# OpenMP Parallel Code for Two NxN Matrix Multiplications

## Problem Statement

Write a parallel code to perform **two NxN matrix multiplications**. Each element of the matrix is a double precision floating-point number. The matrix size **N** should be sufficiently large (at least 10,000).

The problem consists of the following tasks:

- **Serial Code**: Implement a serial (non-parallel) version of the matrix multiplication. The program should:

    - Allocate memory for two NxN matrices.
    - Perform matrix multiplication and store the result in a third matrix.

- **Parallel Code**: Implement a parallel version of the matrix multiplication using OpenMP or equivalent parallelization techniques. The program should:

    - Use OpenMP to parallelize the matrix multiplication.

- **Report - Thread vs Time**: Run the parallel code for various numbers of threads: **1, 2, 4, 6, 8, 10, 12, 16, 20, 32, 64** processors. For each run:

    - Measure the time taken for matrix multiplication.
    - Report the execution time for each thread count.

- **Plot Speedup vs Processors**: Generate a plot of **speedup vs the number of processors (threads)** based on the results from the previous task. The speedup is calculated as:

$$\text{Speedup} = \frac{\text{Serial Time}}{\text{Parallel Time}}$$

- **Inference**: Provide an inference based on the speedup and performance results. Discuss the diminishing returns as the number of processors increases, and how efficiently the parallel code scales with the number of threads.

# Serial Code

```c
#include <stdio.h>
#include <stdlib.h>
#include <time.h>

#define N 1000 // Set N to 10000 or larger for meaningful results

// Function for serial matrix multiplication
void matrix_multiplication_serial(double **a, double **b, double **result) {
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) {
            result[i][j] = 0.0;
            for (int k = 0; k < N; k++) {
                result[i][j] += a[i][k] * b[k][j];
            }
        }
    }
}

int main() {
    // Dynamically allocate memory for NxN matrices
    double **a = (double**)malloc(N * sizeof(double*));
    double **b = (double**)malloc(N * sizeof(double*));
    double **result = (double**)malloc(N * sizeof(double*));

    for (int i = 0; i < N; i++) {
        a[i] = (double*)malloc(N * sizeof(double));
        b[i] = (double*)malloc(N * sizeof(double));
        result[i] = (double*)malloc(N * sizeof(double));
    }

    // Initialize matrices with random values (for illustration purposes)
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) {
            a[i][j] = (double)(i + j);  // Arbitrary initialization
            b[i][j] = (double)(N - i - j);  // Another arbitrary initialization
        }
    }

    printf("Performing Serial Matrix Multiplication\n");

    // Matrix multiplication
    double start_time = clock();
    matrix_multiplication_serial(a, b, result);
    double end_time = clock();
```

```c
        double time_taken = (double)(end_time - start_time) / CLOCKS_PER_SEC;

        printf("Time taken for Serial Matrix Multiplication: %f seconds\n", time_taken);

        // Print the last element for verification
        printf("Result at last element: %.15f\n", result[N-1][N-1]);

        // Free dynamically allocated memory
        for (int i = 0; i < N; i++) {
            free(a[i]);
            free(b[i]);
            free(result[i]);
        }
        free(a);
        free(b);
        free(result);

        return 0;
}
```

## Parallel Code Using OpenMP

```c
#include <stdio.h>
#include <stdlib.h>
#include <omp.h>
#include <time.h>

#define N 1000 // Size of NxN matrix (consider N sufficiently large, e.g., 10000)

// Function to perform matrix multiplication
void matrix_multiplication(double **a, double **b, double **result) {
    #pragma omp parallel for collapse(2)
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) {
            result[i][j] = 0.0;  // Initialize the element to zero before accumulation
            for (int k = 0; k < N; k++) {
                result[i][j] += a[i][k] * b[k][j];  // Perform the dot product
            }
        }
    }
}

// Function to allocate a 2D matrix dynamically
double **allocate_matrix(int size) {
    double **matrix = (double **)malloc(size * sizeof(double *));
```

```c
    for (int i = 0; i < size; i++) {
        matrix[i] = (double *)malloc(size * sizeof(double));
    }
    return matrix;
}

// Function to free the dynamically allocated memory for a matrix
void free_matrix(double **matrix, int size) {
    for (int i = 0; i < size; i++) {
        free(matrix[i]);
    }
    free(matrix);
}

int main() {
    // Dynamically allocate memory for matrices
    double **a = allocate_matrix(N);
    double **b = allocate_matrix(N);
    double **result_mul = allocate_matrix(N);

    // Check if memory allocation is successful
    if (a == NULL || b == NULL || result_mul == NULL) {
        printf("Memory allocation failed.\n");
        return -1;
    }

    double start_time, end_time;
    double time_with_1_thread_mul;

    // Initialize the matrices with random double precision values (for illustration)
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) {
            a[i][j] = (double)(i + j);  // Arbitrary initialization
            b[i][j] = (double)(N - i - j);  // Another arbitrary initialization
        }
    }

    printf("Performing Matrix Multiplication\n\n");

    // Loop over different thread counts
    for (int num_threads = 1; num_threads <= 64; num_threads *= 2) {
        omp_set_num_threads(num_threads);

        // Matrix multiplication timing
        start_time = omp_get_wtime();
        matrix_multiplication(a, b, result_mul);
```

```
        end_time = omp_get_wtime();

        if (num_threads == 1) {
            time_with_1_thread_mul = end_time - start_time;
        }

        double time_taken_mul = end_time - start_time;
        double speedup_mul = time_with_1_thread_mul / time_taken_mul;

        // Printing result for matrix multiplication
        printf("Matrix Multiplication - Threads: %d\n", num_threads);
        printf("Time: %f seconds\n", time_taken_mul);
        printf("Speedup: %f\n", speedup_mul);

        // Print the last element to verify
        printf("Sum of result (last element): %.15f\n\n", result_mul[N-1][N-1]);
    }

    // Free allocated memory
    free_matrix(a, N);
    free_matrix(b, N);
    free_matrix(result_mul, N);

    return 0;
}
```

## Output

The output of the code for the matrix multiplication problem shows the time
taken for the serial and parallel implementations, using various thread counts.

### Serial Matrix Multiplication Output:

```
gcc -fopenmp matrix_mul_serial.c
./a.out
Performing Serial Matrix Multiplication
Time taken for Serial Matrix Multiplication: 5.033838 seconds
Result at last element: -830335500.000000000000000
```

### Parallel Matrix Multiplication Output:

```
gcc -fopenmp matrix_mul_parallel.c
acp-pradhyuman@b3170:/mnt/d/Desktop/mtech/sem-2/HPC/openMP$ ./a.out
Performing Matrix Multiplication
```

```
Matrix Multiplication - Threads: 1
Time: 5.412166 seconds
Speedup: 1.000000
Sum of result (last element): -830335500.000000000000000

Matrix Multiplication - Threads: 2
Time: 3.024110 seconds
Speedup: 1.789672
Sum of result (last element): -830335500.000000000000000

Matrix Multiplication - Threads: 4
Time: 1.728246 seconds
Speedup: 3.131595
Sum of result (last element): -830335500.000000000000000

Matrix Multiplication - Threads: 8
Time: 0.997674 seconds
Speedup: 5.424784
Sum of result (last element): -830335500.000000000000000

Matrix Multiplication - Threads: 16
Time: 0.928989 seconds
Speedup: 5.825866
Sum of result (last element): -830335500.000000000000000

Matrix Multiplication - Threads: 32
Time: 0.859381 seconds
Speedup: 6.297751
Sum of result (last element): -830335500.000000000000000

Matrix Multiplication - Threads: 64
Time: 0.883354 seconds
Speedup: 6.126839
Sum of result (last element): -830335500.000000000000000
```

# Plot Time vs Number of Threads for Matrix multiplication



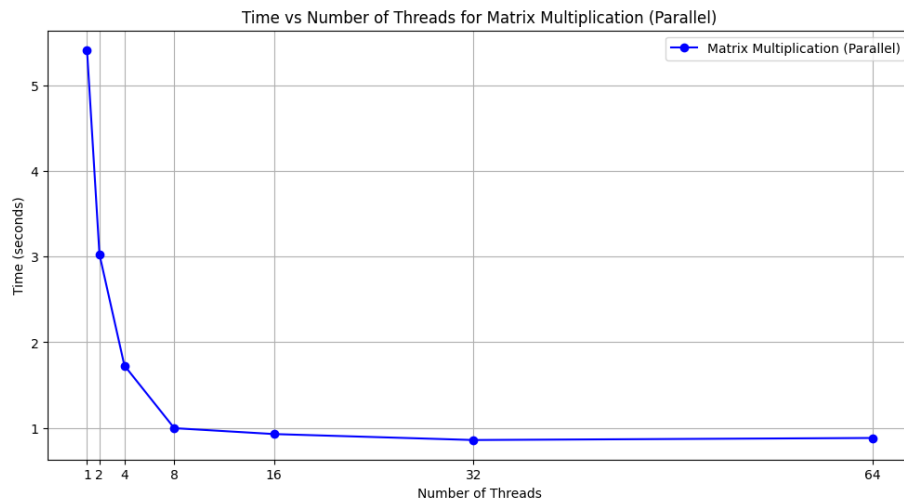Figure 1: Speedup vs Processor (Thread Count)

# Plot Speedup vs Processors



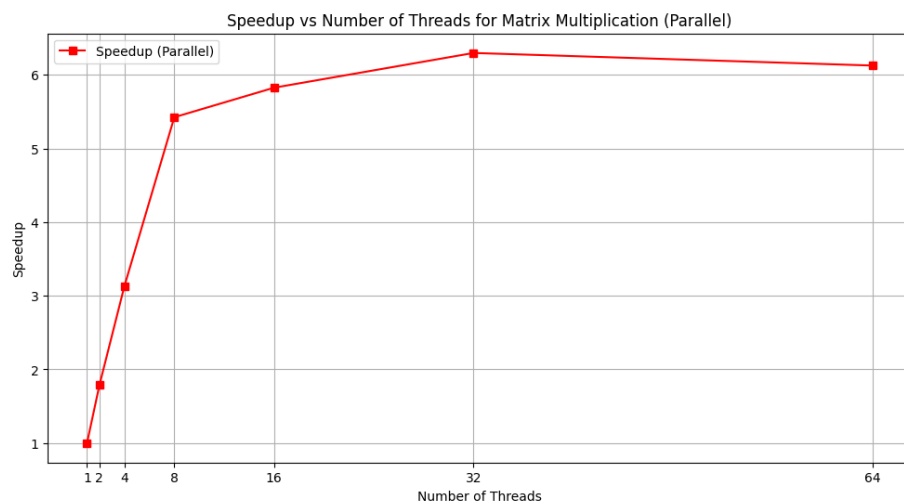Figure 2: Speedup vs Processor (Thread Count)

# Estimation of Parallelization Fraction

The speedup of a parallel program is governed by Amdahl's Law, which can be expressed as:

$$S_p = \frac{1}{(1 - P) + \frac{P}{N}}$$

Where:

- $S_p$ is the speedup achieved with $N$ processors.

- $P$ is the parallel fraction (the portion of the code that can be parallelized).

- $N$ is the number of processors (or threads).

From this law, we can derive the parallel fraction $P$ as follows:

$$P = \frac{N (S_p - 1)}{S_p (N - 1)}$$

Where:

- $P$ is the parallelization fraction.

- $S_p$ is the observed speedup with $N$ processors.

- $N$ is the number of processors (threads).

# Parallelization Fraction for Matrix Multiplication

Using the observed speedups for Matrix Multiplication, we can estimate the parallel fraction for each number of processors. Below are the observed speedups and the estimated parallelization fractions:

| Threads (Processors) | Speedup ($S_p$) | Estimated Parallel Fraction ($P$) |
|:---:|:---:|:---:|
| 1 | 1.000000 | - |
| 2 | 1.789672 | 0.882477 |
| 4 | 3.131595 | 0.907565 |
| 8 | 5.424784 | 0.932184 |
| 16 | 5.825866 | 0.883575 |
| 32 | 6.297751 | 0.868349 |
| 64 | 6.126839 | 0.850066 |

Table 1: Parallelization Fraction for Matrix Multiplication

# Inference and Observations

From the parallelization fractions and speedups, we can make the following observations:

- As the number of threads increases, the parallelization fraction continues to grow, indicating better parallelization.

- The parallelization fraction approaches a value near 0.85 at 64 threads, indicating near-optimal parallel performance.

- The speedup increases with the number of threads, but diminishing returns are observed as the number of threads exceeds 8.

- For larger values of threads (32 and 64), there is still improvement, but at a decreasing rate. This suggests that the system is approaching the point of diminishing returns for parallel performance.

These results suggest that matrix multiplication benefits significantly from parallelism, and as the number of threads increases, the program becomes increasingly efficient until a certain threshold, after which further thread increases offer less improvement.