

PCAP Assignment IV

1.

- A. Write an OpenCL kernel code which converts input matrix A into output matrix B as follows: consider elements of matrix A into 4 equal parts. First part elements should be incremented by 1, Second part elements should be incremented by 2, Third part elements should be incremented by 3 and last part elements should be incremented by 4. Write the kernel code which does this in parallel for all the elements. Example, N = 4

3	8	2	5		4	9	4	7
2	3	5	6		3	4	7	8
2	4	3	1		5	7	7	5
3	2	1	5		6	5	5	9

```
__kernel void operate (__global int *A, __global int *B, int count) {
    int i = get_global_id(0);
    int j = get_global_id(1);
    // 'd' is the increment value
    int d = 1 + (i >= count/2) * 2 + (j >= count/2);
    B[i * count + j] = A[i * count + j] + d;
}
```

- B. Write OpenCL kernel code to calculate the value of Π .

```
__kernel void advance ( __global float *output, const unsigned int
count) {
    int n_rects = get_global_id(0);
    float rect_width = 1.0 / n_rects;
    float rect_left = 0.0;
    float rect_height_squared;
    float rect_height;
    float agg_area = 0.0;
    float pi;
    int i;
    if (n_rects < count) {
        for (i = 0; i < n_rects; i++) {
            rect_left += rect_width;
            rect_height_squared = 1 - rect_left * rect_left;
            if (rect_height_squared < 0.00001) {
                rect_height_squared = 0;
            }
            rect_height = sqrt(rect_height_squared);
            agg_area += rect_width * rect_height;
        }
        pi = 4 * agg_area;
        output[n_rects] = pi;
    }
}
```

2.

- A. Write a parallel program in CUDA to multiply two Matrices A and B of dimensions MxN and NxP resulting in Matrix C of dimension MxP. Create P number of threads, and each column of the resultant matrix is to be computed by one thread. Use 1D Block.

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

__global__ void matrix_mult (int *c, int *a, int *b, int ha, int wa) {
    int j = threadIdx.x;
    int wc = blockDim.x;
    int i, k;
    for (i = 0; i < ha; i++) {
        for (k = 0; k < wa; k++) {
            c[i*wc + j] += a[i*wa + k] * b[k*wc + j];
        }
    }
}

int main () {
    int ha = 2; int wa = 3; int hb = 3; int wb = 2; // Shhhh, it's okay
    int hc = ha;
    int wc = wb;

    int a[6] = { 1, 2, 5, 3, 4, 6 };
    int b[6] = { 7, 10, 8, 11, 9, 12 };
    int c[4] = { 0 };

    int size_a = ha * wa, size_b = hb * wb, size_c = hc * wc;

    // Print A, B

    int *dev_a, *dev_b, *dev_c = NULL;

    cudaMalloc((void*)&dev_a, size_a * sizeof(int));
    cudaMalloc((void*)&dev_b, size_b * sizeof(int));
    cudaMalloc((void*)&dev_c, size_c * sizeof(int));

    cudaMemcpy(dev_a, a, size_a * sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(dev_b, b, size_b * sizeof(int), cudaMemcpyHostToDevice);

    dim3 threadsPerBlock(wc,1); //just for the sake of 2D-Blocks;
    matrixMult <<<1, threadsPerBlock >>>(dev_c, dev_a, dev_b, ha, wa);

    cudaMemcpy(c, dev_c, size_c * sizeof(int), cudaMemcpyDeviceToHost);

    // Print C

    return 0;
}
```

B. Compare Superscalar, VLIW and Hardware Multithreading.

Superscalar

In these designs, the CPU maintains dependence information between instructions in the instruction stream and schedules work onto unused functional units when possible.

The major advantage of out-of-order logic lies with the software developer. By extracting parallelism from the programmer's code automatically within the hardware, serial code performs faster without any extra developer effort.

Out-of-order scheduling logic requires a substantial area of the CPU die to maintain dependence information and queues of instructions to deal with dynamic schedules throughout the hardware.

VLIW

VLIW is a heavily compiler-dependent method for increasing instruction-level parallelism in a processor. Rather than depending entirely on complex out-of-order control logic that maintains dependences in hardware, as we saw when discussing superscalar execution, VLIW moves this dependence analysis work into the compiler.

Instead of providing a scalar instruction stream, each issued instruction in a VLIW processor is a long instruction word comprising multiple instructions intended to be issued in parallel. This instruction will be mapped directly to the execution pipelines of the processor.

Hardware Multithreading

The third common form of parallelism after instruction and data is thread parallelism, or the execution of multiple independent instruction streams, extracting independent instructions from an instruction stream is difficult, in terms of both hardware and compiler work, and it is sometimes impossible. Extracting such instruction parallelism from two independent threads is trivial because those threads already guarantee independence.

There are two main ways to apply on-chip multithreading: Simultaneous multithreading and Temporal multithreading

3.

A. Write an OpenCL Kernel code to perform odd-even transposition sorting in parallel. What is grid?

```
__kernel void even (__global int *a) {
    int idx = get_global_id(0);
    int size = get_global_size(0);
    int temp;
    if ((idx % 2 == 0) && ((idx + 1) < size)) {
        if (a[idx] > a[idx+1]) {
            temp = a[idx];
            a[idx] = a[idx+1];
            a[idx+1] = temp;
        }
    }
}
```

```

__kernel void odd (__global int *a) {
    int idx = get_global_id(0);
    int size = get_global_size(0);
    int temp;
    if ((idx % 2 != 0) && ((idx + 1) < size)) {
        if(a[idx] > a[idx+1]) {
            temp = a[idx];
            a[idx] = a[idx+1];
            a[idx+1] = temp;
        }
    }
}

```

In OpenCL execution model, the kernel instances execute concurrently over a virtual **grid** defined by the host code.

- B. How many grids will be created for odd-even transposition sorting and merge sorting.

Odd Even Transposition: $2 \times \text{elements} / 2 = \text{elements}$

Merge Sort: $\log N$

4.

- A. Write an OpenCL Kernel code to perform merge sort in parallel.

```

void __kernel CompareExchange (__global int* A,int pass,__global int* t) {

    int idx = get_global_id(0);
    int size = get_global_size(0);
    int i = idx * 2 * pass;
    int j = idx * 2 * pass + pass;
    int flag = 0, flag2 = pass;

    for (int k = 0; k <= pass; k += 1) {
        int temp2 = A[i+k];
        flag = 0;

        for (int z = 0; z < pass; z += 1) {
            int temp3 = A[j+z];

            if (A[i+k] >= A[j+z]) {
                int temp = A[i+k];
                A[i+k] = A[j+z];
                A[j+z] = temp;
                flag++;
                k++;
            }
        }
    }
}

```

```

        if (A[i+k] < A[j+z]) {
            flag++;
            k++;
        }
    }

    while (flag > 0) {
        k--;
        flag--;
    }
}
}

```

B. What is data parallelism? Explain With the help of suitable example.

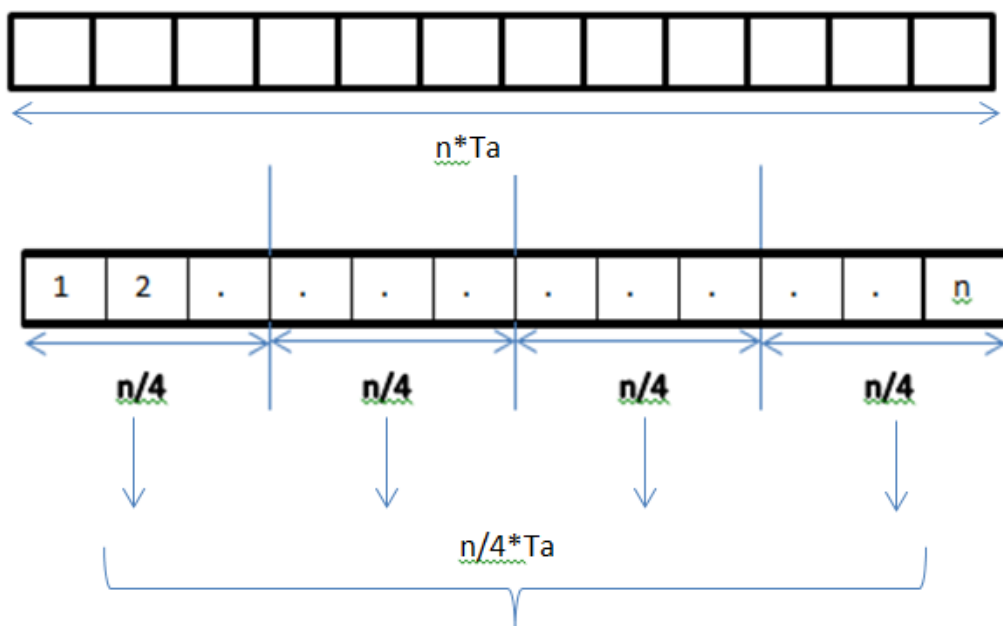
Data parallelism is a form of parallelization across multiple processors in parallel computing environments. It focuses on distributing the data across different nodes, which operate on the data in parallel. It can be applied on regular data structures like arrays and matrices by working on each element in parallel.

Let us assume we want to sum all the elements of the given array and the time for a single addition operation is T_a time units.

In the case of sequential execution, the time taken by the process will be $n \cdot T_a$ time units as it sums up all the elements of an array.

On the other hand, if we execute this job as a data parallel job on 4 processors the time taken would reduce to $(n/4) \cdot T_a + \text{Merging overhead time units}$.

Parallel execution results in a speedup of 4 over sequential execution.



5.

- A. Write a parallel program in CUDA to add two Matrices A and B of dimensions MxN resulting in Matrix C of dimension MxN. Create N number of threads, and each column of the resultant matrix is to be computed by one thread. Handle errors.

```
#include "cuda_runtime.h"
#include "device_launch_parameters.h"

#include <stdio.h>
#include <stdlib.h>

__global__ void matrixAdd(int *c, int *a, int *b, int M) {
    int j = threadIdx.x;
    int k = blockDim.x;
    int i;

    for (i = 0; i < M; i++)
        c[i*k + j] = a[i*k + j] + b[i*k + j];
}

void checkForError(cudaError_t e) {
    if (e != cudaSuccess) {
        printf("Error : %s\n", cudaGetErrorString(e));
        getchar();
        exit(1);
    }
}

int main()
{
    int M = 2, N = 2;

    int a[4] = { 1, 2, 3, 4 };
    int b[4] = { 1, 2, 3, 4 };
    int c[4] = { 0 };

    int *dev_a = NULL, *dev_b = NULL, *dev_c = NULL;

    int size = M * N;

    cudaError_t e;

    e = cudaMalloc((void**)&dev_c, size * sizeof(int));
    checkForError(e);

    e = cudaMalloc((void**)&dev_a, size * sizeof(int));
    checkForError(e);
    e = cudaMalloc((void**)&dev_b, size * sizeof(int));
    checkForError(e);

    e = cudaMemcpy(dev_a, a, size * sizeof(int), cudaMemcpyHostToDevice);
    checkForError(e);
```

```

e = cudaMemcpy(dev_b, b, size * sizeof(int), cudaMemcpyHostToDevice);
checkForError(e);

matrixAdd <<< 1, N >>>(dev_c, dev_a, dev_b, M);
e = cudaGetLastError();
checkForError(e);

e = cudaMemcpy(c, dev_c, size * sizeof(int), cudaMemcpyDeviceToHost);
checkForError(e);

int i;
for (i = 0; i < size; i++) {
    if (i % N == 0)
        printf("\n");
    printf("%d ", c[i]);
}

return 0;
}

```

B. Compare different parallel programming models.

Shared Memory Model

In a shared-memory model, parallel processes share a global address space that they read and write to asynchronously. Asynchronous concurrent access can lead to race conditions and mechanisms such as locks, semaphores and monitors can be used to avoid these.

Message Passing Model

In a message-passing model, parallel processes exchange data through passing messages to one another. These communications can be asynchronous, where a message can be sent before the receiver is ready, or synchronous, where the receiver must be ready.

Implicit Interaction Model

In an implicit model, no process interaction is visible to the programmer and instead the compiler and/or runtime is responsible for performing it. Two examples of implicit parallelism are with domain-specific languages where the concurrency within high-level operations is prescribed, and with functional programming languages because the absence of side-effects allows non-dependent functions to be executed in parallel.