

① Performance of Convolution:-

- CPU with a single core at 2.5 GHz that can perform eight single precision floating-point operations per clock cycle.
- CPU is connected to a DRAM with a peak memory transfer rate at 25.6 GB/s
- Cache size of 256 KB as fast as registers
- Convolution of 2D image I of $n \times n$ size in main memory with a mask M of size $(2K+1) \times (2K+1)$
- Compute an image C of $n \times n$

$$C[x][y] = \sum_{i=-K}^K \sum_{j=-K}^K I[x+i][y+j] \times M[k-i][k-j]$$

for all $0 \leq x, y \leq n-1$

$$I[x+i][y+j] = 0 \quad \text{if } \begin{matrix} x+i < 0 & \text{or} & x+i \geq n \\ y+j < 0 & \text{or} & y+j \geq n \end{matrix}$$

for $x = 0$ to n

for $y = 0$ to n

for $i = -K$ to K

for $j = -K$ to K

if $(x+i) \geq n$ or $(x+i) < 0$ or $(y+j) \geq n$ or $(y+j) < 0$

$I_cell = 0$

else

$I_cell = I[x+i][y+j]$

$C[x][y] += I_cell \times M[k-i][k-j]$

• CPU: ~~3.6 GHz~~ 2.5 GHz

8 single ~~float~~ precision float per cycle

$$\text{FLOPS} = 8 \times 1 \times 2.5 \text{ GHz} = 20 \times 10^9 \text{ FLOPS}$$

• CPU $\Rightarrow 20 \text{ GFLOP/s}$

• DRAM $\Rightarrow 25.6 \text{ GB/s}$

• Cache $\Rightarrow 256 \text{ KB}$

• Size of C $n \times n \Rightarrow 8B \times 256 \times 256 = 512 \text{ KB size of C}$

• 2 ~~B~~ Flops per loop $\times \text{loops} = n \times n \times 2K \times 2K$
 $= 4K^2 n^2 \text{ loops}$

i) $\text{loops} = 4K^2 n^2$
 $= 4 \times 1^2 \times 256^2 = 256K \text{ Loops}$

• Total Flops $= 256K \times 2 = 512K \text{ Flops}$

• CPU computation $= \frac{512K \text{ Flops}}{20 \text{ GFLOPS}} = 24.4 \times 10^{-6} \text{ seconds}$
 $= 24.4 \mu\text{s}$

• Storage $= n \times n \times 8 + K \times K \times 8$
 $= 524296 \text{ Bytes}$
 $\approx 512KB$

• 2 cache cycles $= 2 \times 9.54 \mu\text{s}$
 $= 19.08 \mu\text{s}$

$$t_{\text{exec}} \geq 24.4 \mu\text{s} + 19.1 \mu\text{s} = 43.5 \mu\text{s}$$

$$\frac{512k}{43.5 \mu\text{s}} = 11.22 \text{ GFlops}, \text{ Performance} = \underline{\underline{56.124\%}}$$

$$\textcircled{\text{ii}} \text{ loops} = 4k^2n^2 = 6.25 \text{ Mloops}$$

$$\text{Flops} = 12.5 \text{ MFlops}$$

$$\text{Storage} = 512.2 \text{ KB}$$

$$3 \text{ cache cycle} = 9.54 \times 3 = 28.62 \mu\text{s}$$

$$t_{\text{computation}} = \frac{12.5 \text{ M}}{20 \text{ G}} = 610.35 \mu\text{s}$$

$$t_{\text{exec}} \geq 610.35 + 28.62 \approx 639 \mu\text{s}$$

$$\frac{12.5 \text{ M}}{639 \mu\text{s}} = \cancel{20} 19.10335 \approx 19.1 \text{ GFLOPS}$$

$$\boxed{\text{Performance} = 95\%}$$

$\textcircled{\text{iii}}$ using both Spatial & Temporal Locality

#

② Cache line $L = \text{Size of (float)} \times 16$

NOTES.

$A, B, N \times N, N = 1024$

$C_{ij} = A_{ij} + B_{ij}$ for all i, j

a) Calculate misses of cache

1 miss each new row

Total Misses (1024)

b) 1 miss each cell

Total misses (1024×1024)

$= 1,048,576$ misses

row wise would be faster less cache misses

```

#include <stdio.h>
#include <iostream>
#include <math.h>
#include <pthread.h>
#define N 1024
uint64_t A[N][N];
uint64_t B[N][N];
uint64_t C[N][N];

uint64_t n = 1024;
uint64_t X[1024];

using namespace std;

double Riemann_Zeta(double s, uint64_t k){
    double result = 0.0;
    for (uint64_t i = 1; i < k; i++){
        for (uint64_t j = 1; j < k; j++){
            result += (2 * (i&1)-1)/pow(i+j, s));
        }
    }
    return result*pow(2, s);
}

void *worker_thread_example(void *arg){
    uint64_t k = (uint64_t)arg;
    X[k] = Riemann_Zeta(2,k);
    pthread_exit(NULL);
}

void example(){
    pthread_t threads[n];
    uint64_t Y[1024];

    for(uint64_t k = 0; k < n; k++){
        int ret = pthread_create(&threads[k], NULL, &worker_thread_example, (void*)k);
        if(ret != 0){
            printf("Error: pthread_create() failed\n");
            exit(EXIT_FAILURE);
        }
    }

    for (int i = 0; i < n; i++){
        pthread_join(threads[i], NULL);
    }

    for(uint64_t k = 0; k < n; k++){
        Y[k] = Riemann_Zeta(2, k);
    }

    for(uint64_t k = 0; k < n; k++){
        cout << X[k] << "—" << Y[k] << "=" << X[k]-Y[k] << endl;
    }
}

void row_sum(uint64_t row){
    for(uint64_t i = 0; i < N; i++){
        C[row][i] = A[row][i] + B[row][i];
    }
}

void *worker_thread(void *arg){
    uint64_t row = (uint64_t)arg;
    row_sum(row);
    pthread_exit(NULL);
}

```

main.cpp