

**`index = i ^ mask;`** //Calculates the position of the 2<sup>nd</sup> element

**`if (i <= (size - mask))`**

```
{
    if((i/mask) % 2 != 1) // Check if the 2nd element is already calculated
    {
        output[i] = Umatrix[0] * input[i] + Umatrix[1] * input[index]; //Update the 1st element
        output[index] = Umatrix[2] * input[i] + Umatrix[3] * input[index]; //Update the 2nd element
    }
}
```

The above Kernel function is used for two different memory management methods which are discrete/separate CPU and GPU memories and Unified Virtual Memory with 256 threads in each thread block.

**Timing Report:**

Input (Taken from more examples of Program Assignment 1)	Time taken with Separate Memories (microseconds)	Time taken with Unified Virtual Memory (microseconds)
4-qubit quantum output state after applying a single-qubit gate on qubit 1	26	29
7-qubit quantum output state after applying a single-qubit gate on qubit 4	31	30

**Part-B: GPGPU Sim**

Ran the above CUDA code (for separate CPU and GPU memory management method) with the GPGPU sim and observed the statistics.

quamsimV1.cu is run with input taken from more examples of Program Assignment 1 (7-qubit quantum output state after applying a single-qubit gate on qubit 4) and below are the statistics.

1) IPC of the program

gpu\_sim\_cycle = 5779

gpu\_sim\_insn = 5990

gpu\_ipc = gpu\_sim\_insn / gpu\_sim\_cycle

gpu\_ipc = 1.0365

2) Data Cache miss rate

L1D\_total\_cache\_accesses = 64

L1D\_total\_cache\_misses = 40

L1D\_total\_cache\_miss\_rate = L1D\_total\_cache\_misses / L1D\_total\_cache\_accesses

L1D\_total\_cache\_miss\_rate = 0.625