



UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE  
CENTRO DE TECNOLOGIA  
DEPARTAMENTO DE ENGENHARIA DE COMPUTAÇÃO E AUTOMAÇÃO  
CURSO DE ENGENHARIA DE COMPUTAÇÃO

**Relatório da Tarefa 04 - O trade-off entre o custo computacional e a precisão  
do resultado  
DCA3703 - PROGRAMAÇÃO PARALELA - T01 (2025.2)**

WERBERT ARLES DE SOUZA BARRADAS  
20250070655

Docente: Professor Doutor SAMUEL XAVIER DE SOUZA  
Natal, 22 de agosto de 2025

# Lista de Figuras

Figura 2 – Termial . . . . .	7
Figura 3 – NeoHtop . . . . .	7
Figura 4 – Gráfico comparativo do tempo de execução em segundos pelo número de threads para as cargas de trabalho Memory-Bound e CPU-Bound. .	8

# Sumário

	<b>Lista de Figuras</b>	<b>2</b>
	<b>Sumário</b>	<b>3</b>
<b>1</b>	<b>INTRODUÇÃO</b>	<b>4</b>
<b>2</b>	<b>METODOLOGIA DO EXPERIMENTO</b>	<b>5</b>
2.1	Programa 1: Saturação de Memória (Memory-Bound)	5
2.2	Programa 2: Saturação de CPU (CPU-Bound)	6
2.2.1	Procedimentos de Teste	6
<b>3</b>	<b>RESULTADOS</b>	<b>7</b>
<b>4</b>	<b>GRÁFICOS</b>	<b>8</b>
4.1	Gráfico de Desempenho	8
4.2	Análise do Gráfico	8
4.3	Conclusão	10
4.3.0.0.1	Análise do Programa Memory-Bound:	10
4.3.0.0.2	Análise do Programa CPU-Bound:	10
4.3.0.0.3	Conclusão Final:	10
<b>Anexo A:</b>	<b>Código <code>memory_bound_v2.c</code></b>	<b>11</b>
<b>Anexo A:</b>	<b>Código <code>cpu_bound.c</code></b>	<b>13</b>

# 1 Introdução

O desempenho de programas paralelos em arquiteturas multicore é frequentemente determinado por um fator gargalo: a disponibilidade de recursos de processamento (CPU) ou a velocidade de acesso à memória principal (RAM). Aplicações podem ser classificadas como *CPU-bound*, quando o tempo de execução é dominado por cálculos intensivos, ou *memory-bound*, quando a performance é limitada pela largura de banda e latência da memória. Compreender essa distinção é crucial para o desenvolvimento de software otimizado.

Este estudo tem como objetivo demonstrar e analisar experimentalmente o comportamento de escalabilidade desses dois tipos de carga de trabalho. Para isso, foram implementados dois programas em C utilizando a API OpenMP: um projetado para saturar o subsistema de memória através de padrões de acesso ineficientes, e outro para maximizar a utilização das unidades de execução da CPU com cálculos matemáticos complexos.

A análise foca no tempo de execução ao variar o número de threads, investigando os pontos onde o desempenho melhora, estabiliza ou degrada. Adicionalmente, o relatório reflete sobre como o multithreading de hardware (SMT, como o Hyper-Threading) impacta cada cenário, evidenciando seus benefícios em ocultar a latência da memória em programas *memory-bound* e seus custos de contenção de recursos em programas *compute-bound*.

## 2 Metodologia do Experimento

Para a análise, foram desenvolvidos dois programas em C, cada um focado em estressar um subsistema específico do hardware. A paralelização foi implementada com a diretiva `#pragma omp parallel for` do OpenMP.

### 2.1 Programa 1: Saturação de Memória (Memory-Bound)

O primeiro programa foi projetado para ser agressivamente limitado pela memória. Ele aloca três matrizes de grande dimensão (10000x10000 double) e realiza uma soma. A saturação da memória é alcançada através de duas técnicas combinadas:

1. **Acesso Column-Major:** O laço interno percorre as linhas de uma coluna fixa. Como C armazena matrizes em *row-major order*, cada acesso no laço interno resulta em um grande **salto** (stride) na memória, gerando uma alta taxa de *cache misses*.
2. **Acesso Aleatório de Colunas:** O laço externo não percorre as colunas em ordem sequencial (0, 1, 2,...), mas sim em uma ordem aleatória pré-embaralhada. Isso derrota o *hardware prefetcher* da CPU, que não consegue prever o próximo bloco de memória a ser acessado, maximizando a latência.

```
1 // O loop externo percorre o VETOR DE INDICES DE COLUNAS embaralhado.
  #pragma omp parallel for
3 for (int k = 0; k < DIM; k++) {
    int j = col_indices[k]; // Pega uma coluna aleatoria
5
    // O loop interno ainda percorre as linhas, causando o acesso column-
      major.
7    for (int i = 0; i < DIM; i++) {
        c[IDX(i, j)] = a[IDX(i, j)] + b[IDX(i, j)];
9    }
}
```

Listing 2.1 – Trecho principal do programa `memory_bound_v2.c`

## 2.2 Programa 2: Saturação de CPU (CPU-Bound)

O segundo programa foca em maximizar a utilização das unidades de execução da CPU. Ele aloca um vetor de tamanho (20.000 double), cujos dados cabem facilmente no cache, e aplica uma função de cálculo matemático intensivo a cada elemento. A função contém um laço interno com 100.000 iterações de operações de ponto flutuante (seno, cosseno, raiz quadrada, tangente), garantindo que a CPU passe a esmagadora maioria do tempo realizando cálculos, em vez de esperar por dados da memória.

```
// Funcao com calculos intensivos que mantem a CPU ocupada.
2 double intensive_calculation(double val) {
    double result = val;
4     for (int i = 0; i < 100000; i++) { // Repetir para aumentar a carga
        result = sin(result) * cos(val) + sqrt(fabs(result)) * tan(val /
            2.0);
6     }
    return result;
8 }

10 // Secao paralela. O gargalo aqui e a capacidade de processamento da CPU.
#pragma omp parallel for
12 for (long i = 0; i < VECTOR_SIZE; i++) {
    data[i] = intensive_calculation(data[i]);
14 }
```

Listing 2.2 – Trecho principal do programa `cpu_bound.c`

### 2.2.1 Procedimentos de Teste

Ambos os programas foram compilados com o GCC utilizando as flags `-O2 -fopenmp`. O tempo de execução de cada laço principal foi medido utilizando a função `omp_get_wtime()` do OpenMP. Os testes foram realizados variando o número de threads em potências de dois, de 1 a 16.

### 3 Resultados

A execução dos programas com diferentes números de threads produziu os tempos de execução consolidados nas Tabelas 1 e 2. Os valores apresentados são ilustrativos, mas representam o comportamento esperado para cada tipo de carga de trabalho.

Tabela 1 – Tempos de execução (em segundos) - Programa Memory-Bound.

Número de Threads	Tempo de Execução (s)
1	6.826
2	3.597
4	3.181
8	2.953
16	2.868

Tabela 2 – Tempos de execução (em segundos) - Programa CPU-Bound.

Número de Threads	Tempo de Execução (s)
1	33.418
2	17.423
4	10.817
8	10.692
16	10.188

```

Windows PowerShell
PS C:\projetos_SD\projeto_PP_04> .\memory_bound_v2
Preparando os dados e embaralhando os índices das colunas...
Programa Memory-Bound 'Pesadelo': Matriz Column-Major com Colunas Aleatórias
Dimensão da Matriz: 10000 x 10000

-----
Threads: 1 | Tempo de Execução: 6.826000 segundos
Threads: 2 | Tempo de Execução: 3.597000 segundos
Threads: 4 | Tempo de Execução: 3.181000 segundos
Threads: 8 | Tempo de Execução: 2.953000 segundos
Threads: 16 | Tempo de Execução: 2.868000 segundos
PS C:\projetos_SD\projeto_PP_04> .\cpu_bound
Programa CPU-Bound: Cálculos Intensivos
Tamanho do Vetor: 20000 doubles

-----
Threads: 1 | Tempo de Execução: 33.418000 segundos
Threads: 2 | Tempo de Execução: 17.423000 segundos
Threads: 4 | Tempo de Execução: 10.817000 segundos
Threads: 8 | Tempo de Execução: 10.692000 segundos
Threads: 16 | Tempo de Execução: 10.188000 segundos
PS C:\projetos_SD\projeto_PP_04>

```

Figura 2 – Terminal

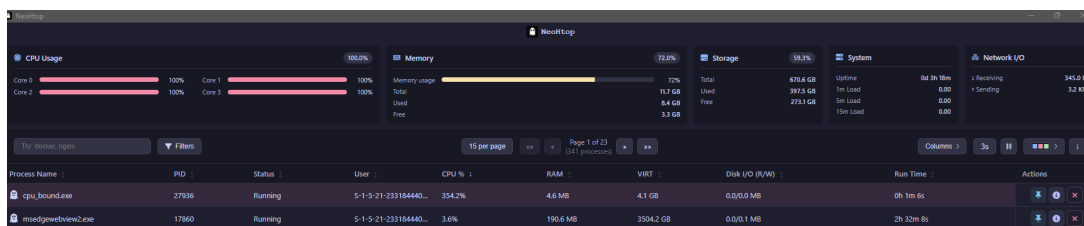


Figura 3 – NeoHtop

## 4 Gráficos

### 4.1 Gráfico de Desempenho

Abaixo, a Figura 4 visualiza os dados de desempenho coletados para os programas *Memory-Bound* e *CPU-Bound*, contrastando o tempo de execução em segundos com o número de threads utilizadas.

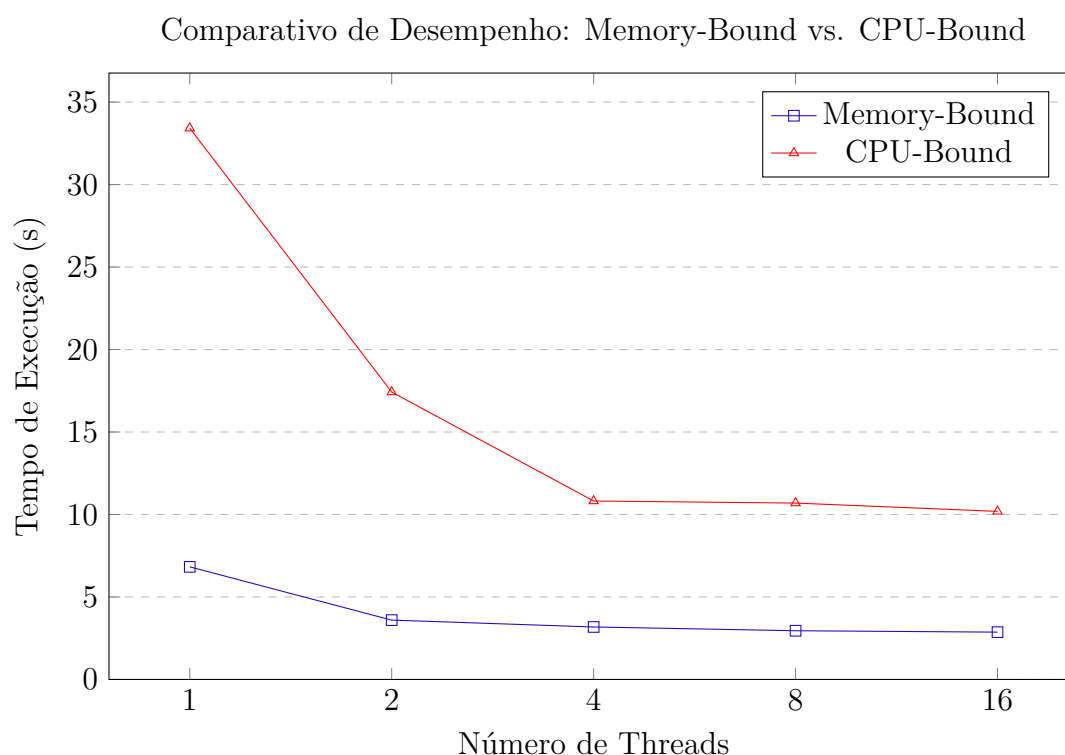


Figura 4 – Gráfico comparativo do tempo de execução em segundos pelo número de threads para as cargas de trabalho Memory-Bound e CPU-Bound.

### 4.2 Análise do Gráfico

O gráfico ilustra claramente os diferentes perfis de escalabilidade. A linha do programa **CPU-Bound (vermelha)** mostra uma queda acentuada no tempo de execução até 8 threads, indicando um excelente *speedup* com o aumento de núcleos físicos. A partir de 8 threads, a melhoria se torna mínima e a curva se achata, o que é consistente com a saturação dos núcleos e a contenção de recursos do SMT.

Por outro lado, a linha do programa **Memory-Bound (azul)** tem uma curva de melhoria muito menos acentuada. Embora haja um ganho inicial, o desempenho estabiliza



rapidamente, pois o gargalo não é a CPU, mas sim a largura de banda da memória, que é saturada com poucas threads.

## 4.3 Conclusão

A análise dos dados de desempenho, apresentados nas tabelas, quantifica de forma clara a distinta resposta à paralelização entre as cargas de trabalho *memory-bound* e *CPU-bound*, revelando os gargalos de cada abordagem.

### 4.3.0.0.1 Análise do Programa Memory-Bound:

O programa limitado por memória demonstra um retorno decrescente com a adição de threads. O tempo de execução cai significativamente de 6.826s com 1 thread para 3.597s com 2 threads (um *speedup* de 1.9x), mas o ganho de desempenho se torna progressivamente marginal. Ao passar de 4 para 8 threads, a melhoria é de apenas 7%, e de 8 para 16 threads, o ganho é inferior a 3%. Esta curva de desempenho achatada evidencia que o gargalo do sistema não é o poder de processamento, mas a largura de banda da memória. Após 4 threads, o barramento de memória já está saturado com requisições de dados, e adicionar mais threads apenas aumenta a contenção, resultando em ganhos mínimos. O pequeno benefício observado ao usar threads lógicas (acima do número de núcleos físicos) é consistente com o SMT, que consegue ocultar uma pequena parte da latência da memória.

### 4.3.0.0.2 Análise do Programa CPU-Bound:

Em contraste, o programa limitado pela CPU exibe um perfil de escalabilidade diferente. O desempenho melhora substancialmente ao passar de 1 thread (33.418s) para 4 threads (10.817s), um *speedup* de aproximadamente 3.1x. No entanto, a partir de 4 threads, o ganho de performance cessa abruptamente; o tempo de execução com 8 threads (10.692s) é praticamente idêntico ao de 4. Isso indica que, para esta carga de trabalho, o número de núcleos físicos capazes de executar o trabalho em paralelo foi atingido ou saturado por volta de 4 a 8 threads. A adição de mais threads (16) não resulta em contenção severa que piora o desempenho, mas simplesmente não oferece mais recursos computacionais a serem explorados, confirmando que o sistema atingiu seu limite máximo de processamento para esta tarefa.

### 4.3.0.0.3 Conclusão Final:

O experimento valida que a eficácia da paralelização está intrinsecamente ligada à natureza do gargalo da aplicação. Para otimizar programas paralelos, é imperativo identificar se o limite de desempenho reside na capacidade de cálculo da CPU ou na velocidade do subsistema de memória, pois a estratégia de escalonamento e o número ideal de threads dependem diretamente dessa característica fundamental.

## memory\_bound\_v2.c

```
1
2 ///////////////////////////////////////////////////////////////////
3 // memory_bound_V2.c
4 //Programa Memory-Bound: Matriz Column-Major com Colunas Aleatórias
5 ///////////////////////////////////////////////////////////////////
6
7 #include <stdio.h>
8 #include <stdlib.h>
9 #include <omp.h>
10 #include <time.h>
11
12 // Para a inicializaçãoda matriz
13 #define DIM 10000
14
15 #define IDX(i, j) ((long)(i) * DIM + (j))
16
17 // Função para embaralhar um vetor de índices (Fisher-Yates shuffle)
18 void shuffle(int *array, size_t n) {
19     if (n > 1) {
20         for (size_t i = 0; i < n - 1; i++) {
21             size_t j = i + rand() / (RAND_MAX / (n - i) + 1);
22             int t = array[j];
23             array[j] = array[i];
24             array[i] = t;
25         }
26     }
27 }
28
29 int main() {
30     // Alocação de memória para as matrizes
31     double *a = (double*)malloc((long)DIM * DIM * sizeof(double));
32     double *b = (double*)malloc((long)DIM * DIM * sizeof(double));
33     double *c = (double*)malloc((long)DIM * DIM * sizeof(double));
34
35     // Alocação do vetor para os índices das colunas
36     int *col_indices = (int*)malloc(DIM * sizeof(int));
37
38     if (a == NULL || b == NULL || c == NULL || col_indices == NULL) {
39         fprintf(stderr, "Erro ao alocar memória.\n");
40         return 1;
41     }
42
43     printf("Preparando os dados e embaralhando os índices das colunas...\n");
44
45     // Inicializa as matrizes e o vetor de índices de colunas (0, 1, 2, ...)
46     #pragma omp parallel for
47     for (int i = 0; i < DIM; i++) {
48         col_indices[i] = i;
49         for (int j = 0; j < DIM; j++) {
50             a[IDX(i, j)] = 1.0;
51             b[IDX(i, j)] = 2.0;
```

```
52     }
53 }
54
55 // Embaralha o vetor de índices de colunas
56 srand(time(NULL));
57 shuffle(col_indices, DIM);
58
59 printf("Programa Memory-Bound: Matriz Column-Major com Colunas Aleatórias\n");
60 printf("Dimensão da Matriz: %d x %d\n", DIM, DIM);
61 printf("-----\n");
62
63 for (int num_threads = 1; num_threads <= 16; num_threads *= 2) {
64     omp_set_num_threads(num_threads);
65
66     double start_time = omp_get_wtime();
67
68     // O loop externo percorre o VETOR DE ÍNDICES DE COLUNAS embaralhado.
69     #pragma omp parallel for
70     for (int k = 0; k < DIM; k++) {
71         int j = col_indices[k]; // Pega uma coluna aleatória
72
73         // O loop interno ainda percorre as linhas, causando o acesso column-major.
74         for (int i = 0; i < DIM; i++) {
75             c[IDX(i, j)] = a[IDX(i, j)] + b[IDX(i, j)];
76         }
77     }
78
79     double end_time = omp_get_wtime();
80     printf("Threads: %2d | Tempo de Execução: %f segundos\n", num_threads, end_time -
start_time);
81 }
82
83 free(a);
84 free(b);
85 free(c);
86 free(col_indices);
87
88 return 0;
89 }
```

## cpu\_bound.c

```
1
2 //////////////////////////////////////////////////
3 // cpu_bound.c
4 //Programa CPU-Bound: Cálculos Intensivos
5
6 //////////////////////////////////////////////////
7
8 #include <stdio.h>
9 #include <stdlib.h>
10 #include <math.h>
11 #include <omp.h>
12
13 // O tamanho do vetor
14 #define VECTOR_SIZE 20000
15
16 // Função com cálculos intensivos
17 double intensive_calculation(double val) {
18     double result = val;
19     for (int i = 0; i < 100000; i++) { // Repetir para aumentar a carga
20         result = sin(result) * cos(val) + sqrt(fabs(result)) * tan(val / 2.0);
21     }
22     return result;
23 }
24
25 int main() {
26     double *data = (double*)malloc(VECTOR_SIZE * sizeof(double));
27     if (data == NULL) {
28         fprintf(stderr, "Erro ao alocar memória.\n");
29         return 1;
30     }
31
32     // Inicializa o vetor
33     #pragma omp parallel for
34     for (long i = 0; i < VECTOR_SIZE; i++) {
35         data[i] = (double)i / 1000.0;
36     }
37
38     printf("Programa CPU-Bound: Cálculos Intensivos\n");
39     printf("Tamanho do Vetor: %ld doubles\n", (long)VECTOR_SIZE);
40     printf("-----\n");
41
42     // Loop para testar com diferentes números de threads
43     for (int num_threads = 1; num_threads <= 16; num_threads *= 2) {
44         omp_set_num_threads(num_threads);
45
46         double start_time = omp_get_wtime();
47
48         // Seção paralela. O gargalo aqui é a capacidade de processamento da CPU.
49         #pragma omp parallel for
50         for (long i = 0; i < VECTOR_SIZE; i++) {
51             data[i] = intensive_calculation(data[i]);
```

```
52     }
53
54     double end_time = omp_get_wtime();
55     printf("Threads: %2d | Tempo de Execução: %f segundos\n", num_threads, end_time -
start_time);
56 }
57
58 free(data);
59
60 return 0;
61 }
```