

Lista Programação Concorrente e distribuída

Resolução feita por Samuel Cavalcanti

Samuel Cavalcanti

Copyright © 2022 Samuel Cavalcanti

Table of contents

1. Lista de Exercícios de Programação Concorrente e distribuída	3
2. Chapter 1	4
2.1 Capítulo 1	4
2.2 Question 1	5
2.3 Question 2	6
2.4 Question 3	8
2.5 Question 4	10
2.6 Question 5	11
2.7 Question 6	14
2.8 Question 9	16
3. Chapter 2	18
3.1 Capítulo 2	19
3.2 Question 01	22
3.3 Question 02	25
3.4 Question 03	26
3.5 Question 05	28
3.6 Question 07	29
3.7 Question 10	31
3.8 Question 15	33
3.9 Question 16	34
3.10 Question 17	38
3.11 Question 19	39
3.12 Question 20	41
3.13 Question 21	42
3.14 Question 24	43
4. Chapter 3	47
4.1 Capítulo 3	48
4.2 Question 02	52

1. Lista de Exercícios de Programação Concorrente e distribuída

- ✓ capítulo 1 (7 questões)
- ✓ capítulo 2 (13 questões)
- capítulo 3 (16 questões)
- capítulo 5 (12 questões)

2. Chapter 1

2.1 Capítulo 1

Capítulo 1: 1-6, 9; (7 questões)

1. Devise formulas for the functions that calculate ***my_first_i*** and ***my_last_i*** in the global sum example. Remember that each core should be assigned roughly the same number of elements of computations in the loop. Hint: First consider the case when ***n*** is evenly divisible by ***p***
 - [Resposta Questão 1](#)
2. We've implicitly assumed that each call to ***Compute_next_value*** requires roughly the same amount of work as the other calls. How would you change your answer to the preceding question if call $i = k$ requires $k + 1$ times as much work as the call with $i = 0$? So if the first call ($i = 0$) requires 2 milliseconds, the second call ($i = 1$) requires 4, the third ($i = 2$) requires 6, and so on.
 - [Resposta Questão 2](#)
3. Try to write pseudo-code for the tree-structured global sum illustrated in Figure 1.1. Assume the number of cores is a power of two (1, 2, 4, 8, ...).
 - [Resposta Questão 3](#)
4. As an alternative to the approach outlined in the preceding problem we can use C's bitwise operators to implement the tree-structured global sum. In order to see how this works, it helps to write down the binary (base 2) representation of each of the core ranks, and note the pairings during each stage
 - [Resposta Questão 4](#)
5. What happens if your pseudo-code in Exercise 1.3 or Exercise 1.4 is run when the number of cores is not a power of two (e.g., 3, 5, 6, 7)? Can you modify the pseudo-code so that it will work correctly regardless of the number of cores?
 - [Resposta Questão 5](#)
6. Derive formulas for the number of receives and additions that core 0 carries out using:
 - a. the original pseudo-code for a global sum
 - b. the tree-structured global sum.
 Make a table showing the numbers of receives and additions carried out by core 0 when the two sums are used with 2, 4, 8, ..., 1024 cores.
 - [Resposta Questão 6](#)
- 7.
- 8.
9. Write an essay describing a research problem in your major that would benefit from the use of parallel computing. Provide a rough outline of how parallelism would be used. Would you use task- or data-parallelism?
 - [Resposta Questão 9](#)

2.2 Question 1

2.2.1 Questão 1

Devise formulas for the functions that calculate ***my_first_i*** and ***my_last_i*** in the global sum example. Remember that each core should be assigned roughly the same number of elements of computations in the loop. Hint: First consider the case when ***n*** is evenly divisible by ***p***

```
struct range
{
    int first;
    int last;
};
struct range new_range(int thread_index, int p, int n)
{
    struct range r;

    int division = n / p;
    int rest = n % p;

    if (rest == 0)
    {
        r.first = thread_index * division;
        r.last = (thread_index + 1) * division;
    }
    else
    {
        r.first = thread_index == 0 ? 0 : thread_index * division + rest;
        r.last = (thread_index + 1) * division + rest;
    }

    if (r.last > n)
        r.last = n;

    return r;
}

struct range new_range_2(int thread_index, int p, int n)
{
    struct range r;

    int division = n / p;
    int rest = n % p;

    if (thread_index < rest)
    {
        r.first = thread_index * (division + 1);
        r.last = r.first + division + 1;
    }
    else
    {
        r.first = thread_index * division + rest;
        r.last = r.first + division;
    }

    return r;
}
```

Saída:

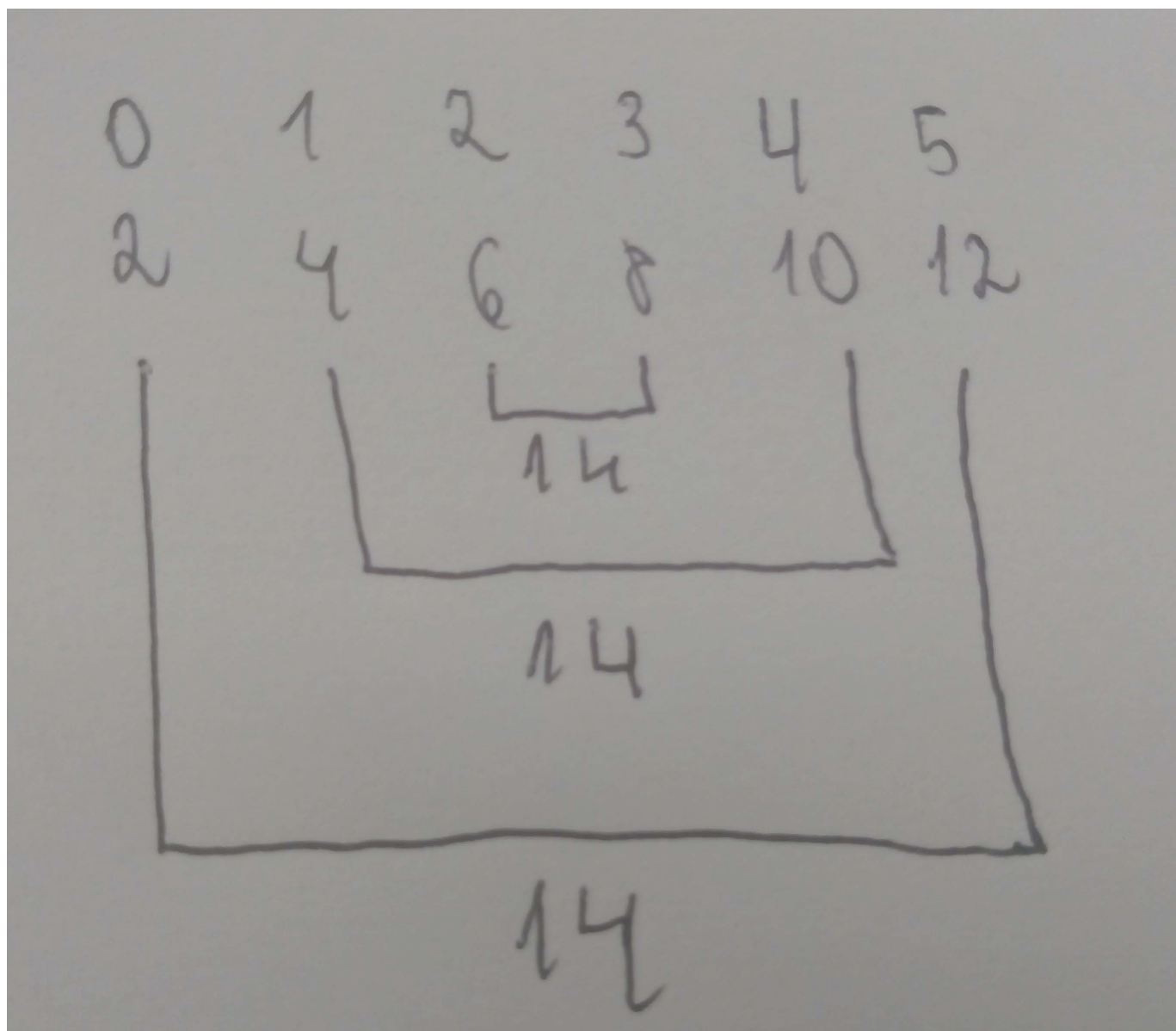
```
First 0 Last 20 m Last - First: 20
First 20 Last 40 m Last - First: 20
First 40 Last 60 m Last - First: 20
First 60 Last 80 m Last - First: 20
First 80 Last 100 m Last - First: 20
First 0 Last 25 m Last - First: 25
First 25 Last 50 m Last - First: 25
First 50 Last 75 m Last - First: 25
First 75 Last 99 m Last - First: 24
First 99 Last 123 m Last - First: 24
Test question 1 success
OLD new range
First 0 Last 27 m Last - First: 27
First 27 Last 51 m Last - First: 24
First 51 Last 75 m Last - First: 24
First 75 Last 99 m Last - First: 24
First 99 Last 123 m Last - First: 24
```

2.3 Question 2

2.3.1 Questão 2

We've implicitly assumed that each call to **Compute_next_value** requires roughly the same amount of work as the other calls. How would you change your answer to the preceding question if call $i = k$ requires $k + 1$ times as much work as the call with $i = 0$? So if the first call ($i = 0$) requires 2 milliseconds, the second call ($i = 1$) requires 4, the third ($i = 2$) requires 6, and so on.

Somatório de Gauss



Exemplo, Supondo que $k = 10$, temos o seguinte vetor de índice:

indices = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

utilizando a lógica da somatório de gauss e organizando os indices temos um array normalizado:

normalized_array = [[0, 9], [1, 8], [2, 7], [3, 6], [4, 5]]

onde o custo de cada índice do `normalized_array` é igual, por tanto podemos usar o algoritmo da questão 1 aplicado ao `normalized_array` resultando:

Thread	normalized_array
1	0,1
2	2,3
3	4

Thread	Compute_next_value	cost
1	0, 9,1, 8	44
2	2, 7,3, 6	44
3	4, 5	22

2.4 Question 3

2.4.1 Questão 3

Try to write pseudo-code for the tree-structured global sum illustrated in Figure 1.1. Assume the number of cores is a power of two (1, 2, 4, 8, ...).

Para criar a árvore, foi considerado que o vetor principal já foi igualmente espaçado entre as p threads, usando o algoritmo da questão 1.

Representação de um Nó

Neste caso foi representado o nó, como uma estrutura que possui um vetor de vizinhos e outro ponteiro para um vetor de inteiros, na prática, o ponteiro para um vetor de inteiros, seria usando o design pattern chamado **Future**, ou um **Option**\<Future>.

Também foi criado dois construtores um construtor que representa, a inicialização do Nó por meio dos seus vizinhos a esquerda e direita, usando na criação de nós intermediários da árvore, e a inicialização do Nó por meio do seu dado, os nós inicializados por meio dos dados abstrai os **núcleos** ou as **threads** que estão sendo usadas para executar o algoritmo de alto custo computacional.

```
class Node
{
public:
    std::vector<Node *> neighborhoods;
    std::vector<int> *data;
    Node(Node *left, Node *right)
    {
        this->neighborhoods.push_back(left);
        this->neighborhoods.push_back(right);
        this->data = nullptr;
    }
    Node(std::vector<int> *data)
    {
        this->data = data;
    }
    ~Node()
    {
        delete this->data;
    }
};
```

Função que cria a Árvore

Para criar a Árvore foi feita uma função recursiva que a partir do nível mais baixo da árvore cria a raiz, ou seja, a partir um vetor com p Nós ,a função vai sendo chamada recursivamente, onde a cada chamada vai-se criando um nível acima da árvore, até que se atinja a raiz, onde a cada nível o número de nós é dividido por 2. **Caso o número de nós inicial não for divisível por 2, o algoritmo não funciona**

```
std::vector<Node *> create_tree_from_core_nodes(std::vector<Node *> nodes)
{
    auto size = nodes.size();

    if (size / 2 == 1)
    {
        auto left = nodes[0];
        auto right = nodes[1];
        receive_value(left, right); // Left receive value from right
        return {left};
    }

    auto new_nodes = std::vector<Node *>{};

    for (auto i = 0; i < size; i += 2)
    {
        auto left = nodes[i];
        auto right = nodes[i + 1];
        receive_value(left, right); // Left receive value from right
        new_nodes.push_back(left);
    }

    return create_tree_from_core_nodes(new_nodes);
}

Node *create_new_tree(std::vector<Node *> nodes)
{
    return create_tree_from_core_nodes(nodes)[0];
}
```


Após criar a árvore basta percorrer-lá recursivamente, lembrando que na prática ***compute_data***, seria um *join*, ou um *await*, de uma thread.

```
int compute_data(std::vector<int> *data)
{
    auto total = 0;
    auto size = data->size();
    for (auto i = 0; i < size; i++)
    {
        total += data->at(i);
    }

    return total;
}

int compute_node(Node &node)
{
    int result_data = node.data == nullptr ? 0 : compute_data(node.data);

    for (auto neighborhood : node.neighborhoods)
        result_data += compute_node(*neighborhood);

    return result_data;
}
```

2.5 Question 4

2.5.1 Questão 4

As an alternative to the approach outlined in the preceding problem we can use C's bitwise operators to implement the tree-structured global sum. In order to see how this works, it helps to write down the binary (base 2) representation of each of the core ranks, and note the pairings during each stage

Semelhante ao questão 3 sendo a diferença utilizar o bitwise << para dividir o tamanho atual da função recursiva:

```
std::vector<Node *> create_new_tree_bitwise(std::vector<Node *> nodes)
{
    auto size = nodes.size();

    if (size >> 1 == 1) // alteração.
    {
        auto left = nodes[0];
        auto right = nodes[1];
        receive_value(left, right);
        return {left};
    }

    auto new_nodes = std::vector<Node *>{};

    for (auto i = 0; i < size; i += 2)
    {
        auto left = nodes[i];
        auto right = nodes[i + 1];
        receive_value(left, right); // Left receive value from right
        new_nodes.push_back(left);
    }

    return create_new_tree_bitwise(new_nodes);
}
```

2.6 Question 5

2.6.1 Questão 5

What happens if your pseudo-code in Exercise 1.3 or Exercise 1.4 is run when the number of cores is not a power of two (e.g., 3, 5, 6, 7)? Can you modify the pseudo-code so that it will work correctly regardless of the number of cores?

Entendendo o que acontece quando os número de núcleos não é da potencia de 2

Se por exemplo, o número de cores, ou Nós for 3 por exemplo, existirá nós que serão "esquecidos" no algoritmo, por tanto o algoritmo não funcionará corretamente.

```
auto size = nodes.size(); // size = 3

if (size >> 1 == 1) // alteração.
{
    auto left = nodes[0];
    auto right = nodes[1];
    // node[2] foi esquecido
    receive_value(left, right); // Left receive value from right
    return {left};
}
```

ou se o por exemplo o número de nós for 7, a última iteração do laço **for**, os índices **i**, **i+1**, serão respectivamente 6 e 7, ou seja será acessado um endereço inválido 7, uma vez que os índices vão de 0 até 6.

```
auto new_nodes = std::vector<Node *>{};

for (auto i = 0; i < size; i += 2)
{
    auto left = nodes[i];
    auto right = nodes[i + 1];
    receive_value(left, right); // Left receive value from right
    new_nodes.push_back(left);
}

return create_new_tree_bitwise(new_nodes);
```

Para isso foram feitas as seguintes modificações:

- adicionado condicionamento para verificar se o tamanho é igual 3 - alterado o bitwise para apenas um comparador **size == 2** - verificado se o tamanho dos nós é par, caso não seja adicionado uma logica extra.

```
std::vector<Node *> create_new_tree_bitwise(std::vector<Node *> nodes)
{
    auto size = nodes.size();

    if (size == 2)
    {
        auto left = nodes[0];
        auto right = nodes[1];
        receive_value(left, right);
        return {left}; // Construtor C++ Moderno.
    }

    if (size == 3)
    {
        auto left = nodes[0];
        auto middle = nodes[1];
        auto right = nodes[2];
        receive_value(left, middle); // Left receive value from middle
        receive_value(left, right); // Left receive value from right
        return {left}; // Construtor C++ Moderno.
    }

    auto new_nodes = std::vector<Node *>{};

    if (size % 2 != 0) // lógica extra.
    {
        size = size - 1;
        new_nodes.push_back(nodes[size]);
    }

    for (auto i = 0; i < size; i += 2)
    {
        auto left = nodes[i];
        auto right = nodes[i + 1];
        receive_value(left, right); // Left receive value from right
    }

    return new_nodes;
}
```

```

        new_nodes.push_back(left);
    }

    return create_new_tree_bitwise(new_nodes);
}

```

Explicando a lógica extra

Além de adicionar uma verificação para saber se o tamanho é par, foi adicionado dois comandos extras, o primeiro é alterar o tamanho (size), para um valor menor, uma vez que estávamos acessando um índice maior que o permitido. Segundo foi adicionar o nó que não será percorrido pelo laço para o vetor ***new_nodes*** que será a entrada da próxima função recursiva

```

if (size % 2 != 0) // verificação se é par.
{
    size = size - 1; //1
    new_nodes.push_back(nodes[size]); // 2
}

```

Explicando o novo caso base

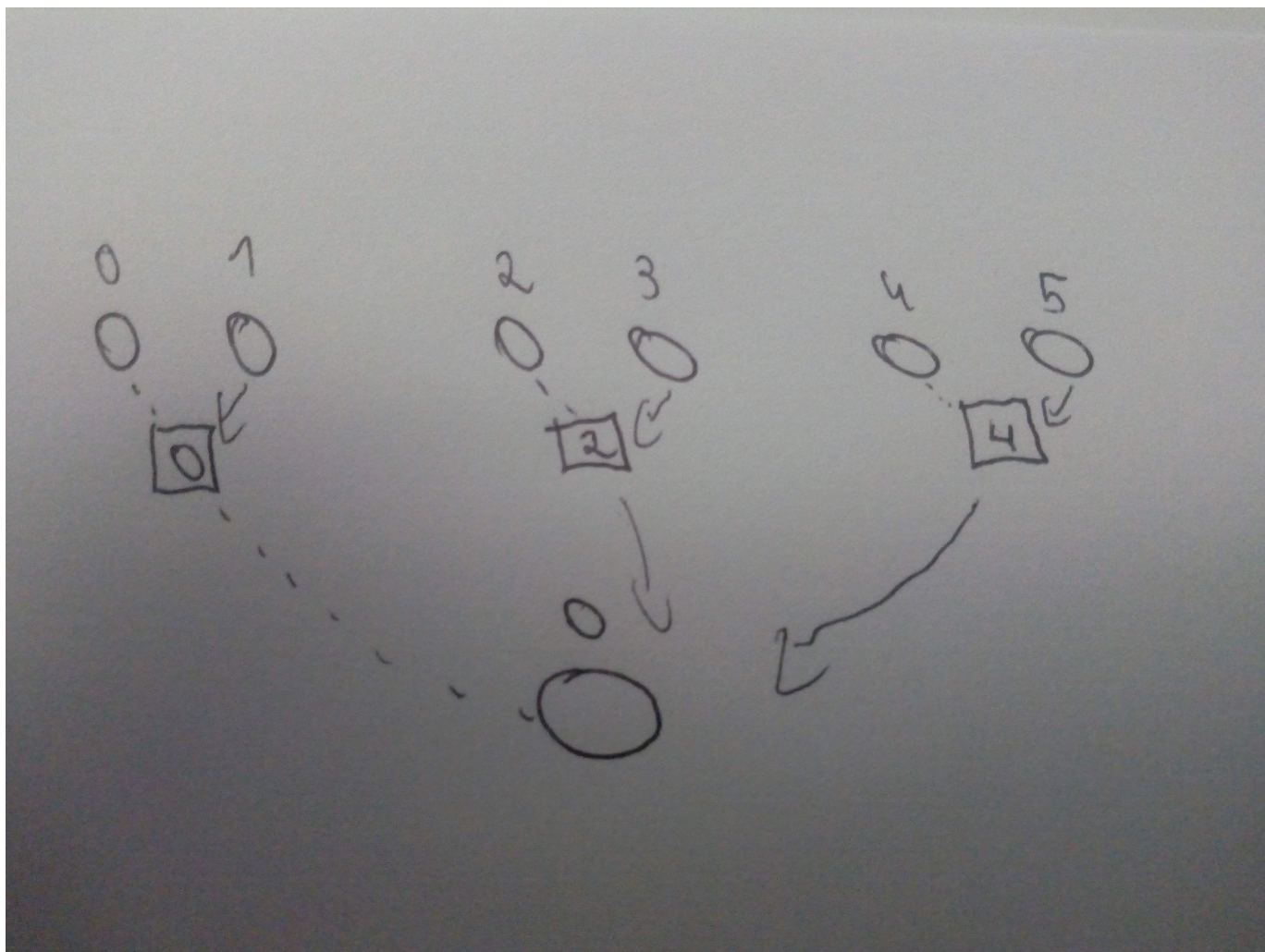
Percebemos que além do $2/2 == 1$, a divisão inteira de $3/2$ também é igual 1. Por tanto além do caso base de quando o tamanho do vetor de nós ser igual a 2, temos que tratar também quando o número de nós ser igual a 3.

```

if (size == 2)
{
    auto left = nodes[0];
    auto right = nodes[1];
    receive_value(left, right);
    return {left}; // Construtor C++ Moderno.
}
if (size == 3)
{
    auto left = nodes[0];
    auto middle = nodes[1];
    auto right = nodes[2];
    receive_value(left, middle); // Left receive value from middle
    receive_value(left, right); // Left receive value from right
    return {left}; // Construtor C++ Moderno.
}

```

Como no exemplo abaixo, onde a segunda iteração do algoritmo o número de nós é 3.



2.7 Question 6

2.7.1 Questão 6

Derive formulas for the number of receives and additions that core 0 carries out using:

- the original pseudo-code for a global sum
- the tree-structured global sum.

Make a table showing the numbers of receives and additions carried out by core 0 when the two sums are used with 2, 4, 8, ..., 1024 cores.

Tabela mostrando o número de receives, ou joins que o core 0 terá

- A coluna Cores representa o número de núcleos, que estão em potência de 2.
- A coluna Naive, é o número de receives que o core 0 terá caso utilizado a abordagem ingenua do core 0 esperar todos os outros e somar todo de uma vez.
- A coluna Tree é o número de receives que o core 0 terá utilizando a abordagem em árvore.

Cores	Naive	Tree
2	1	1
4	3	2
8	7	3
16	15	4
32	31	5
64	63	6
128	127	7
256	255	8
512	512	9
1024	1023	10

Podemos observar claramente que a abordagem ingênua segue a formula, $p - 1$ e quando usamos a árvore, percebemos que a cada 2 núcleos, o número de ligações aumentar em 1, ou seja, $\log(p)$ de base 2. Podemos ver o número de ligações crescendo linearmente com cada

dois núcleos na imagem abaixo

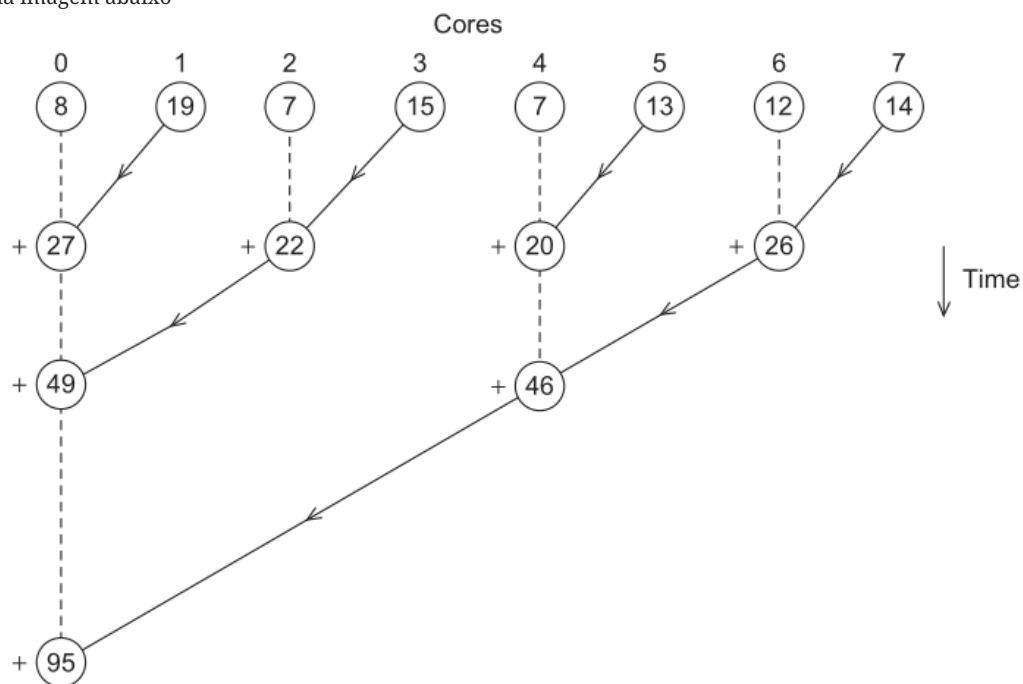


FIGURE 1.1

Multiple cores forming a global sum

2.8 Question 9

2.8.1 Questão 9

Write an essay describing a research problem in your major that would benefit from the use of parallel computing. Provide a rough outline of how parallelism would be used. Would you use task- or data-parallelism ?

TinyML in Context of web browser

Em contexto de aplicações web, especificamente a camada front-end, vem se popularizando frameworks que utilizam o WebAssembly (Wasm). O WebAssembly é um novo tipo de código que pode ser executado em browsers modernos — se trata de uma linguagem de baixo nível como assembly, com um formato binário compacto que executa com performance quase nativa e que fornece um novo alvo de compilação para linguagens como C/C++, para que possam ser executadas na web developer.mozilla.org. Através do Wasm é possível implementar algoritmos de Tiny machine learning (TinyML) para classificação, análise de dados sem a necessidade de comunicação com backend. TinyML é amplamente definido como uma rápida e crescente área de aprendizado de máquina e suas aplicações que inclui hardware, algoritmos e software capazes de performar análise de dados em dispositivos de baixo consumo de energia, tipicamente em milliwatts, assim habilitado variados casos de uso e dispositivos que possuem bateria. Em um navegador, ou durante a navegação o usuário está a todo momento produzindo dados que muitas vezes estão sendo enviados de forma bruta ou quase bruta para a camada de aplicação ou back-end, onde nela é gasto processamento e memória para a primeira etapa de classificação ou análise dos dados. Uma vez analisado desempenho de técnicas e algoritmos de TinyML utilizando WebAssembly, pode ser possível transferir a responsabilidade da análise dos dados para o front-end. Em contexto de navegador quase tudo é paralelo ou distribuído, uma aba ou *tab* em inglês é um processo diferente. Criar um uma extensão, que faça aquisição e análise dos dados de diferentes abas, seria criar um sistema que se comunica com diferentes processos por meio de mensagens e o algoritmo de aprendizado pode fazer uso de idealmente uma ou duas threads para realizar a análise rapidamente. Por tanto é um sistema que é task-and-data parallel.

3. Chapter 2

3.1 Capítulo 2

Capítulo 2: 1-3, 5, 7, 10, 15-17, 19-21, 24; (13 questões)

1. When we were discussing floating point addition, we made the simplifying assumption that each of the functional units took the same amount of time. Suppose that fetch and store each take 2 nanoseconds and the remaining operations each take 1 nanosecond.
 - a. How long does a floating point addition take with these assumptions ?
 - b. How long will an unpipelined addition of 1000 pairs of floats take with these assumptions ?
 - c. How long will a pipelined addition of 1000 pairs of floats take with these assumptions ?
 - d. The time required for fetch and store may vary considerably if the operands/results are stored in different levels of the memory hierarchy. Suppose that a fetch from a level 1 cache takes two nanoseconds, while a fetch from a level 2 cache takes five nanoseconds, and a fetch from main memory takes fifty nanoseconds. What happens to the pipeline when there is a level 1 cache miss on a fetch of one of the operands? What happens when there is a level 2 miss ?

• [Resposta questão 1](#)

2. Explain how a queue, implemented in hardware in the CPU, could be used to improve the performance of a write-through cache.

• [Resposta questão 2](#)

3. Recall the example involving cache reads of a two-dimensional array (page 22). How does a larger matrix and a larger cache affect the performance of the two pairs of nested loops? What happens if $MAX = 8$ and the cache can store four lines ? How many misses occur in the reads of A in the first pair of nested loops ? How many misses occur in the second pair ?

• [Resposta questão 3](#)

4.

5. Does the addition of cache and virtual memory to a von Neumann system change its designation as an SISD system ? What about the addition of pipelining? Multiple issue? Hardware multithreading ?

• [Resposta questão 5](#)

6.

7. Discuss the differences in how a GPU and a vector processor might execute the following code:

```
sum = 0.0;
for (i = 0; i < n; i++) {
    y[i] += a*x[i];
    sum += z[i]*z[i];
}
```

• [Resposta questão 7](#)

8.

9.

10. Suppose a program must execute 10^{12} instructions in order to solve a particular problem. Suppose further that a single processor system can solve the problem in 10^6 seconds (about 11.6 days). So, on average, the single processor system executes 10^6 or a million instructions per second. Now suppose that the program has been parallelized for execution on a distributed-memory system. Suppose also that if the parallel program uses p processors, each processor will execute $10^{12}/p$ instructions and each processor must send $10^9 (p - 1)$ messages. Finally, suppose that there is no additional overhead in executing the parallel program. That is, the program will complete after each processor has executed all of its instructions and sent all of its messages, and there won't be any delays due to things such as waiting for messages.

- a. Suppose it takes 10^{-9} seconds to send a message. How long will it take the program to run with 1000 processors, if each processor is as fast as the single processor on which the serial program was run ?
- b. Suppose it takes 10^{-3} seconds to send a message. How long will it take the program to run with 1000 processors ? [Resposta questão 10](#)

11.

12.

13.

14.

15. Suppose a shared-memory system uses snooping cache coherence and write-back caches. Also suppose that core 0 has the variable x in its cache, and it executes the assignment $x = 5$. Finally suppose that core 1 doesn't have x in its cache, and after core 0's update to x , core 1 tries to execute $y = x$. What value will be assigned to y ? Why?
- b. Suppose that the shared-memory system in the previous part uses a directory-based protocol. What value will be assigned to y ? Why?
- c. Can you suggest how any problems you found in the first two parts might be solved?

• [Resposta questão 15](#)

16. Suppose the run-time of a serial program is given by $T_{serial} = n^2$, where the units of the run-time are in microseconds. Suppose that a parallelization of this program has run-time $T_{parallel} = \frac{n^2}{p} + \log_2(p)$. Write a program that finds the speedups and efficiencies of this program for various values of n and p . Run your program with $n = 10, 20, 40, \dots, 320$, and $p = 1, 2, 4, \dots, 128$. What happens to the speedups and efficiencies as p is increased and n is held fixed? What happens when p is fixed and n is increased?
- b. Suppose that $T_{parallel} = \frac{T_{serial}}{p} + T_{overhead}$. Also suppose that we fix p and increase the problem size.
- Show that if $T_{overhead}$ grows more slowly than T_{serial} , the parallel efficiency will increase as we increase the problem size.
 - Show that if, on the other hand, $T_{overhead}$ grows faster than T_{serial} , the parallel efficiency will decrease as we increase the problem size.

• [Resposta questão 16](#)

17. A parallel program that obtains a speedup greater than p —the number of processes or threads—is sometimes said to have **superlinear speedup**. However, many authors don't count programs that overcome "resource limitations" as having superlinear speedup. For example, a program that must use secondary storage for its data when it's run on a single processor system might be able to fit all its data into main memory when run on a large distributed-memory system. Give another example of how a program might overcome a resource limitation and obtain speedups greater than p .

• [Resposta questão 17](#)

18.

19. Suppose $T_{serial} = n$ and $T_{parallel} = \frac{n}{p} + \log_2(p)$, where times are in microseconds. If we increase p by a factor of k , find a formula for how much we'll need to increase n in order to maintain constant efficiency. How much should we increase n by if we double the number of processes from 8 to 16? Is the parallel program scalable?

• [Resposta questão 19](#)

20. Is a program that obtains linear speedup strongly scalable? Explain your answer.

• [Resposta questão 20](#)

21. Bob has a program that he wants to time with two sets of data, *input_data1* and *input_data2*. To get some idea of what to expect before adding timing functions to the code he's interested in, he runs the program with two sets of data and the Unix shell command *time*:

```
$ time ./bobs prog < input_data1
real 0m0.001s
user 0m0.001s
sys 0m0.000s
$ time ./bobs prog < input_data2
real 1m1.234s
user 1m0.001s
sys 0m0.111s
```

The timer function Bob is using has millisecond resolution. Should Bob use it to time his program with the first set of data? What about the second set of data? Why or why not?

• [Resposta questão 21](#)

22.

23.

24. If you haven't already done so in Chapter 1, try to write pseudo-code for our tree-structured global sum, which sums the elements of *loc_bin_cts*. First consider how this might be done in a shared-memory setting. Then consider how this might be done in a distributed-memory setting. In the shared-memory setting, which variables are shared and which are private?

• [Resposta questão 24](#)

3.2 Question 01

3.2.1 Questão 1

When we were discussing floating point addition, we made the simplifying assumption that each of the functional units took the same amount of time. Suppose that fetch and store each take 2 nanoseconds and the remaining operations each take 1 nanosecond.

- How long does a floating point addition take with these assumptions ?
- How long will an unpipelined addition of 1000 pairs of floats take with these assumptions ?
- How long will a pipelined addition of 1000 pairs of floats take with these assumptions ?
- The time required for fetch and store may vary considerably if the operands/results are stored in different levels of the memory hierarchy. Suppose that a fetch from a level 1 cache takes two nanoseconds, while a fetch from a level 2 cache takes five nanoseconds, and a fetch from main memory takes fifty nanoseconds. What happens to the pipeline when there is a level 1 cache miss on a fetch of one of the operands? What happens when there is a level 2 miss ?

1.a

Instructions	Time in nanosecond
Fetch	2
Store	2
Functional OP	1

"As an alternative, suppose we divide our floating point adder into seven separate pieces of hardware or functional units. The first unit will fetch two operands, the second will compare exponents, and so on." (Página 26)

O Author do livro considera que existe sete operações, considerando que duas delas são fetch e store custa 2 nanosegundos e o restante 1 nanosegundo.

$$1 \cdot 5 + 2 \cdot 2 = 9 \text{ nanosegundos}$$

1.b

Considerando que existem 1000 pares de valores vão serem somados:

$$1000 \cdot 9 = 9000 \text{ nanosegundos}$$

1.c

foi pensado o seguinte: Nó memento que o dado passa pelo fetch, e vai para a próxima operação já é realizado o fetch da segunda operação. Executando o pipeline:

Tempo em nanosegundos	Fetch	OP1	OP2	OP3	OP4	OP5	Store
0	1	wait	wait	wait	wait	wait	wait
2	2	1	wait	wait	wait	wait	wait
3	2	wait	1	wait	wait	wait	wait
4	3	2	wait	1	wait	wait	wait
5	3	wait	2	wait	1	wait	wait
6	4	3	wait	2	wait	1	wait
7	4	wait	3	wait	2	wait	1
8	5	4	wait	3	wait	2	wait
9	5	wait	4	wait	3	wait	2
10	6	5	wait	4	wait	3	2
11	6	wait	5	wait	4	wait	3

Percebe-se que a primeira instrução irá ser finalizada ou sumir na tabela quanto for 9 segundos ou seja, a primeira instrução dura 9 segundos, no entanto, no momento em que a primeira instrução é finalizada, a segunda já começa a ser finalizada ou seja, demora apenas 2 nanosegundos até segunda operação ser finalizada e mais 2 nanosegundos para a terceira ser finalizada e assim por diante. Por tanto para executar todos os 1000 dados, o custo total fica:

$$9 + 999 \cdot 2 = 2007$$

1.d

No caso, considerando que a cache nível não falhe a tabela continua mesma, pois o fetch e store custam o mesmo 2 nanosegundos:

Tempo em nanosegundos	Fetch	OP1	OP2	OP3	OP4	OP5	Store
10	6	5	wait	4	wait	3	2
11	6	wait	5	wait	4	wait	3

mas se imaginarmos que na 12 iteração o Fetch e Store passa a custar 5 nanosegundos:

Tempo em nanosegundos	Fetch	OP1	OP2	OP3	OP4	OP5	Store
10	6	5	wait	4	wait	3	2
11	6	wait	5	wait	4	wait	3
12	6	wait	wait	5	wait	4	3
13	6	wait	wait	wait	5	4	3
14	6	wait	wait	wait	5	4	3
15	7	6	wait	wait	5	4	3
16	7	wait	6	wait	wait	5	4

Quando mais lento fica a transferência para a memória principal, mais nítido fica o gargalo de Von Neumann, ou seja, percebe-se que a performance do processador fica limitado a taxa de transferência de dados com a memória principal.

3.3 Question 02

3.3.1 Questão 2

Explain how a queue, implemented in hardware in the CPU, could be used to improve the performance of a write-through cache.

Como observado na [questão 1](#) cada momento que a escrita é debilitada, fica nítido o gargalo de Von Neuman se considerarmos que uma escrita na cache é uma escrita na memória principal, então cada Store iria demorar **50** nano segundos. Colocando uma fila e supondo que ela nunca fique cheia, a CPU não irá gastar tanto tempo no Store, mas uma vez a fila cheia, a CPU terá que aguardar uma escrita na memória principal.

Tabela com cache miss questão 1

Tempo em nanosegundos	Fetch	OP1	OP2	OP3	OP4	OP5	Store
10	6	5	wait	4	wait	3	2
11	6	wait	5	wait	4	wait	3
12	6	wait	wait	5	wait	4	3
13	6	wait	wait	wait	5	4	3
14	6	wait	wait	wait	5	4	3
15	7	6	wait	wait	5	4	3
16	7	wait	6	wait	wait	5	4

3.4 Question 03

3.4.1 Questão 3

Recall the example involving cache reads of a two-dimensional array (page 22). How does a larger matrix and a larger cache affect the performance of the two pairs of nested loops? What happens if $MAX = 8$ and the cache can store four lines ? How many misses occur in the reads of A in the first pair of nested loops ? How many misses occur in the second pair ?

Exemplo envolvendo leituras na cache em um array bidimensional

```
double A[MAX][MAX], x[MAX], y[MAX];

//. . .
// Initialize A and x, assign y = 0 */
//. . .

/* First pair of loops */
for (i = 0; i < MAX; i++)
    for (j = 0; j < MAX; j++)
        y[i] += A[i][j]*x[j];

//. . .
// Assign y = 0 */
//. . .

/* Second pair of loops */
for (j = 0; j < MAX; j++)
    for (i = 0; i < MAX; i++)
        y[i] += A[i][j]*x[j];
```

Cache line	Elements of A
0	A[0][0] A[0][1] A[0][2] A[0][3]
1	A[1][0] A[1][1] A[1][2] A[1][3]
2	A[2][0] A[2][1] A[2][2] A[2][3]
3	A[3][0] A[3][1] A[3][2] A[3][3]

How does a larger matrix and a larger cache affect the performance of the two pairs of nested loops ?

Supondo que a cache tenha a mesma proporção do que a Matrix, o número de cache miss seria igual ao número de linhas da matriz, como apontado no exemplo dado no livro, quando o processador pede o valor $A[0][0]$, baseado na ideia de vizinhança, a cache carrega todas as outras colunas da linha 0, portanto é plausível pensar que o número de miss é igual ao número de linhas, ou seja, o número miss é igual a **MAX**, pois a cache tem o mesmo número de linhas que a matrix A, suponto que não preciso me preocupar com x e y.

if $MAX = 8$ and the cache can store four lines ? How many misses occur in the reads of A in the first pair of nested loops ? How many misses occur in the second pair ?

FIRST PAIR OF LOOPS

Tendo a a cache armazenando metade dos valores de uma linha da Matriz A, então para cada linha da Matriz, vai haver duas cache miss, a primeira no $A[i][0]$ e a segunda no $A[i][4]$. Outro ponto é que como a cache só possui 4 linhas, então após ocorrer os cache misses $A[0][0]$, $A[0][4]$ e $A[1][0]$, $A[1][4]$ toda a cache terá sido preenchida, ou seja, Tendo a matriz 8 linhas e para cada linha tem 2 cache miss por tanto:

```
8*2 =16 cache miss
```

como tanto a primeira parte quando na segunda parte, percorre-se todas as linhas irá haver 16 cache miss, suponto que não preciso me preocupar com x e y .

Cache line	Elements of A
0	A[0][0] A[0][1] A[0][2] A[0][3]
1	A[0][4] A[0][5] A[0][6] A[0][7]
2	A[1][0] A[1][1] A[1][2] A[1][3]
3	A[1][4] A[1][5] A[1][6] A[1][7]

SECOND PAIR OF LOOPS

No segundo par de loops, vemos que o segundo laço for, percorre os valores: A[0][0], A[1][0], A[2][0] ... A[7][0], para quando $j=0$. Isso faz com que todo hit seja miss, ou seja iremos ter miss para cada acesso em A, portanto:

8*8 = 64 cache miss

```
/* Second pair of loops */
for (j = 0; j < MAX; j++)
    for (i = 0; i < MAX; i++)
        y[i] += A[i][j]*x[j]; // cada acesso um miss.
```

3.5 Question 05

3.5.1 Questão 5

Does the addition of cache and virtual memory to a von Neumann system change its designation as an SISD system ? What about the addition of pipelining? Multiple issue? Hardware multithreading ?

Um SISD system ou Single Instruction Single Data system, são sistemas que executam uma única instrução por vez e sua taxa de transferência de dados é de um item por vez também.

the addition of cache and virtual memory

Adicionar um cache e memória virtual, pode ajudar a reduzir o tempo que **única** instrução é lida da memória principal, mas não aumenta o número de instruções buscadas na operação Fetch ou na Operação Store, por tanto o sistema continua sendo Single Instruction Single Data.

the addition of pipelining

Como demonstrado na [questão 1](#), ao adicionar um pipeline, podemos realizar a mesma instrução complexa em múltiplos dados, ou seja, Single Instruction Multiple Data System, portanto sim.

the addition of Multiple issue and Hardware multithreading

No momento em que possibilitamos uma máquina executar antecipadamente uma instrução ou possibilitamos a execução de múltiplas threads, nesse momento então a máquina está executando várias instruções em vários dados ao mesmo tempo, por tanto o sistema se torna Multiple Instruction Multiple Data system, ou seja, a designação muda.

3.6 Question 07

3.6.1 Questão 7

Discuss the differences in how a GPU and a vector processor might execute the following code:

```
sum = 0.0;
for (i = 0; i < n; i++) {
    y[i] += a*x[i];
    sum += z[i]*z[i];
}
```

CPU com vetorização

Um processo de vetorização em cima desse laço for dividiria as entradas em chunks ou blocos de dados e executaria em paralelo a instrução complexa. Algo como:

```
//executando o bloco paralelamente.
y[0] += a*x[0];
y[1] += a*x[1];
y[2] += a*x[2];
y[3] += a*x[3];

z[0]*z[0] // executando em paralelo
z[1]*z[1] // executando em paralelo
z[2]*z[2] // executando em paralelo
z[3]*z[3] // executando em paralelo
// somando tudo depois
sum+= z[0]*z[0] + z[1]*z[1] + z[2]*z[2] + z[3]*z[3];
```

GPU

Atualmente essa operação em GPU é muito mais interessante, pois hoje podemos compilar ou gerar instruções complexas que podem ser executadas em GPU. A primeira vantagem seria separar o cálculo do sum:

```
sum += z[i]*z[i];
```

do cálculo do y

```
y[i] += a*x[i];
```

ficando assim:

SHADERS GLSL EQUIVALENTE SUPONTO QUE N = 4

```
# version 330

layout(location = 0) in vec4 x;
layout(location = 1) in mat4 a;
layout(location = 2) in vec4 y;

/* o buffer gl_Position é re-inserido no y */
void main()
{
    gl_Position = a*x + y;
}
//
```

```
# version 330

layout(location = 0) in mat4 z;
uniform float sum;

/* transpose é uma função que calcula a transposta já existem no Glsl */
void main()
{
    mat4 temp = transpose(z) * z;

    sum = 0;
    for (int i = 0; i < 4; i++)
        // desde que o laço for seja baseado em constantes ou variáveis uniformes
        // esse laço for é possível.
```

```
{  
    sum += temp[i];  
}  
  
// recupera o valor no index 0  
gl_Position = vec4(sum, 0.0, 0.0, 0.0, 0.0);  
}
```

A grande vantagem de usar os shaders seria dependendo do tamanho do vetor de dados, executar as instruções de uma só vez em todos os dados, na prática assim como a vetorização envia em blocos, na GPU você também enviaria em blocos, comumente chamados de buffers, a grande diferença seria justamente no fato que um bloco na GPU possui um tamanho muito maior que o bloco da vetorização.

3.7 Question 10

3.7.1 Questão 10

Suppose a program must execute 10^{12} instructions in order to solve a particular problem. Suppose further that a single processor system can solve the problem in 10^6 seconds (about 11.6 days). So, on average, the single processor system executes 10^6 or a million instructions per second. Now suppose that the program has been parallelized for execution on a distributed-memory system. Suppose also that if the parallel program uses p processors, each processor will execute $\frac{10^{12}}{p}$ instructions and each processor must send $10^9(p - 1)$ messages. Finally, suppose that there is no additional overhead in executing the parallel program. That is, the program will complete after each processor has executed all of its instructions and sent all of its messages, and there won't be any delays due to things such as waiting for messages.

1. Suppose it takes 10^{-9} seconds to send a message. How long will it take the program to run with 1000 processors, if each processor is as fast as the single processor on which the serial program was run ?
2. Suppose it takes 10^3 seconds to send a message. How long will it take the program to run with 1000 processors ?

main.py

```
import datetime

NUMBER_OF_INSTRUCTIONS = 10**12
NUMBER_OF_MESSAGES = 10**9
AVERAGE_SECOND_PER_INSTRUCTIONS = (10**6) / NUMBER_OF_INSTRUCTIONS

def cost_time_per_instruction(instructions: int) -> float:
    return AVERAGE_SECOND_PER_INSTRUCTIONS * instructions

def number_of_instructions_per_processor(p: int) -> int:
    return NUMBER_OF_INSTRUCTIONS/p

def number_of_messages_per_processor(p: int) -> int:
    return NUMBER_OF_MESSAGES * (p-1)

def simulate(time_per_message_in_seconds: float, processors: int):
    print(
        f'time to send a message: {time_per_message_in_seconds} processors: {processors}')
    instructions = number_of_instructions_per_processor(processors)
    number_of_messages = number_of_messages_per_processor(processors)
    each_process_cost_in_seconds = cost_time_per_instruction(instructions)

    total_messages_in_seconds = time_per_message_in_seconds * number_of_messages

    result = total_messages_in_seconds + each_process_cost_in_seconds
    result_date = datetime.timedelta(seconds=result)

    print(f'executing instructions is {instructions}')
    print(f'spend sending messages is {total_messages_in_seconds}')

    print(f'total time in seconds: {result}')
    print(f'total time in HH:MM:SS {result_date}')

def a():
    time_per_message_in_seconds = 1e-9
    processors = 1e3
    simulate(time_per_message_in_seconds, processors)

def b():
    time_per_message_in_seconds = 1e-3
    processors = 1e3
    simulate(time_per_message_in_seconds, processors)

def main():
    print('A:')
    a()
    print('B:')
    b()
```

```
if __name__ == '__main__':  
    main()
```

Executando:

```
python chapter_2/question_10/main.py  
A:  
time to send a message: 1e-09 processors: 1000.0  
executing instructions is 1000000000.0  
spend sending messages is 999.0000000000000001  
total time in seconds: 1999.0  
total time in HH:MM:SS 0:33:19  
B:  
time to send a message: 0.001 processors: 1000.0  
executing instructions is 1000000000.0  
spend sending messages is 999000000.0  
total time in seconds: 999001000.0  
total time in HH:MM:SS 11562 days, 12:16:40
```

11562 dias são 32 anos.

3.8 Question 15

3.8.1 Questão 15

1. Suppose a shared-memory system uses snooping cache coherence and write-back caches. Also suppose that core 0 has the variable x in its cache, and it executes the assignment $x = 5$. Finally suppose that core 1 doesn't have x in its cache, and after core 0's update to x , core 1 tries to execute $y = x$. What value will be assigned to y ? Why?
2. Suppose that the shared-memory system in the previous part uses a directory-based protocol. What value will be assigned to y ? Why?
3. Can you suggest how any problems you found in the first two parts might be solved?

1 What value will be assigned to y ? Why?

Não é possível determinar qual valor será atribuído ao y independentemente se for write-back ou write-through, uma vez que não houve uma sincronização entre os cores sobre o valor de x . A atribuição de $y = x$ do core 1 pode ocorrer antes ou depois das operações no core 0.

shared-memory system in the previous part uses a directory-based protocol What value will be assigned to y ? Why?

Com o sistema de arquivos, ao core 0 irá notificar o a memória principal que a consistência dos dados foi comprometida, no entanto, ainda não dá para saber qual o valor de y , uma vez que a atribuição de $y = x$ do core 1 pode ocorrer antes ou depois das operações no core 0.

Can you suggest how any problems you found in the first two parts might be solved?

Existe dois problemas, o problema da consistência do dados, temos que garantir que ambos os cores façam alterações que ambas sejam capaz de ler e o segundo é um mecanismo de sincronização, onde por exemplo, o core 1 espera o core 0 finalizar o seu processamento com a variável x para ai sim começar o seu. Podemos utilizar por exemplo um mutex, onde inicialmente o core 0 faria o lock e ao finalizar ele entrega a chave a qual, o core 1 pegaria.

3.9 Question 16

3.9.1 Questão 16

a. Suppose the run-time of a serial program is given by $T_{serial} = n^2$, where the units of the run-time are in microseconds. Suppose that a parallelization of this program has run-time $T_{parallel} = \frac{n^2}{p} + \log_2(p)$. Write a program that finds the speedups and efficiencies of this program for various values of n and p . Run your program with $n = 10, 20, 40, \dots, 320$, and $p = 1, 2, 4, \dots, 128$. What happens to the speedups and efficiencies as p is increased and n is held fixed? What happens when p is fixed and n is increased?

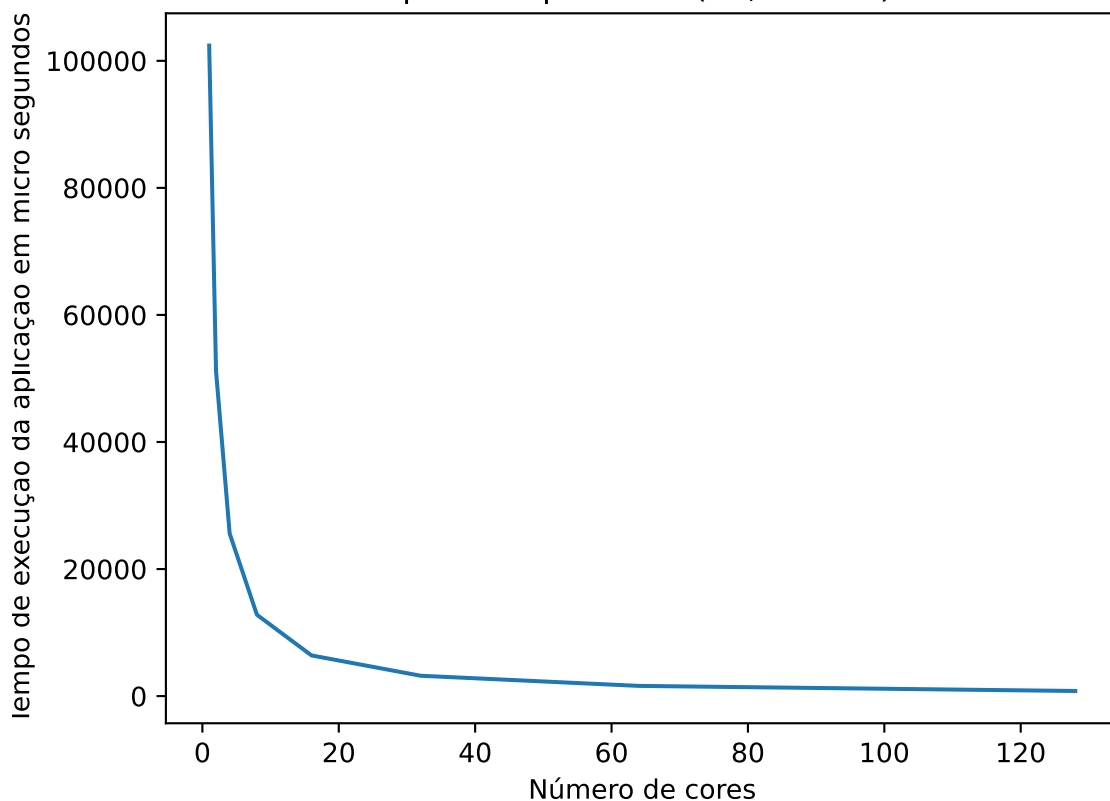
b. Suppose that $T_{parallel} = \frac{T_{serial}}{p} + T_{overhead}$. Also suppose that we fix p and increase the problem size.

- Show that if $T_{overhead}$ grows more slowly than T_{serial} , the parallel efficiency will increase as we increase the problem size.
- Show that if, on the other hand, $T_{overhead}$ grows faster than T_{serial} , the parallel efficiency will decrease as we increase the problem size.

16 a

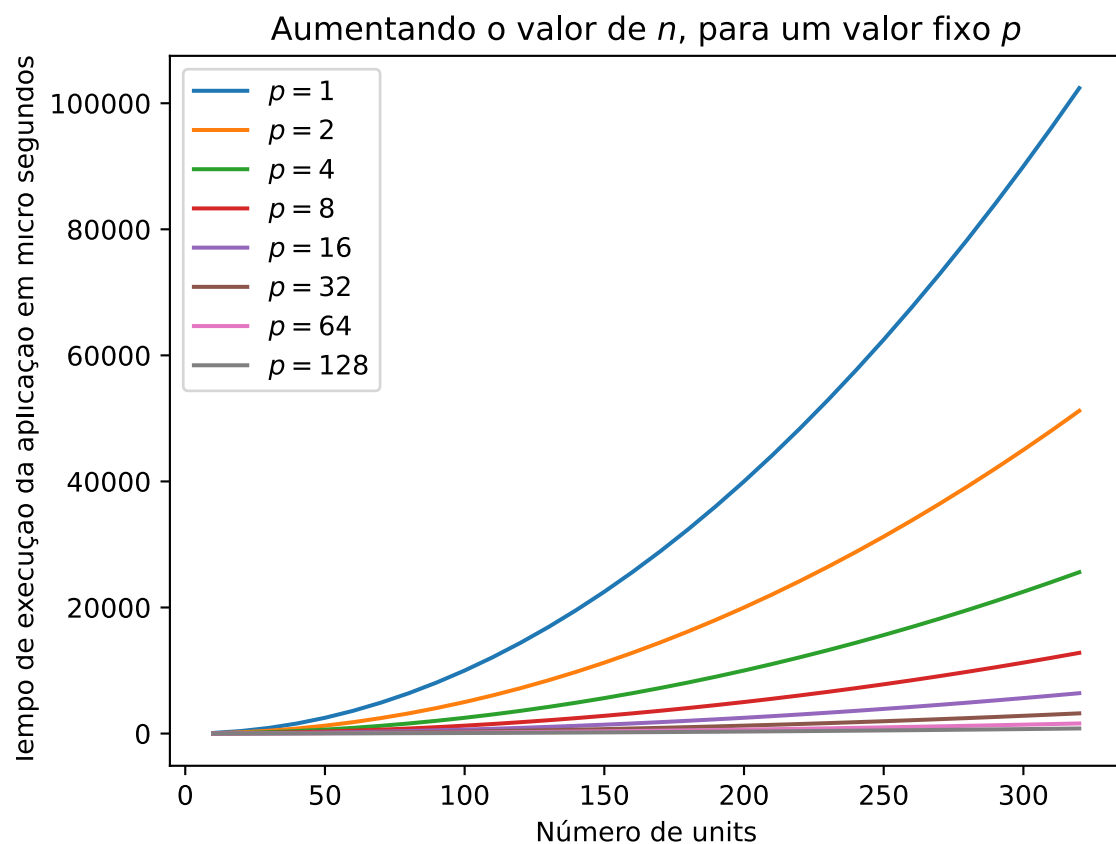
WHAT HAPPENS TO THE SPEEDUPS AND EFFICIENCIES AS P IS INCREASED AND N IS HELD FIXED ?

Aumentando o valor de p , para um $n = 320$,
possível platô em (64, 1606.0)



Podemos observar a Amadahl's law, "a lei de Amadahl diz: a menos que virtualmente todo o programa serial seja paralelizado, o possível speedup, ou ganhado de performance será muito limitado, independentemente dos números de cores disponíveis." No exemplo em questão, podemos observar que a partir de 32 cores, o tempo da aplicação fica estagnado por volta dos 1606 micro segundos.

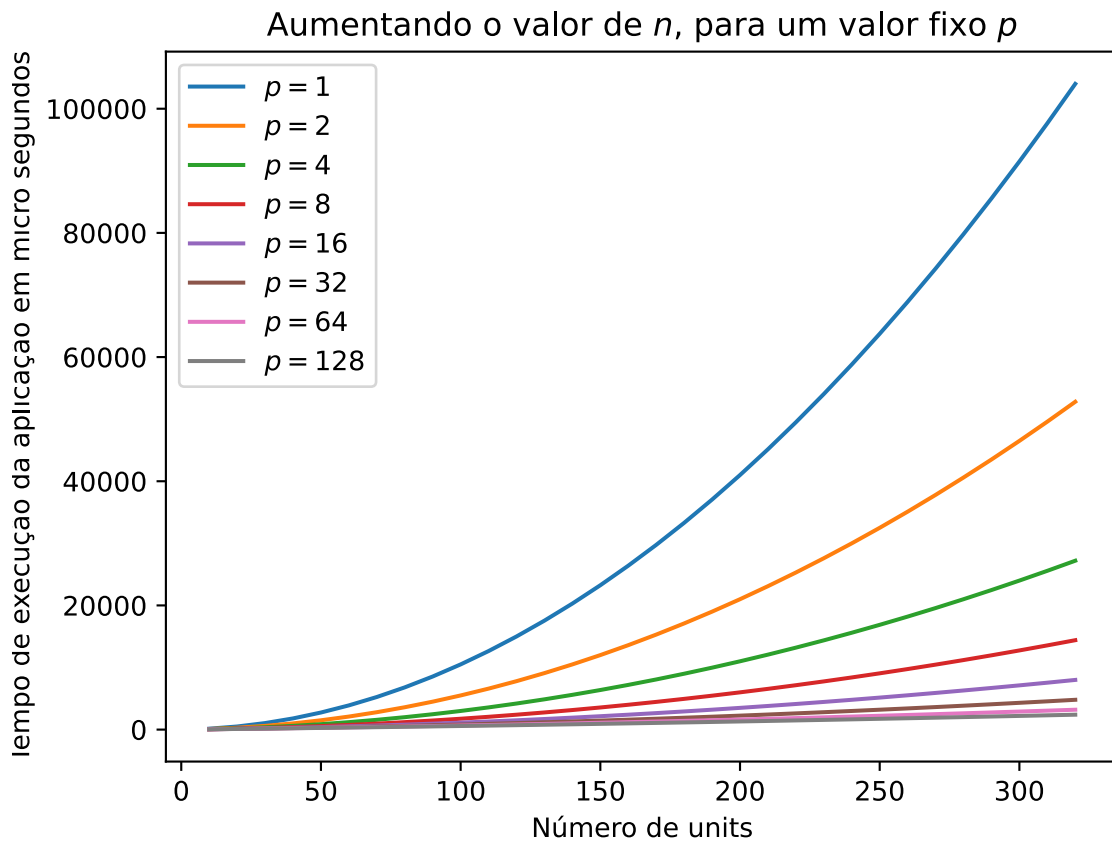
WHAT HAPPENS WHEN P IS FIXED AND N IS



Mesmo com a lei de Amdahl, podemos observar que o aumento de performance é bastante significativo, valendo a pena paralelizar a aplicação, também pelo fato que os hardwares atualmente possuem mais de um core.

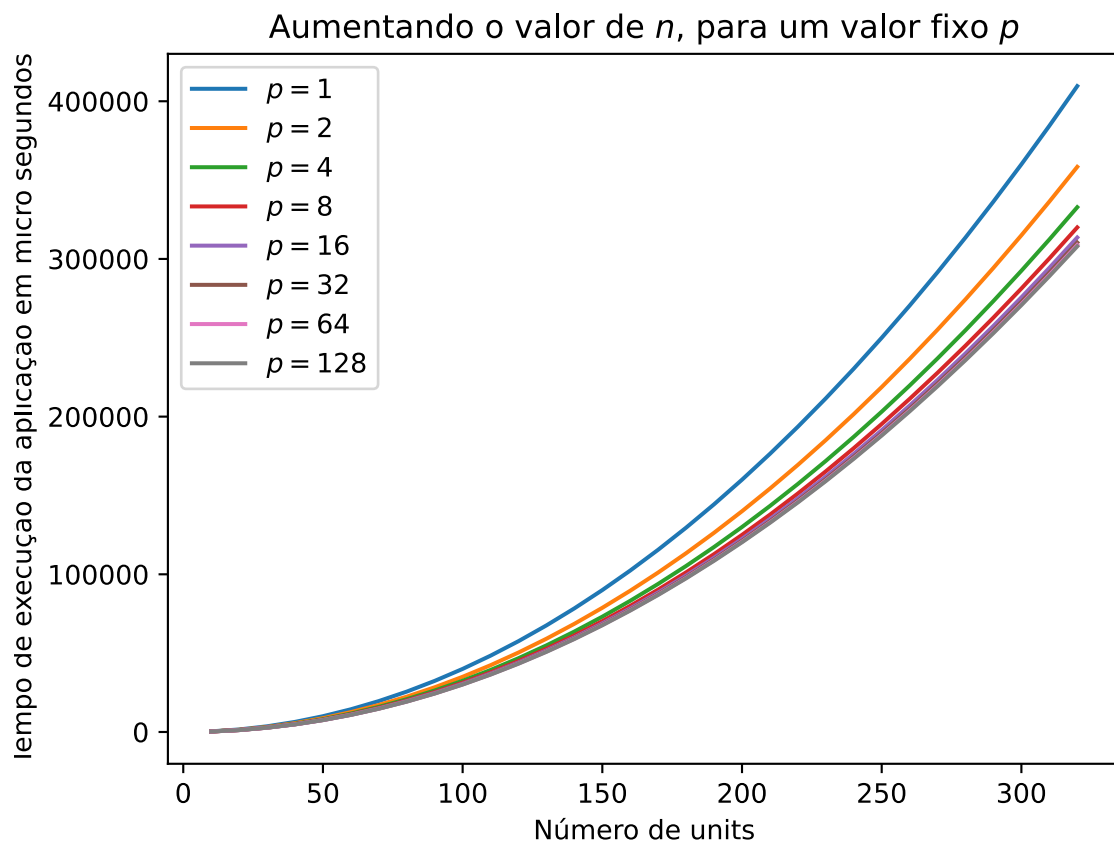
16 b

SHOW THAT IF $T_{overhead}$ GROWS MORE SLOWLY THAN T_{serial} , THE PARALLEL EFFICIENCY WILL INCREASE AS WE INCREASE THE PROBLEM SIZE.



Sabendo que o tempo serial é o quadrático da entrada, considerei o tempo de overhead sendo: $T_{overhead} = 5n$, ou seja, uma função linear. Comparando com o gráfico da letra **a** dificilmente é notado uma diferença entre os gráficos, podendo até ser desconsiderado.

SHOW THAT IF, ON THE OTHER HAND, $T_{overhead}$ GROWS FASTER THAN T_{serial} , THE PARALLEL EFFICIENCY WILL DECREASE AS WE INCREASE THE PROBLEM SIZE.



Sabendo que o tempo serial é o quadrático da entrada, considerei o tempo de overhead sendo: $T_{overhead} = 3T_{serial} = 3n^2$, ou seja, o overhead cresce 3 vezes mais que o serial. Comparando com o gráfico da letra **a** podemos observar que paralelizar não é uma boa opção, pois nem com $p = 128$ consegue ser melhor do que com $p = 1$, ou seja, não paralelizar, ou seja, a solução serial sem o overhead.

Observações

Para essa atividade foi utilizado Python na sua versão **3.10.4** para instalar as dependências do script e criar os gráficos:

```
# para criar um ambiente virtual, suponto que tenha python 3.10.4 instalado
python -menv .env
source .env/bin/activate
# instalando as dependências e executando a aplicação
pip install -r requirements.txt
python main.py
```

3.10 Question 17

3.10.1 Questão 17

A parallel program that obtains a speedup greater than p —the number of processes or threads—is sometimes said to have **superlinear speedup**. However, many authors don't count programs that overcome "resource limitations" as having superlinear speedup. For example, a program that must use secondary storage for its data when it's run on a single processor system might be able to fit all its data into main memory when run on a large distributed-memory system. Give another example of how a program might overcome a resource limitation and obtain speedups greater than p

HTTP server em ESP32

Sistemas embarcados com CPUs multi-core. Atualmente já existem microcontroladores como **ESP32** que dependendo do modelo pode possuir mais de um núcleo de processamento. Sabendo que todo o programa fica carregado na memória, então uma aplicação como um servidor HTTP, pode ter mais que o dobro de performance, quando observado o seu requests por segundo.

Fazemos o seguinte cenário:

Temos um desenvolvedor que sabe que executar operações de escrita em hardware é uma operação insegura e utiliza estruturas de dados para sincronização dessas operações, cada dispositivo tem seu tempo de sincronização. Temos que o dispositivo A custa 5 milissegundos, e dispositivo B custa 4 milissegundos. Sabendo que se criarmos a aplicação em single core, temos que esperar a sincronização de A e B e que modificar o tempo de sincronização de um dispositivo custa 3 milissegundos. Dado que se tem 2 cores, 2 conversores AD, se delegarmos cada dispositivo para um core, eliminaremos 3 milissegundos do processo de escrita no seu pior cenário. Supondo que o tempo de um request fique 30% na escrita de um dispositivo e que ele gasta em outras operações **8 milissegundos** temos dois cenários,

CENÁRIO 1

$$R_e = 0.7(8) + 0.3(5 + 3)$$

$R_e = 8$, ou seja uma Request de escrita custa 8 milissegundos

CENÁRIO 2

$$R_e = 0.7(8) + 0.3(5)$$

$R_e = 7.1$, ou seja uma Request de escrita custa 7.1 milissegundos, além do fato do **ESP32** ser capaz de realizar o dobro de requests

3.11 Question 19

3.11.1 Questão 19

Suppose $T_{serial} = n$ and $T_{parallel} = \frac{n}{p} + \log_2(p)$, where times are in microseconds. If we increase p by a factor of k , find a formula for how much we'll need to increase n in order to maintain constant efficiency. How much should we increase n by if we double the number of processes from 8 to 16? Is the parallel program scalable?

Encontre uma formula para o quanto nós teremos que aumentar n para obter uma eficiência constante.

$$E(p) = \frac{T_{serial}}{pT_{parallel}}$$

$$E(p) = \frac{n}{p(\frac{n}{p} + \log_2(p))}$$

$$E(p) = \frac{n}{n + p\log_2(p)}$$

$$E(kp) = \frac{n}{n + kp\log_2(kp)}$$

Se considerarmos a constante A o número de vezes que temos que aumentar n para obter uma eficiência constante, logo:

$$E_a(kp) = \frac{An}{An + kp\log_2(kp)}$$

$$E_a(kp) = E(p)$$

$$\frac{An}{An + kp\log_2(kp)} = \frac{n}{n + p\log_2(p)}$$

$$An = \frac{n(nA + kp\log_2(kp))}{n + p\log_2(p)}$$

$$A = \frac{nA + kp\log_2(kp)}{n + p\log_2(p)}$$

$$A = \frac{nA}{n + p\log_2(p)} + \frac{kp\log_2(kp)}{n + p\log_2(p)}$$

$$A - \frac{nA}{n + p\log_2(p)} = \frac{kp\log_2(kp)}{n + p\log_2(p)}$$

$$A[1 - \frac{n}{n + p\log_2(p)}] = \frac{kp\log_2(kp)}{n + p\log_2(p)}$$

$$A[\frac{n + p\log_2(p)}{n + p\log_2(p)} - \frac{n}{n + p\log_2(p)}] = \frac{kp\log_2(kp)}{n + p\log_2(p)}$$

$$A \frac{p\log_2(p)}{n + p\log_2(p)} = \frac{kp\log_2(kp)}{n + p\log_2(p)}$$

$$Ap\log_2(p) = kp\log_2(kp)$$

$$A = \frac{kp\log_2(kp)}{p\log_2(p)}$$

$$A = \frac{k\log_2(kp)}{\log_2(p)}$$

Quantas vezes nós devemos aumentar n se nós dobrarmos o número de cores de 8 para 16

$$A(k, p) = \frac{k\log_2(kp)}{\log_2(p)} \text{ se } k = 2 \text{ e } p = 8 \text{ então:}$$

$$A = \frac{2\log_2(16)}{\log_2(8)}$$

$$A = \frac{2(4)}{3}$$

$$A = \frac{8}{3}$$

O programa paralelo é escalável ?

Dado a definição do autor sim. Para o autor escalável é quando a eficiência de um programa paralelo se mantém constante, ou seja, se existe uma taxa que relaciona o crescimento do tamanho do problema com o crescimento do número de threads/processos, então o programa paralelo é fracamente escalável (***weakly scalable***), .

3.12 Question 20

3.12.1 Questão 20

Is a program that obtains linear speedup **strongly** scalable ? Explain your answer.

Dado a definição do autor sim. Para o autor escalável é quando a eficiência de um programa paralelo se mantém constante. Linear speedup pode ser escrito pela seguinte expressão:

$$S = \frac{T_{serial}}{T_{parallel}} = p, \text{ onde } p \text{ é número de cores e } S \text{ é o speedup.}$$

Portanto dado que eficiência é dado pela seguinte expressão:

$$E = \frac{T_{serial}}{pT_{parallel}}, \text{ onde } T_{serial} \text{ é o tempo da aplicação em serial e } T_{parallel} \text{ o tempo da aplicação em paralelizada.}$$

se o speedup for linear, ou seja, $S = p$, temos que

$$E = \frac{S}{p}, \text{ portanto}$$

$$E = \frac{p}{p} = 1, \text{ Como a eficiência é constante, logo, por definição a aplicação é fortemente escalável (**strongly scalable**).$$

3.13 Question 21

3.13.1 Questão 21

Bob has a program that he wants to time with two sets of data, *input_data1* and *input_data2*. To get some idea of what to expect before adding timing functions to the code he's interested in, he runs the program with two sets of data and the Unix shell command *time*:

```
$ time ./bobs prog < input data1
real 0m0.001s
user 0m0.001s
sys 0m0.000s
$ time ./bobs prog < input data2
real 1m1.234s
user 1m0.001s
sys 0m0.111s
```

The timer function Bob is using has millisecond resolution. Should Bob use it to time his program with the first set of data ? What about the second set of data ? Why or why not ?

Segundo a referência o comando *time*, retorna três valores:

- real, que o tempo total desde a inicialização até a terminação do programa
- user, o total de cpu usado pelo programa em modo usuário, ou seja, o a quantidade de cpu usada, eliminando chamadas do sistema e o tempo que o sistema ficou bloqueado ou aguardando outros processos.
- Sys, o tempo de cpu usado pelo kernel/ SO para esse processo em específico.

Primeiro timer

Na primeira chamada observamos que o tempo coletado é praticamente 0, ou seja, o tempo levado para executar o programa está fora da resolução do relógio do sistema, por tanto não podemos concluir nada sobre a primeira chamada e se essa for a primeira chamada é bem provável que a próxima também dê praticamente 0, uma vez que a aplicação pode ter pouco tamanho de entrada se comparado a máquina em questão.

Segundo timer

Já no segundo timer podemos observar informações sobre a aplicação, estão dentro da resolução do relógio e que maior parte da aplicação foi gasta em user mode. Bob pode fazer proveito dessas informações.

3.14 Question 24

3.14.1 Questão 24

If you haven't already done so in Chapter 1, try to write pseudo-code for our tree-structured global sum, which sums the elements of `loc_bin_cts`. First consider how this might be done in a shared-memory setting. Then consider how this might be done in a distributed-memory setting. In the shared-memory setting, which variables are shared and which are private ?

Um código em C++, mas sem o uso de estruturas de dados de programação paralela, pode ser observado na resposta da questão 5 [Question_5.cpp](#). Também foi implementado o **mesmo** algoritmo em rust e nessa implementação foi utilizado threads, smart pointers e mutex para resolver o problema. O código pode ser observado aqui: [main.rs](#)

Main.rs

```
use std::{
    sync::{Arc, Mutex},
    thread::JoinHandle,
};

#[derive(Debug)]
struct Node {
    data: Vec<i32>,
    neighborhoods: Mutex<Vec<JoinHandle<i32>>>,
}

impl Node {
    fn from_value(data: i32) -> Node {
        Node {
            data: vec![data],
            neighborhoods: Mutex::new(vec![]),
        }
    }

    fn compute(self) -> i32 {
        /*
         * Em termos de memória essa função desaloca toda memória
         * usada pela estrutura Node e retorna um inteiro de 32bits.
         * Dado que foi utilizado safe rust e o código compila, logo esse
         * código está livre de race e como não referências cíclicas
         * também está livre de memory leak.
         */
        let result: i32 = self.data.iter().sum();

        let neighborhoods = self.neighborhoods.into_inner().unwrap();

        let neighborhoods_sum: i32 = neighborhoods
            .into_iter()
            .map(|handle| handle.join().expect("Unable to lock neighborhood"))
            .sum();

        result + neighborhoods_sum
    }
}

fn start_to_compute_node(node: Node) -> JoinHandle<i32> {
    std::thread::spawn(move || {
        let result = node.compute();
        std::thread::sleep(std::time::Duration::from_micros(500));

        result
    })
}

fn receive_value(left: Arc<Node>, right: Arc<Node>) {
    let right = Arc::try_unwrap(right).unwrap();

    let mut left_neighborhoods = left
        .neighborhoods
        .lock()
        .expect("Unable to lock neighborhood");

    left_neighborhoods.push(start_to_compute_node(right))
}

fn create_new_tree_bitwise(mut nodes: Vec<Arc<Node>>) -> Vec<Arc<Node>> {
    let size = nodes.len();

    match size {
        2 => {
            let left = nodes.remove(0);
            let right = nodes.remove(0);
```

```

        receive_value(left.clone(), right);

        vec![left]
    }

    3 => {
        let left = nodes.remove(0);
        let middle = nodes.remove(0);
        let right = nodes.remove(0);

        receive_value(left.clone(), middle);
        receive_value(left.clone(), right);

        vec![left]
    }

    - => {
        let mut new_nodes = vec![];

        let mut size = size;

        if size % 2 != 0 {
            size = size - 1;

            new_nodes.push(nodes.remove(size - 1));
        }

        for i in (0..size).step_by(2) {
            let left = nodes.remove(0);
            let right = nodes.remove(0);
            println!("i: {} left: {:?} right: {:?}", i, left, right);
            receive_value(left.clone(), right);
            new_nodes.push(left);
        }
        println!("Next iteration");

        create_new_tree_bitwise(new_nodes)
    }
}

fn main() {
    let data = vec![8, 19, 7, 15, 7, 13, 12, 14];

    let nodes: Vec<Arc<Node>> = data
        .clone()
        .iter()
        .map(|&v| Node::from_value(v))
        .map(|node| Arc::new(node))
        .collect();

    let root = create_new_tree_bitwise(nodes)[0].clone();

    let root = Arc::try_unwrap(root).unwrap();

    let total = root.compute();

    assert!(total == data.iter().sum::<i32>());
    println!(
        "total: {} data sum {} ",
        total,
        data.iter().sum::<i32>()
    );
}

#[test]
fn tree_sum() {
    /*
     * Teste com várias entradas
     */
    let data_tests = vec![
        vec![8, 19, 7, 15, 7, 13, 12, 14],
        vec![8, 19, 7, 15, 7, 13, 12],
        vec![8, 19, 7, 15, 7, 13],
        vec![8, 19, 7, 15, 7],
        vec![8, 19, 7, 15],
        vec![8, 19, 7],
        vec![8, 19],
    ];

    for data in data_tests {
        let nodes: Vec<Arc<Node>> = data
            .clone()
            .iter()
            .map(|&v| Node::from_value(v))
            .map(|node| Arc::new(node))
            .collect();

        let root = create_new_tree_bitwise(nodes)[0].clone();

        let root = Arc::try_unwrap(root).unwrap();

        assert_eq!(root.compute(), data.iter().sum::<i32>());
    }
}

```

```
}  
}
```

Observação

caso tenha o rust instalado você pode observar a execução do caso de teste com o comando **cargo**.

```
cargo test
```

O teste é mesmo teste feito em c++.

4. Chapter 3

4.1 Capítulo 3

Capítulo 3: 2, 4, 6, 9, 11, 12, 13, 16, 17, 19, 20, 22, 23, 27 e 28 (16 questões);

- 1.
2. Modify the trapezoidal rule so that it will correctly estimate the integral even if *comm_sz* doesn't evenly divide *n*. (You can still assume that $n \geq \text{comm_sz}$).
- [Resposta questão 2](#)
- 3.
4. Modify the program that just prints a line of output from each process (*mpi_output.c*) so that the output is printed in process rank order: process 0s output first, then process 1s, and so on.
- 5.
6. Suppose *comm_sz* = 4 and suppose that **x** is a vector with $n = 14$ components.
 - a. How would the components of **x** be distributed among the processes in program that used a block distribution ?
 - b. How would the components of **x** be distributed among the processes in a program that used a cyclic distribution ?
 - c. How would the components of **x** be distributed among the processes in a program that used a block-cyclic distribution with blocksize $b = 2$?
- 7.
- 8.
9. Write an MPI program that implements multiplication of a vector by a scalar and dot product. The user should enter two vectors and a scalar, all of which are read in by process 0 and distributed among the processes. The results are calculated and collected onto process 0, which prints them. You can assume that n , the order of the vectors, is evenly divisible by *comm_sz*

- 10.
11. Finding **prefix sums** is a generalization of global sum. Rather than simply finding the sum of n values,

$x_0 + x_1 + \dots + x_{n-1}$
 the prefix sums are the n partial sums
 $x_0, x_0 + x_1, x_0 + x_1 + x_2, \dots, x_0 + x_1 + \dots + x_{n-1}$

- a. Devise a serial algorithm for computing the n prefix sums of an array with n elements.
- b. Parallelize your serial algorithm for a system with n processes, each of which is storing one of the x_i s.
- c. Suppose $n = 2^k$ for some positive integer k . Can you devise a serial algorithm and a parallelization of the serial algorithm so that the parallel algorithm requires only k communication phases ?
- d. MPI provides a collective communication function, MPI Scan, that can be used to compute prefix sums:

```
int MPI_Scan(
    void*      sendbuf_p /* in */,
    void*      recvbuf_p /* out */,
    int        count     /* in */,
    MPI_Datatype datatype /* in */,
    MPI_Op     op        /* in */,
    MPI_Comm   comm      /* in */);
```

It operates on arrays with *count* elements; both *sendbuf_p* and *recvbuf_p* should refer to blocks of *count* elements of type *datatype*. The *op* argument is the same as *op* for *MPI_Reduce*. Write an MPI program that generates a random array of *count* elements on each MPI process, finds the prefix sums, and print the results.

12. An alternative to a butterfly-structured allreduce is a **ring-pass** structure. In a ring-pass, if there are p processes, each process q sends data to process $q + 1$, except that process $p - 1$ sends data to process 0. This is repeated until each process has the desired result. Thus, we can implement allreduce with the following code:

```
sum = temp_val = my_val;

for (i = 1; i < p; i++) {
    MPI_Sendrecv_replace(&temp_val, 1, MPI_INT, dest,
        sendtag, source, recvtg, comm, &status);
    sum += temp_val;
}
```

- a. Write an MPI program that implements this algorithm for allreduce. How does its performance compare to the butterfly-structured allreduce ?
- b. Modify the MPI program you wrote in the first part so that it implements prefix sums.

13. MPI Scatter and MPI Gather have the limitation that each process must send or receive the same number of data items. When this is not the case, we must use the MPI functions MPI Gatherv and MPI Scatterv. Look at the man pages for these functions, and modify your vector sum, dot product program so that it can correctly handle the case when n isn't evenly divisible by $comm_sz$.
- 14.
- 15.
16. Suppose $comm_sz = 8$ and the vector $\mathbf{x} = (0, 1, 2, \dots, 15)$ has been distributed among the processes using a block distribution. Draw a diagram illustrating the steps in a butterfly implementation of allgather of \mathbf{x}
17. *MPI_Type_contiguous* can be used to build a derived datatype from a collection of contiguous elements in an array. Its syntax is

```
int MPI_Type_contiguous(
    int          count,          /* in */
    MPI_Datatype old_mpi_t,      /* in */
    MPI_Datatype* new_mpi_t_p    /* out */);
```

Modify the *Read_vector* and *Print_vector* functions so that they use an MPI datatype created by a call to *MPI_Type_contiguous* and a *count* argument of 1 in the calls to *MPI_Scatter* and *MPI_Gather*.

- 18.
19. *MPI_Type_indexed* can be used to build a derived datatype from arbitrary array elements. Its syntax is

```
int MPI_Type_indexed(
    int          count,          /* in */
    int          array_of_blocklengths[], /* in */
    int          array_of_displacements[], /* in */
    MPI_Datatype old_mpi_t,      /* in */
    MPI_Datatype* new_mpi_t_p    /* out */);
```

Unlike *MPI_Type_create_struct*, the displacements are measured in units of *old_mpi_t* --not bytes. Use *MPI_Type_indexed* to create a derived datatype that corresponds to the upper triangular part of a square matrix. For example in the 4×4 matrix.

$$\begin{pmatrix} 0 & 1 & 2 & 3 \\ 4 & 5 & 6 & 7 \\ 8 & 9 & 10 & 11 \\ 12 & 13 & 14 & 15 \end{pmatrix}$$

the upper triangular part is the elements 0, 1, 2, 3, 5, 6, 7, 10, 11, 15. Process 0 should read in an $n \times n$ matrix as a one-dimensional array, create the derived datatype, and send the upper triangular part with a single call to *MPI_Send*. Process 1 should receive the upper triangular part with a single call of *MPI_Recv* and then print the data it received.

20. The functions *MPI_Pack* and *MPI_Unpack* provide an alternative to derived datatypes for grouping data. *MPI_Pack* copies the data to be sent, one block at a time, into a user-provided buffer. The buffer can then be sent and received. After the data is received, *MPI_Unpack* can be used to unpack it from the receive buffer. The syntax of *MPI_Pack* is

```
int MPI_Pack(
    void*      in_buf,          /* in */
    int        in_buf_count,    /* in */
    MPI_Datatype datatype,      /* in */
    void*      pack_buf,        /* out */
    int        pack_buf_sz,     /* in */
    int*       position_p,      /* in/out */
    MPI_Comm   comm            /* in */);
```

We could therefore pack the input data to the trapezoidal rule program with the following code:

```
char pack_buf[100];
int position = 0;
MPI_Pack(&a, 1, MPI_DOUBLE, pack_buf, 100, &position, comm);
MPI_Pack(&b, 1, MPI_DOUBLE, pack_buf, 100, &position, comm);
MPI_Pack(&n, 1, MPI_INT, pack_buf, 100, &position, comm);
```

The key is the position argument. When *MPI_Pack* is called, position should refer to the first available slot in *pack_buf*. When *MPI_Pack* returns, it refers to the first available slot after the data that was just packed, so after process 0 executes this code, all the processes can call *MPI_Bcast*:

```
MPI_Bcast(pack_buf, 100, MPI_PACKED, 0, comm);
```

Note that the MPI datatype for a packed buffer is `MPI_PACKED`. Now the other processes can unpack the data using: *MPI_Unpack*:

```
int MPI_Unpack(
    void*      pack_buf,          /* in */
    int        pack_buf_sz,       /* in */
    int*       position_p,        /* in/out */
    void*      out_buf,          /* out */
    int        out_buf_count,     /* in */
    MPI_Datatype datatype,       /* in */
    MPI_Comm   comm              /* in */);
```

This can be used by “reversing” the steps in *MPI_Pack*, that is, the data is unpacked one block at a time starting with position = 0. Write another Get input function for the trapezoidal rule program. This one should use *MPI_Pack* on process 0 and *MPI_Unpack* on the other processes.

- 21.
22. Time our implementation of the trapezoidal rule that uses *MPI_Reduce*. How will you choose n , the number of trapezoids ? How do the minimum times compare to the mean and median times ? What are the speedups ? What are the efficiencies ? On the basis of the data you collected, would you say that the trapezoidal rule is scalable ?
23. Although we don't know the internals of the implementation of *MPI_Reduce*, we might guess that it uses a structure similar to the binary tree we discussed. If this is the case, we would expect that its run-time would grow roughly at the rate of $\log_2(p)$, since there are roughly $\log_2(p)$ levels in the tree. (Here, $p = \text{comm_sz}$.) Since the run-time of the serial trapezoidal rule is roughly proportional to n , the number of trapezoids, and the parallel trapezoidal rule simply applies the serial rule to $\frac{n}{p}$ trapezoids on each process, with our assumption about *MPI_Reduce*, we get a formula for the overall run-time of the parallel trapezoidal rule that looks like $T_{\text{parallel}}(n, p) \approx a \times \frac{n}{p} + b \log_2(p)$ for some constants a and b .
 - a. Use the formula, the times you've taken in Exercise 3.22, and your favorite program for doing mathematical calculations to get a least-squares estimate of the values of a and b .
 - b. Comment on the quality of the predicted run-times using the formula and the values for a and b computed in part (a).
- 24.
- 25.
- 26.
27. Find the speedups and efficiencies of the parallel odd-even sort. Does the program obtain linear speedups? Is it scalable ? Is it strongly scalable ? Is it weakly scalable ?
28. Modify the parallel odd-even transposition sort so that the Merge functions simply swap array pointers after finding the smallest or largest elements. What effect does this change have on the overall run-time ?

4.1.1 Questão 16

Quando o compilador não é capaz de vetorizar automaticamente, ou vetoriza de forma ineficiente, o OpenMP provê a diretiva **omp simd**, com a qual o programador pode indicar um laço explicitamente para o compilador vetorizar. No código abaixo, a inclusão da cláusula **reduction** funciona de forma similar a flag **-ffast-math**, indicando que a redução na variável soma é segura e deve ser feita.

```
#pragma omp simd reduction(+:soma)
for (i=0; i<n; i++) {
    x = (i+0.5)*h;
    soma += 4.0/(1.0+x*x);
}
```

Por que não é necessário usar a cláusula **private(x)** neste caso mas seria caso a diretiva **omp simd** fosse combinada com a diretiva **omp parallel for** ?

4.2 Question 02

4.2.1 Questão 2

Modify the trapezoidal rule so that it will correctly estimate the integral even if *comm_sz* doesn't evenly divide *n*. (You can still assume that $n \geq \text{comm_sz}$).

Código antes da alteração

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>

double trap(double a, double b, long int n);

double function(double x);

int main(int argc, char *argv[])
{
    int my_rank;
    int comm_sz;

    int message_tag = 0;
    // receber argumentos globais via linha de comando
    MPI_Init(&argc, &argv);

    const double A = atof(argv[1]);
    const double B = atof(argv[2]);
    const int N = atoi(argv[3]);

    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
    MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);

    double h = (B - A) / N;

    int local_n = N / comm_sz;
    /*
     * Como podemos observar que quando a divisão
     * inteira N/comm_sz o n local, não corresponde
     * ao número total de trapézios desejados.
     * Ex: 1024/3 == 341, sendo que 341*3 == 1023
     */
    double local_a = A + my_rank * local_n * h;
    double local_b = local_a + local_n * h;
    double result_integral = trap(local_a, local_b, local_n);

    if (my_rank != 0)
    {
        MPI_Send(&result_integral, 1, MPI_DOUBLE, 0, message_tag, MPI_COMM_WORLD);
    }
    else // my_rank == 0
    {
        for (int source = 1; source < comm_sz; source++)
        {
            double result_integral_source;
            MPI_Recv(&result_integral_source, 1, MPI_DOUBLE, source, message_tag, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
            result_integral += result_integral_source;
        }
        printf("with n = %d trapezoids, our estimate\n", N);
        printf("of the integral from %f to %f = %.15e\n", A, B, result_integral);
    }

    MPI_Finalize();

    return 0;
}

double trap(double a, double b, long int n)
{
    double h = (b - a) / n;

    double approx = (function(a) + function(b)) / 2.0;
    for (int i = 1; i < n - 1; i++)
    {
        // printf("f_%i \n", i);
        double x_i = a + i * h;

        approx += function(x_i);
    }
}
```

```

    }

    return h * approx;
}

double function(double x) { return x * x; }

```

compilação e execução

```

mpicc trap_rule.c -o trap
# a =0, b =3, n = 1024
mpiexec -n 4 ./trap 0 3 1024

```

Alteração no código

Sabendo que pode ocorrer divisões cujo o número não pertence aos números racionais, por tanto se faz necessário utilizar a seguinte expressão:

$$\frac{N}{D} = a \times D + \text{resto}$$

exemplo:

$$\frac{1024}{3} = 341 \times 3 + 1$$

Podemos resolver o problema simplesmente despejando os trapézios restantes em um processo ou podemos dividir a "carga" entre os processos. Dividir a carga entre os processos de forma mais igualitária possível, foi resolvido no [exercício 1 capítulo 1](#).

```

struct range
{
    int first;
    int last;
};

struct range new_range_2(int thread_index, int p, int n)
{
    struct range r;

    int division = n / p;
    int rest = n % p;

    if (thread_index < rest)
    {
        r.first = thread_index * (division + 1);
        r.last = r.first + division + 1;
    }
    else
    {
        r.first = thread_index * division + rest;
        r.last = r.first + division;
    }

    return r;
}
...

int main(int argc, char *argv[])
{
    ...
    /*
    onde thread_index é o rank,
    o número de cores é o número de processos,
    o tamanho do vetor é o número de trapézios
    */
    struct range r = new_range_2(my_rank, comm_sz, N);

    double h = (B - A) / N;

    /*
    perceba que o número local de trapézios é
    o tamanho do intervalo calculado
    */
    int local_n = r.last - r.first;
    double local_a = A + r.first * h;
    double local_b = A + r.last * h;
    double result_integral = trap(local_a, local_b, local_n);

    printf("local n: %i local a: %f local b %f \n", local_n, local_a, local_b);
}

```

Simulações feitas offline:

```
mpicc trap_rule.c -o trap; mpiexec -n 3 ./trap 1 3 1024
local n: 342 local a: 1.000000 local b 1.667969
local n: 341 local a: 1.667969 local b 2.333984
local n: 341 local a: 2.333984 local b 3.000000
with n = 1024 trapezoids, our estimate
of the integral from 1.000000 to 3.000000 = 8.633069768548012e+00
```

```
mpicc trap_rule.c -o trap; mpiexec -n 4 ./trap 1 3 2024
local n: 506 local a: 1.000000 local b 1.500000
local n: 506 local a: 1.500000 local b 2.000000
local n: 506 local a: 2.000000 local b 2.500000
local n: 506 local a: 2.500000 local b 3.000000
with n = 2024 trapezoids, our estimate
of the integral from 1.000000 to 3.000000 = 8.645439504647232e+00
```

```
mpicc trap_rule.c -o trap; mpiexec -n 3 ./trap 0 3 1024
local n: 342 local a: 0.000000 local b 1.001953
local n: 341 local a: 1.001953 local b 2.000977
local n: 341 local a: 2.000977 local b 3.000000
with n = 1024 trapezoids, our estimate
of the integral from 0.000000 to 3.000000 = 8.959068736061454e+00
```

código final

```
#include <mpi.h>
#include <stdio.h>
#include <stdlib.h>

struct range
{
    int first;
    int last;
};

struct range new_range_2(int thread_index, int p, int n)
{
    struct range r;

    int division = n / p;
    int rest = n % p;

    if (thread_index < rest)
    {
        r.first = thread_index * (division + 1);
        r.last = r.first + division + 1;
    }
    else
    {
        r.first = thread_index * division + rest;
        r.last = r.first + division;
    }

    return r;
}

double trap(double a, double b, long int n);

double function(double x);

int main(int argc, char *argv[])
{
    int my_rank;
    int comm_sz;

    int message_tag = 0;

    MPI_Init(&argc, &argv);

    const double A = atof(argv[1]);
    const double B = atof(argv[2]);
    const int N = atoi(argv[3]);

    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
    MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);

    struct range r = new_range_2(my_rank, comm_sz, N);

    double h = (B - A) / N;

    int local_n = r.last - r.first;
    double local_a = A + r.first * h;
    double local_b = A + r.last * h;
    double result_integral = trap(local_a, local_b, local_n);

    printf("local n: %i local a: %f local b %f \n", local_n, local_a, local_b);
```

```

if (my_rank != 0)
{
    MPI_Send(&result_integral, 1, MPI_DOUBLE, 0, message_tag, MPI_COMM_WORLD);
}
else // my_rank ==0
{
    for (int source = 1; source < comm_sz; source++)
    {
        double result_integral_source;
        MPI_Recv(&result_integral_source, 1, MPI_DOUBLE, source, message_tag, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
        result_integral += result_integral_source;
    }
    printf("with n = %d trapezoids, our estimate\n", N);
    printf("of the integral from %f to %f = %.15e\n", A, B, result_integral);
}

MPI_Finalize();

return 0;
}

double trap(double a, double b, long int n)
{
    double h = (b - a) / n;
    // printf("H: %f\n",h);

    double approx = (function(a) + function(b)) / 2.0;
    for (int i = 1; i < n - 1; i++)
    {
        // printf("f_%i \n", i);
        double x_i = a + i * h;

        approx += function(x_i);
    }

    return h * approx;
}

double function(double x) { return x * x; }

```