
HIGH PERFORMANCE COMPUTING

Project - CUDA

Topic - Using Discrete Fourier Transform for finding the normalized frequency of a random wave and parallelizing the code using CUDA.

CONTENTS

S.no.	Topic	Page no.
1.	Introduction	2
2.	Formulae used	4
3.	Functions used	4
4.	Procedure	5
5.	Sample inputs and outputs	6
6.	Execution time, speed-up and parallel fraction	9
7.	Code	11
8.	Profiling	13

Bhaavanaa Thumu - CED17I021

1. Introduction

About discrete fourier transform :

The discrete fourier transform converts a finite sequence of equally-spaced samples of a function into a same-length sequence of equally-spaced samples of the discrete-time Fourier transform (DTFT), which is a complex valued function of frequency. The DFT is therefore said to be a frequency domain representation of the original input sequence.

The DFT is the most important discrete transform, used to perform Fourier analysis in many practical applications. In digital signal processing, the function is any quantity or signal that varies over time, such as the pressure of a sound wave, a radio signal, or daily temperature readings, sampled over a finite time interval. In image processing, the samples can be the values of pixels along a row or column of a raster image. The DFT is also used to efficiently solve partial differential equations, and to perform other operations such as convolutions or multiplying large integers.

Since it deals with a finite amount of data, it can be implemented in computers by numerical algorithms or even dedicated hardware. The advantage of DFT over FFT is the restriction of 2^n samples and fixed/discrete sample rates for the FFT.

About power spectrum :

For a given signal, the power spectrum gives a plot of the portion of a signal's power (energy per unit time) falling within given frequency bins. The most common way of generating a power spectrum is by using the discrete fourier transform, and other techniques such as the maximum entropy method can also be used.

About normalized frequency :

Normalised frequency is frequency in Hz (or more generically cycles/second or some other unit) divided by the sample frequency of your signal in Hz (or the same units as your original frequency).

2. Formulae used

$$\text{Speed-up} = T(1) / T(P)$$

$$\text{Parallel fraction, } f = (1 - T(P)/T(1)) / (1 - (1/P))$$

where,

$T(1)$ - time taken for serial execution

$T(P)$ - time taken for parallel execution

P - number of processes/threads/processors

Formula for finding DFT of a sample -

$$X(k) = \sum_{n=0}^{N-1} x(n) e^{-j2\pi kn/N}$$

3. Functions used

cudaMalloc - allocates a block of memory

cudaMemcpy - copies 'n' characters from source to destination memory area.

4. Procedure

The steps for performing the task of finding the normalized frequency of a random wave using discrete fourier transform are -

- a. Produce an array of N random numbers within a range, e.g. -1 to 1 using the rand() function. This would be our array of samples.
- b. Define a struct which contains real and imaginary double precision numbers, so that it can store the real and imaginary coefficients of the output of DFT. Initialize the real and imaginary values to 0.
- c. Find the DFT of all the samples present in the array.

$$X(k) = \sum_{n=0}^{N-1} x(n)e^{-j2\pi kn/N}$$

- d. For every element in the array of structures containing the coefficients of the DFT obtained above, find the power spectrum of it. It is the sum of the square of the real coefficient and the square of the imaginary coefficient of the DFT.
- e. Find the sum of the product of the frequency and the power spectrum value of all the elements. Here,

frequency of an element = $2\pi k/N$

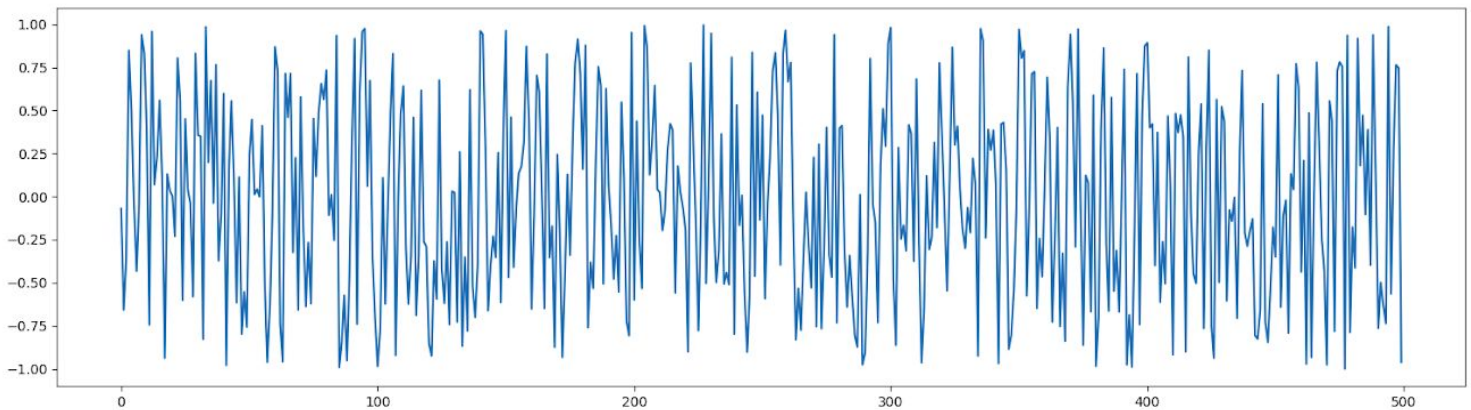
- f. Find the sum of the power spectrum values.
- g. The normalized frequency of the random wave generated above is the sum of the product of the frequency and the power spectrum value of all the elements (obtained in step e) divided by the sum of the power spectrum values (obtained in step f).

5. Sample inputs and outputs

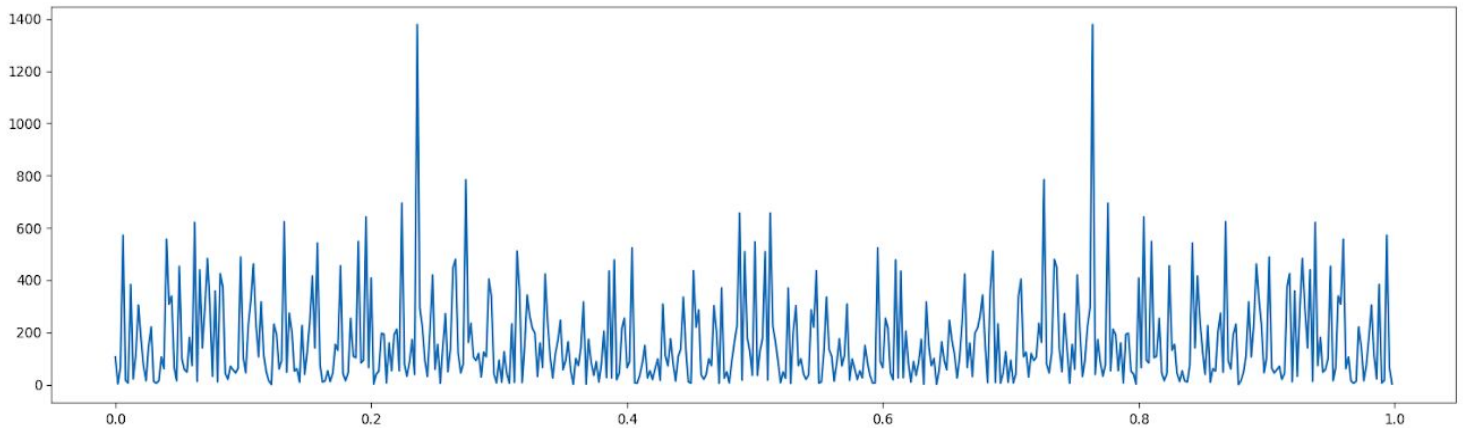
Example 1:

Input :- No. of random samples : 500

Range of the random samples : -1 to 1



Random input signal



Power spectrum vs k/N of the random input signal

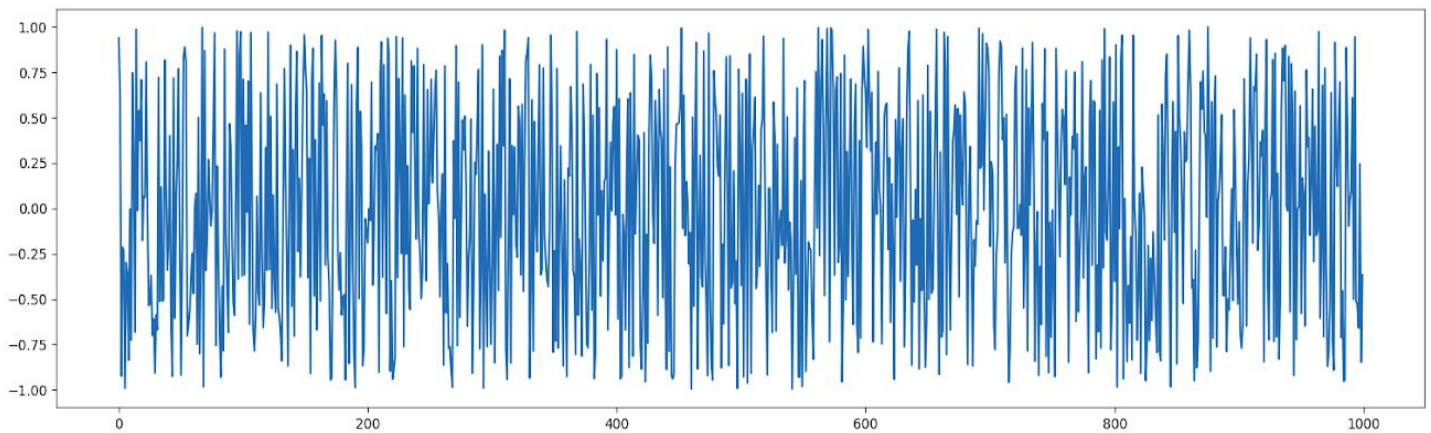
Output :- Normalized frequency of the wave : 3.137617 Hz

Time taken for the program for the given number of processes.

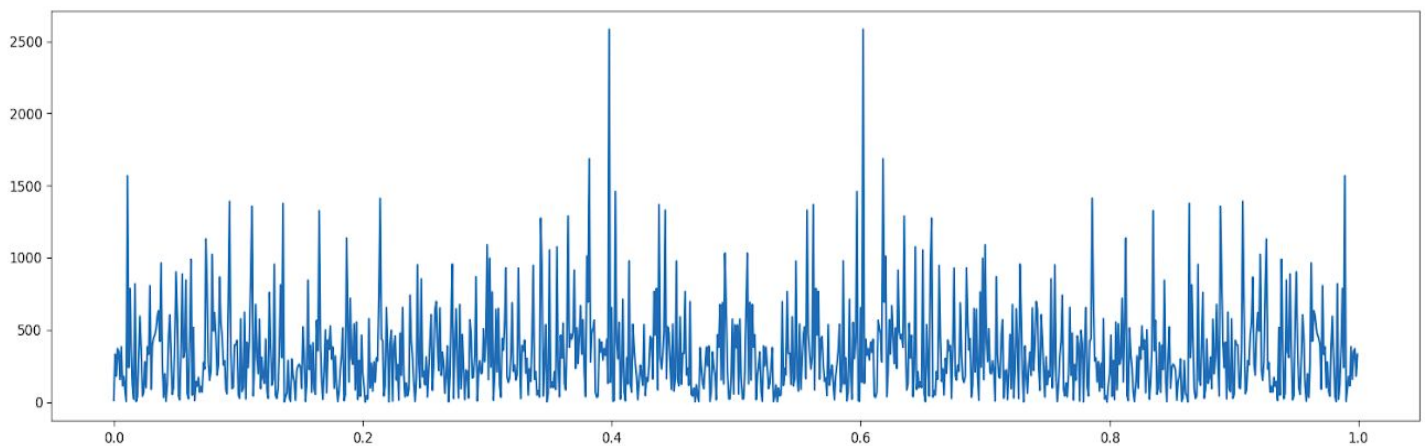
Example 2:

Input :- No. of random samples : 1,000

Range of the random samples : -1 to 1



Random input signal



Power spectrum vs k/N of the random input signal

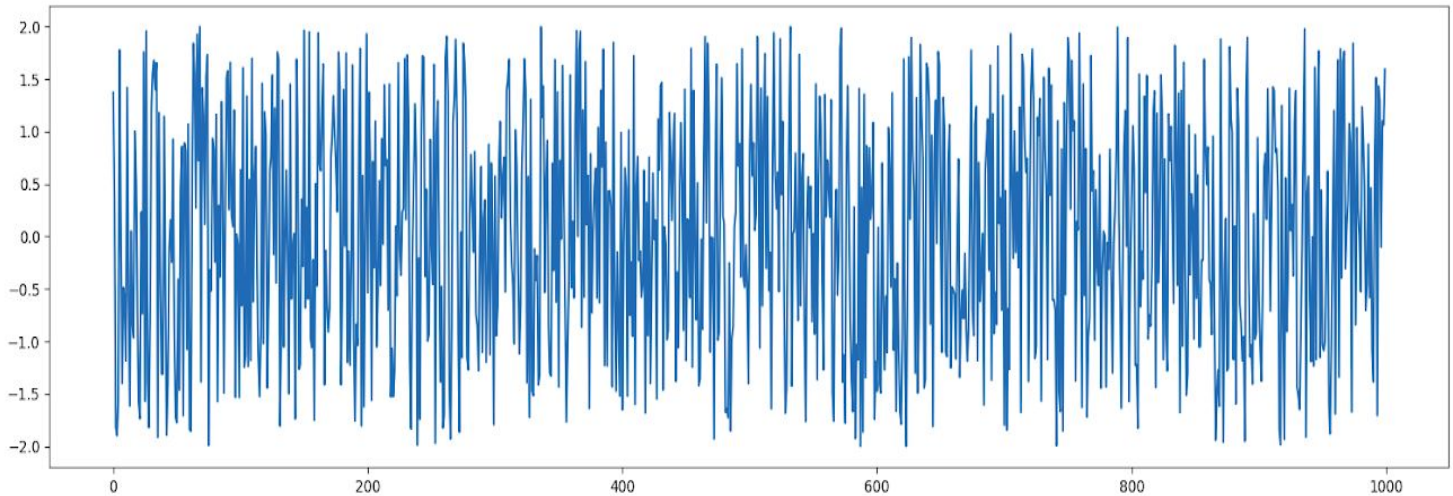
Output :- Normalized frequency of the wave : 3.141485 Hz

Time taken for the program for the given number of processes.

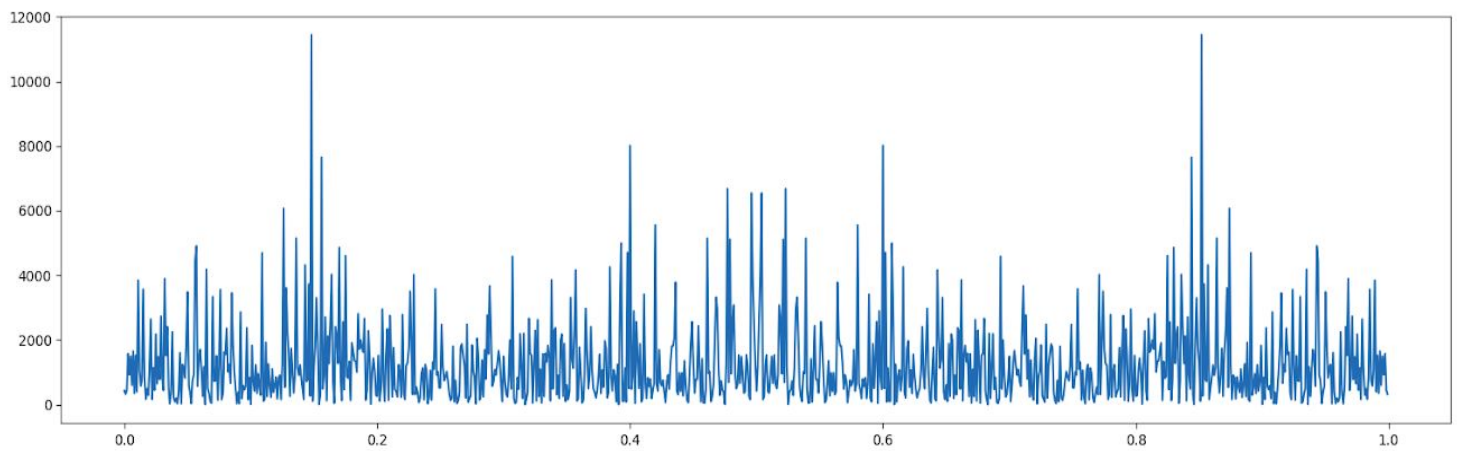
Example 3:

Input :- No. of random samples : 1000

Range of the random samples : -2 to 2



Random input signal



Power spectrum vs k/N of the random input signal

Output :- Normalized frequency of the wave : 3.140520 Hz

Time taken for the program for the given number of processes.

6. Execution time, speed-up and parallel fraction

For the below analysis, the program contains 10,00,000 sample points in the random wave. The range for these random points, considered here, is -1 to 1.

The number of processors initialized and the corresponding execution time, speed-up and parallel fraction is as shown in the table below.

No. of processes	Execution time	Speed-up	Parallel fraction
1	131.711246	1	#DIV/0!
2	132.350666	0.9951687436	-0.009709421472
4	65.922643	1.997966708	0.6659882128
8	33.052886	3.984863712	0.856057587
16	16.692128	7.890620417	0.93148507
32	8.578441	15.35375087	0.9650264106
64	4.625175	28.47702973	0.9801995968
128	2.949827	44.65049849	0.9853015014
256	2.435393	54.08213212	0.985358662
500	2.015343	65.35425781	0.9866721224

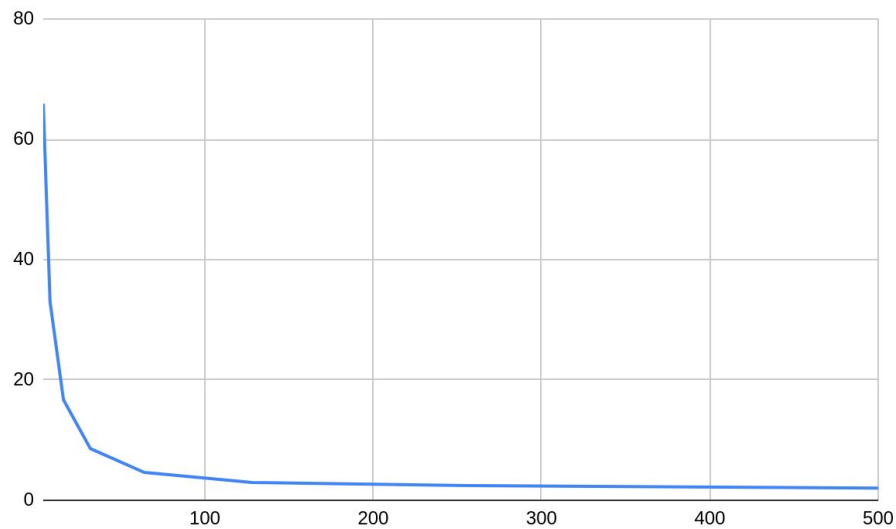
Hence, in this case, the optimal number of processes is 500. In this case-

Execution time = 2.015343

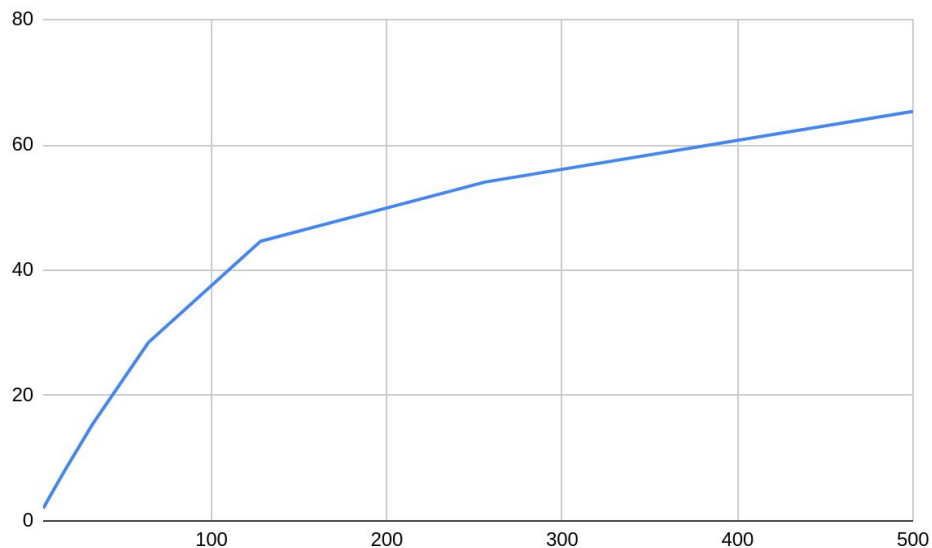
Speed-up = 65.35425781

Parallel fraction = 0.9866721224

The plot below is the graph between the **number of processes** and **execution time**.



The plot below is the graph between the **number of processes** and **speed-up**.



In this case, finding the DFT of the elements is parallelized. The load is equally distributed among all the threads and the leftover load is also done by the master process.

7. Code

```
[ ] 1 %tcu
2
3 #include<stdio.h>
4 #include<stdlib.h>
5 #include<math.h>
6 #include <sys/time.h>
7
8 #define PI 3.14159265
9
10 #define N 10000 // the number of random samples
11 #define thread_count 500 // the number of threads initialized
12
13
14 global__ void dft(double *d_samples, double *d_dft_real, double *d_dft_img, double *d_k_by_N, double *d_power_spectrum, double *d_sum_power_spectrum, double *d_sum_freq_power)
15 {
16     int load = N/thread_count; // calculating the load per thread
17
18     __shared__ double partial_sum_power_spectrum[thread_count]; // initializing a partial sum variable which will find the partial sums of the power_spectrum for the load assigned to it
19     __shared__ double partial_sum_freq_power[thread_count]; // initializing a partial sum variable which will find the partial sums of the freq_power for the load assigned to it
20
21     partial_sum_power_spectrum[threadIdx.x] = 0; // initializing the partial sum variables to 0
22     partial_sum_freq_power[threadIdx.x] = 0;
23
24     if(threadIdx.x != thread_count-1) // if it is not the last thread, it will do load amt of work
25     {
26         for(int k=threadIdx.x*load; k<(threadIdx.x+1)*load; k++) // for every element in the load
27         {
28             for(int n=0; n<N; n++)
29             {
30                 d_dft_real[k] += d_samples[n] * cos((2*PI*k*n)/N); // find the real and img part of the dft
31                 d_dft_img[k] += d_samples[n] * sin((2*PI*k*n)/N);
32             }
33
34             d_k_by_N[k] = (double)k/N; // find the k/N value
35
36             d_power_spectrum[k] = d_dft_real[k]*d_dft_real[k] + d_dft_img[k]*d_dft_img[k]; // find the power spectrum value
37
38             partial_sum_power_spectrum[threadIdx.x] += d_power_spectrum[k]; // obtain the partial power_spectrum sum of the load assigned to the thread
39             partial_sum_freq_power[threadIdx.x] += 2*PI*d_k_by_N[k]*d_power_spectrum[k]; // obtain the partial freq_power sum of the load assigned to the thread
40         }
41     }
42     else // if it is last thread, then it will do load+leftover work
43     {
44         for(int k=threadIdx.x*load; k<N; k++)
45         {
46             for(int n=0; n<N; n++)
47             {
48                 d_dft_real[k] += d_samples[n] * cos((2*PI*k*n)/N);
49                 d_dft_img[k] += d_samples[n] * sin((2*PI*k*n)/N);
50             }
51
52             d_k_by_N[k] = (double)k/N;
53
54             d_power_spectrum[k] = d_dft_real[k]*d_dft_real[k] + d_dft_img[k]*d_dft_img[k];
55
56             partial_sum_power_spectrum[threadIdx.x] += d_power_spectrum[k];
57             partial_sum_freq_power[threadIdx.x] += 2*PI*d_k_by_N[k]*d_power_spectrum[k];
58         }
59     }
60
61     __syncthreads();
62
63     if(threadIdx.x == 0) // if it is the master process
64     {
65         *d_sum_power_spectrum = 0; // initialize the total sums to 0
66         *d_sum_freq_power = 0;
67
68         for(int i=0; i<thread_count; i++) // find total sum = sum of partial sums
69         {
70             *d_sum_power_spectrum += partial_sum_power_spectrum[i];
71             *d_sum_freq_power += partial_sum_freq_power[i];
72         }
73     }
74 }
```

```

77 int main()
78 {
79     struct timeval t1, t2;
80
81     double *samples, *d_samples;
82     double *dft_real, *d_dft_real, *dft_img, *d_dft_img;
83     double *k_by_N, *d_k_by_N, *power_spectrum, *d_power_spectrum;
84     double sum_power_spectrum, *d_sum_power_spectrum, sum_freq_power, *d_sum_freq_power;
85
86     samples = (double*)malloc(sizeof(double) * N);
87     dft_real = (double*)malloc(sizeof(double) * N);
88     dft_img = (double*)malloc(sizeof(double) * N);
89     k_by_N = (double*)malloc(sizeof(double) * N);
90     power_spectrum = (double*)malloc(sizeof(double) * N);
91
92     cudaMalloc((void**)&d_samples, sizeof(double) * N);
93     cudaMalloc((void**)&d_dft_real, sizeof(double) * N);
94     cudaMalloc((void**)&d_dft_img, sizeof(double) * N);
95     cudaMalloc((void**)&d_k_by_N, sizeof(double) * N);
96     cudaMalloc((void**)&d_power_spectrum, sizeof(double) * N);
97
98     cudaMalloc((void**)&d_sum_power_spectrum, sizeof(double));
99     cudaMalloc((void**)&d_sum_freq_power, sizeof(double));
100
101     for(int i=0; i<N; i++)
102     {
103         samples[i] = ((float)rand()/(RAND_MAX))*(float)(2.0) - 1;
104         // samples[i] = i;
105         dft_real[i] = 0;
106         dft_img[i] = 0;
107         k_by_N[i] = 0;
108         power_spectrum[i] = 0;
109     }
110
111     gettimeofday(&t1, 0);
112
113     cudaMemcpy(d_samples, samples, sizeof(double) * N, cudaMemcpyHostToDevice);
114     cudaMemcpy(d_dft_real, dft_real, sizeof(double) * N, cudaMemcpyHostToDevice);
115     cudaMemcpy(d_dft_img, dft_img, sizeof(double) * N, cudaMemcpyHostToDevice);
116     cudaMemcpy(d_k_by_N, k_by_N, sizeof(double) * N, cudaMemcpyHostToDevice);
117     cudaMemcpy(d_power_spectrum, power_spectrum, sizeof(double) * N, cudaMemcpyHostToDevice);
118
119     dft<<<1,thread_count>>>(d_samples, d_dft_real, d_dft_img, d_k_by_N, d_power_spectrum, d_sum_power_spectrum, d_sum_freq_power); // function call
120
121     cudaMemcpy(&sum_power_spectrum, d_sum_power_spectrum, sizeof(double), cudaMemcpyDeviceToHost);
122     cudaMemcpy(&sum_freq_power, d_sum_freq_power, sizeof(double), cudaMemcpyDeviceToHost);
123
124     gettimeofday(&t2, 0);
125
126     float normalized_freq = (float)sum_freq_power/sum_power_spectrum;
127
128     printf("The normalized frequency of the random wave is %f\n\n", normalized_freq);
129
130     double time = (1000000.0*(t2.tv_sec-t1.tv_sec) + t2.tv_usec-t1.tv_usec)/1000000.0;
131     printf("time taken : %lf sec\n", time);
132
133     return 0;
134 }

```

// for keeping track of the start and end time
// array which contains the sample values
// array which stores the real and img part of the DFT of the sample
// arrays which store the k/N and power spectrum
// for storing the final sum of the power spectrum and the total frequency

// allocating space for samples in the host
// allocating space for dft_real in the host
// allocating space for dft_img in the host
// allocating space for k by N in the host
// allocating space for power_spectrum in the host

// allocating space for d_samples in the device
// allocating space for d_dft_real in the device
// allocating space for d_dft_img in the device
// allocating space for d_k by N in the device
// allocating space for d_power_spectrum in the device

// allocating space for sum_power_spectrum in the device
// allocating space for sum_freq_power in the device

// assigning random values to samples
// initializing the elements of the other arrays to 0

// obtaining the start time
// copying the host variables to the device variables using cudaMemcpy

// copying the results from the device variables to the host variables
// obtaining the end time
// obtaining the normalized frequency
// calculating the time taken

8. Profiling

Gprof

Profiling tool is gprof from the GNU binutils package. It uses both sampling and instrumentation and gives us a flat profile and call graph which gives various insights in the users code.

Commands used:

- gcc -pg sample.c
- ./a.out
- gprof -b a.out gmon.out

The above commands give us a flat profile and call graph.

Flat profile -

```
dell@bhaavana:~/semester_7/High_performance_computing/Lab/exp_10$ gprof -b ./a.out gmon.out
Flat profile:

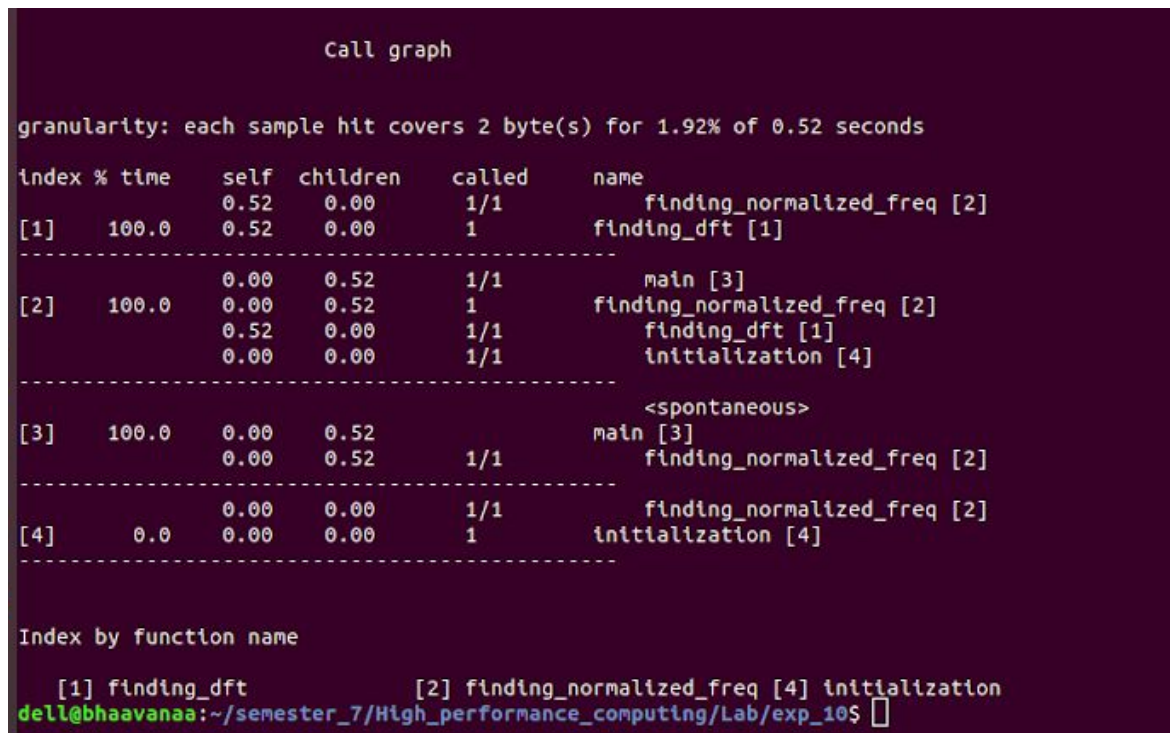
Each sample counts as 0.01 seconds.
%   cumulative   self           calls   ms/call  ms/call  name
time seconds    seconds                   ms/call  ms/call
98.38      0.52      0.52             1      521.40    521.40  finding_dft
 0.00      0.52      0.00             1       0.00    521.40  finding_normalized_freq
 0.00      0.52      0.00             1       0.00     0.00  initialization
```

The flat profile contains information about execution times of all the program's functions and how often they were called.

Inference : By correlating flat profile and call graph, we see that 98% of the time is spent in find_dft() function, (called by the finding_normalized_freq() function, which is called by the main() function), out of the three functions (initialization(), finding_dft() and finding_normalized_freq()), even though it

is called once (which rules out the possibility that more calls take more time). So there can be a potential hotspot inside this function.

Call graph -



Call graph is in depth analysis of flat profile. It basically says what are all the functions which called it and what all functions it calls and other insights.

Inference : We can say that the finding_dft() function seems to be the hotspot in the code and can be a potential bottleneck if the input is large in size.

Likwid

A mini tool suite with various performance analysis measures. LIKWID stands for “*Like I Knew What I’m Doing*”. It is an easy to use yet powerful command line performance tool suite for the GNU/Linux operating system. While the focus of LIKWID is on x86 processors, some of the tools are portable and not limited to any specific architecture.

```
dell@bhaavana:~/semester_7/High_performance_computing/Lab/exp_10$ sudo likwid-perfctr -C 2 -g BRANCH ./a.out
[sudo] password for dell:
-----
CPU name:      Intel(R) Core(TM) i7-7500U CPU @ 2.70GHz
CPU type:      Intel Kabylake processor
CPU clock:     2.90 GHz
-----
The normalized freq of the random wave : 3.141577 Hz
-----
Group 1: BRANCH
+-----+-----+-----+
| Event                | Counter | Core 2 |
+-----+-----+-----+
| INSTR_RETIRED_ANY     | FIXC0   | 32006453996 |
| CPU_CLK_UNHALTED_CORE | FIXC1   | 15865962968 |
| CPU_CLK_UNHALTED_REF  | FIXC2   | 13244435787 |
| BR_INST_RETIRED_ALL_BRANCHES | PMC0   | 3684275188 |
| BR_MISP_RETIRED_ALL_BRANCHES | PMC1   | 21168073 |
+-----+-----+-----+
+-----+-----+
| Metric                | Core 2 |
+-----+-----+
| Runtime (RDTSC) [s]   | 4.5731 |
| Runtime unhaltd [s]   | 5.4649 |
| Clock [MHz]           | 3477.8997 |
| CPI                   | 0.4957 |
| Branch rate           | 0.1151 |
| Branch misprediction rate | 0.0007 |
| Branch misprediction ratio | 0.0057 |
| Instructions per branch | 8.6873 |
+-----+-----+
```

Likwid-perfctr:

Count hardware performance events. It can be used as a wrapper application, which does not require modification of the code to be analyzed, or with a marker API, which restricts the event counting to parts of the code.

So this gives us the overall performance of the program like clocks per instruction, branch rate, branch misprediction rate, branch misprediction ratio.

Inference : We can infer that since branch misprediction rate is not even 1%, there is very less of miss penalty added to the program. Given the case that there is 11% branching happening it is safe to say branch misprediction does not have an effect on the program.