

Tarefa 15: Otimização de Comunicação em MPI

Escondendo Latência com Sobreposição de Computação

Werbert Arles de Souza Barradas

DCA3703 - Programação Paralela
UFRN - Universidade Federal do Rio Grande do Norte

Setembro de 2025

Objetivo e Estratégia do Experimento

Tese Central

Provar que a **Latência** (τ) da rede pode ser completamente mitigada realizando trabalho útil (computação) simultaneamente.

Metodologia:

- **Problema:** Difusão 1D (Troca de *Halos* a cada passo).
- **Configuração:** $N_{\text{Global}} = 1.000.000$, $\text{STEPS} = 10.000$, Coeficiente de difusão = 0.1.
- **Otimização:** Comparar V1 (Bloqueante), V2(MPI_wait) com V3 (MPI_Test) em cenários de alta concorrência ($P = 8, 16, 32, 64$).

V1 & V2: Bloqueio Total ou Parcial

- ① **V1 (Baseline):** MPI_Send / MPI_Recv.
Processo espera a rede.
- ② **V2 (Wait):** MPI_Isend/Irecv + MPI_Wait.
Overhead de requisição, mas ainda espera bloqueando.

V3: Sobreposição Otimizada

- **Primitivas:** MPI_Isend/Irecv + MPI_Test.
- **Mecanismo:** Inicia comunicação e imediatamente computa os pontos **não dependentes** dos halos, usando o tempo de latência da rede de forma produtiva.

Análise: Início da Eficiência ($P = 8$)

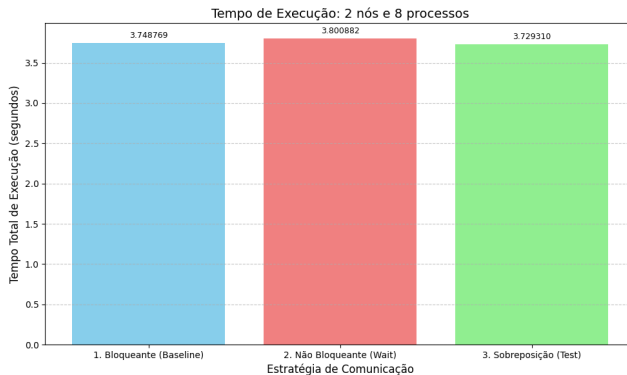


Figura: Gráfico $P = 8$

- $P = 8$: A V3 não é a mais rápida. O T_{comp} é muito longo, e a Latência é insignificante, resultando em *overhead*.

Análise: Início da Eficiência ($P = 8$ vs. $P = 16$)

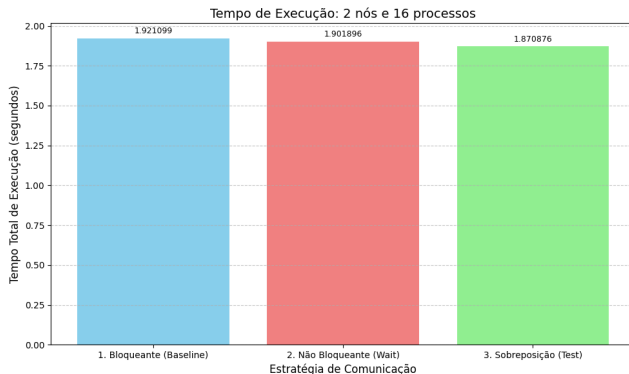


Figura: Gráfico $P = 16$

- $P = 16$: A V3 (1.070876 s) se torna, marginalmente, a melhor. A redução da carga local permite o início da **sobreposição efetiva**.

Prova Final: Domínio da Latência ($P = 32$)

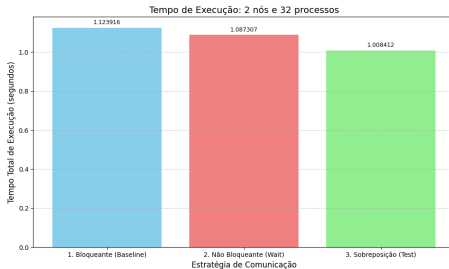


Figura: Desempenho no Cenário de Máxima Concorrência ($P = 32$)

Ganho da Sobreposição

- **V1 (Baseline):** 1.123916 s
- **V3 (Otimizada):** 1.008412 s

Ganho Total: $\approx 10.28\%$

- **Princípio Validado:** O ganho de **10.28%** demonstra que, para otimizar a escalabilidade, é preciso fazer a latência da rede coincidir com o tempo de computação útil.
- **Condição Crítica:** A otimização máxima (**V3**) só se manifestou no cenário de **Latência Dominante** ($P = 32$), onde o trabalho de computação local é baixo e o custo da comunicação é alto.

difusao_bloqueante.c

```
1  #include <mpi.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <math.h>
5
6  // Parâmetros da Simulação
7  #define GLOBAL_N 100000 // Tamanho total da barra
8  #define STEPS 5000      // Número de passos de tempo
9  #define ALPHA 0.1       // Coeficiente de difusão (precisa ser < 0.5 para
                          // estabilidade)
10
11
12 #define TAG_LEFT_TO_RIGHT 0
13
14 #define TAG_RIGHT_TO_LEFT 1
15
16 void compute_inner(double* u_new, double* u, int size) {
17
18     for (int i = 1; i < size - 1; i++) {
19
20         u_new[i] = u[i] + ALPHA * (u[i-1] - 2.0 * u[i] + u[i+1]);
21     }
22 }
23
24 int main(int argc, char** argv) {
25     MPI_Init(&argc, &argv);
26
27     int rank, size;
28     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
29     MPI_Comm_size(MPI_COMM_WORLD, &size);
30
31     if (size < 2) {
32         fprintf(stderr, "Este programa requer pelo menos 2 processos.\n");
33         MPI_Finalize();
34         return 1;
35     }
36
37
38     int local_data_size = GLOBAL_N / size;
39
40     int local_size = local_data_size + 2;
41
42     // Alocação dos arrays: u (atual) e u_new (próxima iteração)
43     double* u = (double*)calloc(local_size, sizeof(double));
44     double* u_new = (double*)calloc(local_size, sizeof(double));
45
46
47     int left = (rank > 0) ? rank - 1 : MPI_PROC_NULL;
48     int right = (rank < size - 1) ? rank + 1 : MPI_PROC_NULL;
49
```



```
50
51     if (rank == 0) {
52         // Inicializa uma seção com um valor alto para simular calor
53         for(int i = 1; i < local_data_size/2; i++) {
54             u[i] = 10.0;
55         }
56     }
57
58     // --- Loop Principal ---
59     double start_time = MPI_Wtime();
60
61     for (int t = 0; t < STEPS; t++) {
62
63         if (right != MPI_PROC_NULL) {
64             MPI_Send(&u[local_size - 2], 1, MPI_DOUBLE, right,
65 TAG_RIGHT_TO_LEFT, MPI_COMM_WORLD);
66         }
67
68         if (right != MPI_PROC_NULL) {
69             MPI_Recv(&u[local_size - 1], 1, MPI_DOUBLE, right,
70 TAG_LEFT_TO_RIGHT, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
71         }
72
73         if (left != MPI_PROC_NULL) {
74             MPI_Send(&u[1], 1, MPI_DOUBLE, left, TAG_LEFT_TO_RIGHT,
75 MPI_COMM_WORLD);
76         }
77
78         if (left != MPI_PROC_NULL) {
79             MPI_Recv(&u[0], 1, MPI_DOUBLE, left, TAG_RIGHT_TO_LEFT,
80 MPI_COMM_WORLD, MPI_STATUS_IGNORE);
81         }
82
83         // --- FIM DA COMUNICAÇÃO BLOQUEANTE ---
84
85         compute_inner(u_new, u, local_size);
86
87         // 3. Trocar Ponteiros para o próximo passo de tempo
88         double *temp = u;
89         u = u_new;
90         u_new = temp;
91     }
92
93     double total_time = MPI_Wtime() - start_time;
94
95     if (rank == 0) {
96         printf("Versao 1 (Bloqueante - Send/Recv) | N=%d, STEPS=%d: %.6f s\n",
97 GLOBAL_N, STEPS, total_time);
98     }
```

```
98  
99     // Limpeza  
100    free(u);  
101    free(u_new);  
102    MPI_Finalize();  
103    return 0;  
104 }  
105
```

difusão_Nao_bloqueante_wait.c

```
1  #include <mpi.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <math.h>
5
6  // Parâmetros da Simulação (Ajuste para o seu teste de desempenho)
7  #define GLOBAL_N 1000000
8  #define STEPS 10000
9  #define ALPHA 0.1
10
11 #define TAG_LEFT_TO_RIGHT 0
12 #define TAG_RIGHT_TO_LEFT 1
13
14
15 void compute_inner(double* u_new, double* u, int size) {
16
17     for (int i = 1; i < size - 1; i++) {
18
19         u_new[i] = u[i] + ALPHA * (u[i-1] - 2.0 * u[i] + u[i+1]);
20     }
21 }
22
23 int main(int argc, char** argv) {
24     MPI_Init(&argc, &argv);
25
26     int rank, size;
27     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
28     MPI_Comm_size(MPI_COMM_WORLD, &size);
29
30     if (size < 2) {
31         if (rank == 0) fprintf(stderr, "Este programa requer pelo menos 2
32 processos.\n");
33         MPI_Finalize();
34         return 1;
35     }
36
37     int local_data_size = GLOBAL_N / size;
38     int local_size = local_data_size + 2;
39
40     double* u = (double*)calloc(local_size, sizeof(double));
41     double* u_new = (double*)calloc(local_size, sizeof(double));
42
43     int left = (rank > 0) ? rank - 1 : MPI_PROC_NULL;
44     int right = (rank < size - 1) ? rank + 1 : MPI_PROC_NULL;
45
46     // Declarar Request e Status para comunicação não bloqueante
47     MPI_Request requests[4];
48     MPI_Status status;
49
50     // Inicialização (Ponto quente no primeiro processo)
```

```
50     if (rank == 0) {
51         for(int i = 1; i < local_data_size/2; i++) {
52             u[i] = 10.0;
53         }
54     }
55
56     MPI_Barrier(MPI_COMM_WORLD);
57     double start_time = MPI_Wtime();
58
59     for (int t = 0; t < STEPS; t++) {
60
61         // 1. Inicia Comunicação Não Bloqueante (4 chamadas)
62         // [0] e [2]: ISend (saídas) | [1] e [3]: IRecv (entradas)
63
64         // Envio/Recebimento na Direita
65         MPI_Isend(&u[local_size - 2], 1, MPI_DOUBLE, right, TAG_RIGHT_TO_LEFT,
66 MPI_COMM_WORLD, &requests[0]);
67         MPI_Irecv(&u[local_size - 1], 1, MPI_DOUBLE, right, TAG_LEFT_TO_RIGHT,
68 MPI_COMM_WORLD, &requests[1]);
69
70         // Envio/Recebimento na Esquerda
71         MPI_Isend(&u[1], 1, MPI_DOUBLE, left, TAG_LEFT_TO_RIGHT,
72 MPI_COMM_WORLD, &requests[2]);
73         MPI_Irecv(&u[0], 1, MPI_DOUBLE, left, TAG_RIGHT_TO_LEFT,
74 MPI_COMM_WORLD, &requests[3]);
75
76         // Espera pelos ISend (índices 1 e 3)
77         MPI_Wait(&requests[1], &status);
78         MPI_Wait(&requests[3], &status);
79
80         // Espera pelos ISend (índices 0 e 2)
81         MPI_Wait(&requests[0], &status);
82         MPI_Wait(&requests[2], &status);
83
84         // 3. Computação Interna (Apenas após o Wait/chegada das bordas)
85         compute_inner(u_new, u, local_size);
86
87         // 4. Trocar Ponteiros
88         double *temp = u;
89         u = u_new;
90         u_new = temp;
91     }
92
93     double total_time = MPI_Wtime() - start_time;
94     MPI_Barrier(MPI_COMM_WORLD);
95
96     if (rank == 0) {
97         printf("Versao 2 (Nao Bloqueante - Wait): %.6f s\n", total_time);
98     }
99
100     free(u);
101     free(u_new);
102     MPI_Finalize();
```

```
99 |         return 0;  
100 |     }  
101 |
```

difusão_Nao_bloqueante_test.c

```
1  #include <mpi.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4  #include <math.h>
5
6  // Parâmetros da Simulação (Os mesmos para todas as versões)
7  #define GLOBAL_N 1000000
8  #define STEPS 10000
9  #define ALPHA 0.1
10
11 #define TAG_LEFT_TO_RIGHT 0
12 #define TAG_RIGHT_TO_LEFT 1
13
14
15 void compute_inner(double* u_new, double* u, int size) {
16
17     for (int i = 1; i < size - 1; i++) {
18         u_new[i] = u[i] + ALPHA * (u[i-1] - 2.0 * u[i] + u[i+1]);
19     }
20 }
21
22 int main(int argc, char** argv) {
23     MPI_Init(&argc, &argv);
24
25     int rank, size;
26     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
27     MPI_Comm_size(MPI_COMM_WORLD, &size);
28
29     if (size < 2) {
30         if (rank == 0) fprintf(stderr, "Este programa requer pelo menos 2
31 processos.\n");
32         MPI_Finalize();
33         return 1;
34     }
35
36     int local_data_size = GLOBAL_N / size;
37     int local_size = local_data_size + 2;
38
39     double* u = (double*)calloc(local_size, sizeof(double));
40     double* u_new = (double*)calloc(local_size, sizeof(double));
41
42     int left = (rank > 0) ? rank - 1 : MPI_PROC_NULL;
43     int right = (rank < size - 1) ? rank + 1 : MPI_PROC_NULL;
44
45     MPI_Request requests[4];
46     MPI_Status status;
47
48
49     int inner_overlap_start = 2;
```

```
50     int inner_overlap_end = local_size - 3;
51
52
53     if (rank == 0) {
54         for(int i = 1; i < local_data_size/2; i++) {
55             u[i] = 10.0;
56         }
57     }
58
59     MPI_Barrier(MPI_COMM_WORLD);
60     double start_time = MPI_Wtime();
61
62     for (int t = 0; t < STEPS; t++) {
63
64
65
66         // Envio/Recebimento na Direita
67         MPI_Isend(&u[local_size - 2], 1, MPI_DOUBLE, right, TAG_RIGHT_TO_LEFT,
68 MPI_COMM_WORLD, &requests[0]);
69         MPI_Irecv(&u[local_size - 1], 1, MPI_DOUBLE, right, TAG_LEFT_TO_RIGHT,
70 MPI_COMM_WORLD, &requests[1]);
71
72         // Envio/Recebimento na Esquerda
73         MPI_Isend(&u[1], 1, MPI_DOUBLE, left, TAG_LEFT_TO_RIGHT,
74 MPI_COMM_WORLD, &requests[2]);
75         MPI_Irecv(&u[0], 1, MPI_DOUBLE, left, TAG_RIGHT_TO_LEFT,
76 MPI_COMM_WORLD, &requests[3]);
77
78         // 2. Computação da Zona Interna (SOBREPOSIÇÃO)
79         // 0 processador agora calcula os pontos internos que não dependem da
80 comunicação.
81         for (int i = inner_overlap_start; i <= inner_overlap_end; i++) {
82             u_new[i] = u[i] + ALPHA * (u[i-1] - 2.0 * u[i] + u[i+1]);
83         }
84
85         int flag = 0;
86
87         while (!flag) {
88
89             int flag_recv_right = 0;
90             int flag_recv_left = 0;
91
92             MPI_Test(&requests[1], &flag_recv_right, &status);
93             MPI_Test(&requests[3], &flag_recv_left, &status);
94
95             flag = flag_recv_right && flag_recv_left;
96
97         }
98
99         u_new[1] = u[1] + ALPHA * (u[0] - 2.0 * u[1] + u[2]);
100         u_new[local_size - 2] = u[local_size - 2] + ALPHA * (u[local_size - 3]
101 - 2.0 * u[local_size - 2] + u[local_size - 1]);
102     }
```

```
97         // Certifica-se que os envios também terminaram antes do próximo passo
(Importante para o buffer)
98         MPI_Wait(&requests[0], &status);
99         MPI_Wait(&requests[2], &status);
100
101         // 5. Trocar Ponteiros
102         double *temp = u;
103         u = u_new;
104         u_new = temp;
105     }
106
107     double total_time = MPI_Wtime() - start_time;
108     MPI_Barrier(MPI_COMM_WORLD);
109
110     if (rank == 0) {
111         printf("Versao 3 (Sobreposicao - Test): %.6f s\n", total_time);
112     }
113
114     free(u);
115     free(u_new);
116     MPI_Finalize();
117     return 0;
118 }
119
```