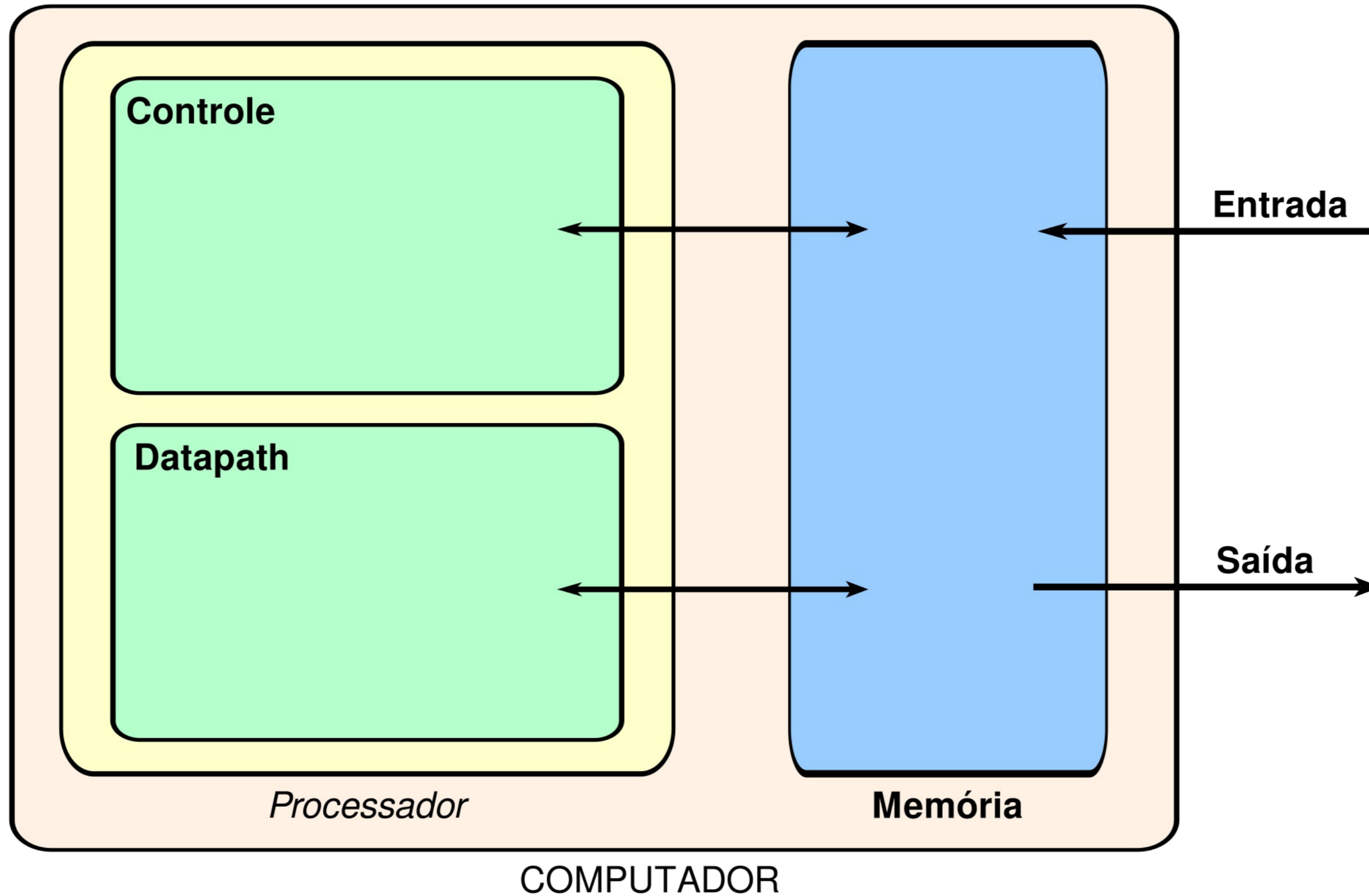


MOTIVAÇÃO



MOTIVAÇÃO



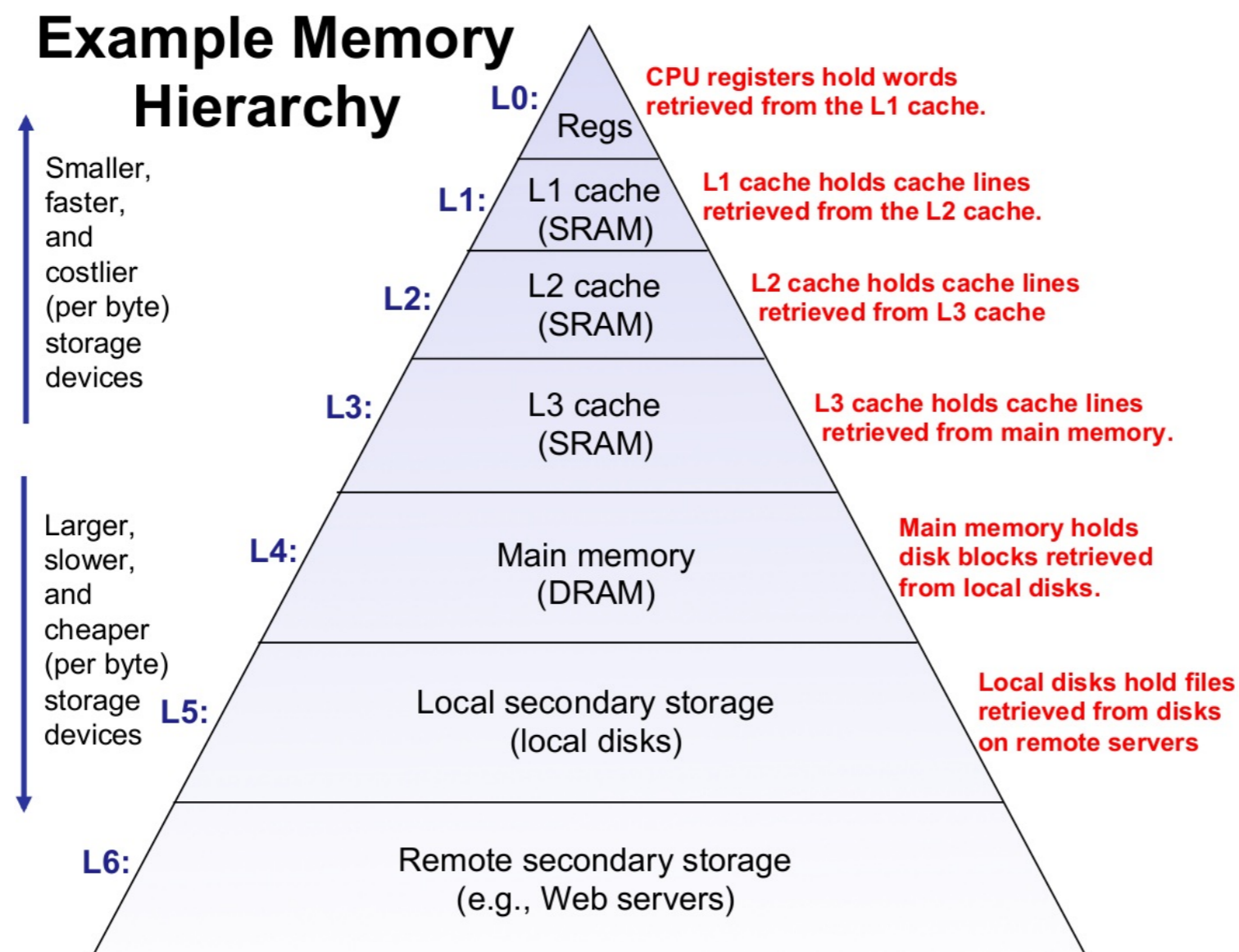
Entender a memória importa!

- **Muitos programas dependem basicamente da memória RAM, mas a memória RAM não existe, é uma abstração limitada. O que existe é um complexo sistema hierárquico de memórias diferentes que tenta nos dar a impressão de alta capacidade de armazenamento e alta performance.**
 - * **Deve ser gerenciada: alocada e desalocada**
- **Bugs no uso da memória são perniciosos!**
 - * **Efeitos são distantes no tempo e no espaço**
 - * **Difícil de debugar**
- **A performance da memória não é uniforme**
 - * **Memória virtual e cache afetam a performance**
 - * **Adaptar um programa às características da memória pode levar a ganhos expressivos de performance**

A memória RAM é uma abstração!

- Muitos programas dependem basicamente da memória RAM, mas a memória RAM não existe, é uma abstração limitada. O que existe é um complexo sistema hierárquico de memórias diferentes que tenta nos dar a impressão de alta capacidade de armazenamento e alta performance.

* Deve ser gerenciada: alocada e desalocada



Bugs no referenciamento da memória são perniciosos!

- Bugs no uso da memória são perniciosos!

- * Efeitos são distantes no tempo e no espaço
- * Difícil de debugar

```
1 /**
2  * CR6.190A: Arquitetura de Computadores
3  * https://cursos.computacaoraiz.com.br
4  *
5  * erro_de_memoria.c
6  * 0 array não tem fim?
7  */
8
9 #include <stdio.h>
10
11 int main(void)
12 {
13     int a[2] = {22, 33};
14
15     for (int i = 0; i <= 10; ++i)
16     {
17         a[i] = i + 1;
18         printf("a[%d] = %d\n", i, a[i]);
19     }
20
21     return 0;
22 }
```

```
[abrantestasf@cosmos ~/cr6.190a/introducao]$ ./erro_de_memoria
a[0] = 1
a[1] = 2
a[2] = 3
a[3] = 4
a[4] = 5
a[5] = 6
a[6] = 7
a[7] = 8
a[8] = 9
a[9] = 10
a[10] = 11
*** stack smashing detected ***: terminated
Aborted (core dumped)
```


Bugs no referenciamento da memória são perniciosos!

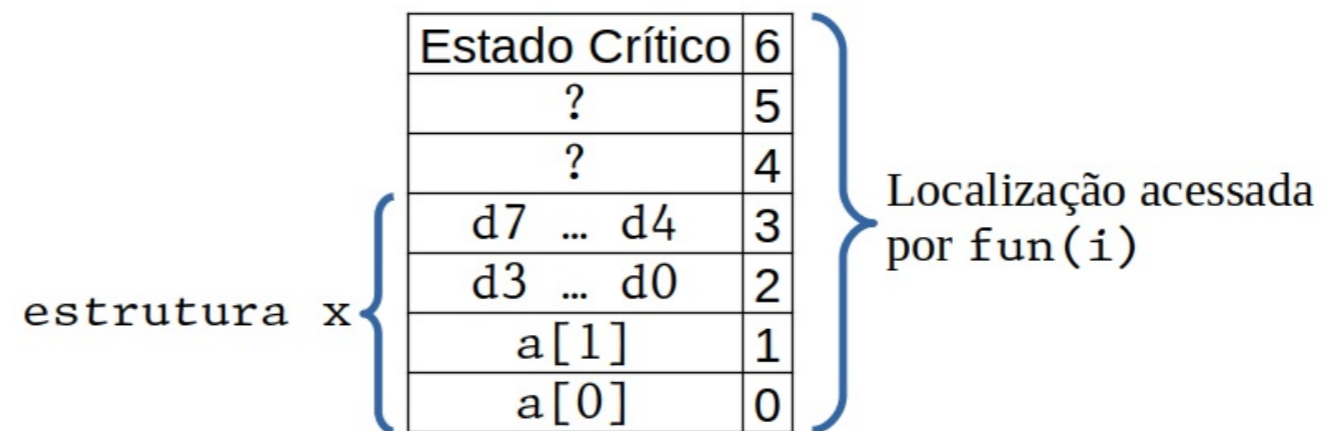
- Bugs no uso da memória são perniciosos!

- * Efeitos são distantes no tempo e no espaço

- * Difícil de debugar

```
1 /**
2  * CR6.190A: Arquitetura de Computadores
3  * https://cursos.computacaoraiz.com.br
4  *
5  * erro_de_memoria2.c
6  * Erro intermitente é o pior.
7  *
8  * Código retirado de:
9  * Computer Systems: A Programmer's Perspective, 3ª ed (CS:APP3e)
10 * https://csapp.cs.cmu.edu/3e/instructors.html
11 */
12
13 #include <stdio.h>
14
15 typedef struct
16 {
17     int a[2];
18     double d;
19 } estrutura;
20
21 double fun(int i)
22 {
23     volatile estrutura x;
24     x.d = 3.14;
25     x.a[i] = 1073741824;
26     return x.d;
27 }
28
29 int main(void)
30 {
31     for (int i = 0; i < 10; ++i)
32         printf("fun(%d) = %.2f\n", i, fun(i));
33 }
```

```
[abrantesasf@cosmos ~/cr6.190a/introducao]$ ./erro_de_memoria2
fun(0) = 3.14
fun(1) = 3.14
fun(2) = 3.14
fun(3) = 2.00
fun(4) = 3.14
fun(5) = 3.14
*** stack smashing detected ***: terminated
Aborted (core dumped)
```



Bugs no referenciamento da memória são perniciosos!

- Bugs no uso da memória são perniciosos!
- C e C++ não fornecem nenhuma proteção de memória
 - * não verificam limites de arrays
 - * não verificam ponteiros inválidos
 - * é fácil abusar de malloc/free
- Levam a bugs estranhos e difíceis de encontrar
 - * se um bug terá ou não efeito depende do seu sistema e do compilador, tornando tudo mais difícil (erros intermitentes)
 - * ação à distância (tempo e espaço): o objeto corrompido pode não ter relação lógica nenhuma com o objeto sendo acessado
- Como evitar?
 - * entender como a memória funciona e como prevenir
 - * usar ferramentas para diagnóstico (ex.: valgrind)
 - * usar Python, Java, Ruby, ML... 🙄

A performance da memória não é uniforme!

- **Programadores geralmente estão mais preocupados em saber se o algoritmo é correto, se as estruturas de dados são adequadas, e com a complexidade assintótica de seus programas (big-O). Isso é perfeito e deve ser realmente a preocupação inicial, mas:**
 - * **É necessário também otimizar para o baixo nível, para o uso de memória**
 - * **É importante entender o que o computador faz, como ele faz, e o que faz ele rodar mais rápido ou mais devagar**
 - * **É fácil obter até 10x mais performance dependendo apenas de como o código é escrito**
 - * **Otimizar: algoritmo, estrutura de dados, procedimentos, loops**
- **A performance da memória não é uniforme**
 - * **Memória virtual e cache afetam a performance**
 - * **Adaptar um programa às características da memória pode levar a ganhos expressivos de performance**
 - * **Fatores constantes também importam**
 - * **É necessário entender o funcionamento para otimizar corretamente**

A performance da memória não é uniforme!

$$A_{4096 \times 4096} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & \cdots & a_{0,4095} \\ a_{1,0} & a_{1,1} & a_{1,2} & \cdots & a_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{4095,0} & a_{4095,1} & a_{4095,2} & \cdots & a_{4095,4095} \end{pmatrix}$$

$$B_{4096 \times 4096} = \begin{pmatrix} b_{0,0} & b_{0,1} & b_{0,2} & \cdots & b_{0,4095} \\ b_{1,0} & b_{1,1} & b_{1,2} & \cdots & b_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{4095,0} & b_{4095,1} & b_{4095,2} & \cdots & b_{4095,4095} \end{pmatrix}$$

$$B = A$$

(16.777.216 operações de cópia)

A performance da memória não é uniforme!

$$\begin{pmatrix} b_{0,0} & b_{0,1} & b_{0,2} & \cdots & b_{0,4095} \\ b_{1,0} & b_{1,1} & b_{1,2} & \cdots & b_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{4095,0} & b_{4095,2} & b_{4095,3} & \cdots & b_{4095,4095} \end{pmatrix} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & \cdots & a_{0,4095} \\ a_{1,0} & a_{1,1} & a_{1,2} & \cdots & a_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{4095,0} & a_{4095,2} & a_{4095,3} & \cdots & a_{4095,4095} \end{pmatrix} \quad \begin{pmatrix} b_{0,0} & b_{0,1} & b_{0,2} & \cdots & b_{0,4095} \\ b_{1,0} & b_{1,1} & b_{1,2} & \cdots & b_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{4095,0} & b_{4095,2} & b_{4095,3} & \cdots & b_{4095,4095} \end{pmatrix} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & \cdots & a_{0,4095} \\ a_{1,0} & a_{1,1} & a_{1,2} & \cdots & a_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{4095,0} & a_{4095,2} & a_{4095,3} & \cdots & a_{4095,4095} \end{pmatrix}$$

```
1 /**
2  * CR6.190A: Arquitetura de Computadores
3  * https://cursos.computacaoraiz.com.br
4  *
5  * copiar_matrizes2.c
6  * B = A
7  *
8  */
9
10 #include <stdio.h>
11
12 #define N 4096
13
14 int main(void)
15 {
16     static int matrizA[N][N];
17     static int matrizB[N][N];
18
19     for (int l = 0; l < N; ++l)
20         for (int c = 0; c < N; ++c)
21             matrizA[l][c] = 1;
22
23     for (int c = 0; c < N; ++c)
24         for (int l = 0; l < N; ++l)
25             matrizB[l][c] = matrizA[l][c];
26 }
```

$$B = A$$
$$O(n^2)$$

```
1 /**
2  * CR6.190A: Arquitetura de Computadores
3  * https://cursos.computacaoraiz.com.br
4  *
5  * copiar_matrizes1.c
6  * B = A
7  *
8  */
9
10 #include <stdio.h>
11
12 #define N 4096
13
14 int main(void)
15 {
16     static int matrizA[N][N];
17     static int matrizB[N][N];
18
19     for (int l = 0; l < N; ++l)
20         for (int c = 0; c < N; ++c)
21             matrizA[l][c] = 1;
22
23     for (int l = 0; l < N; ++l)
24         for (int c = 0; c < N; ++c)
25             matrizB[l][c] = matrizA[l][c];
26 }
```

```
[abrantestasf@cosmos ~/cr6.190a/introducao]$ time ./copiar_matrizes2
```

```
real    0m0,543s
user    0m0,489s
sys     0m0,052s
```

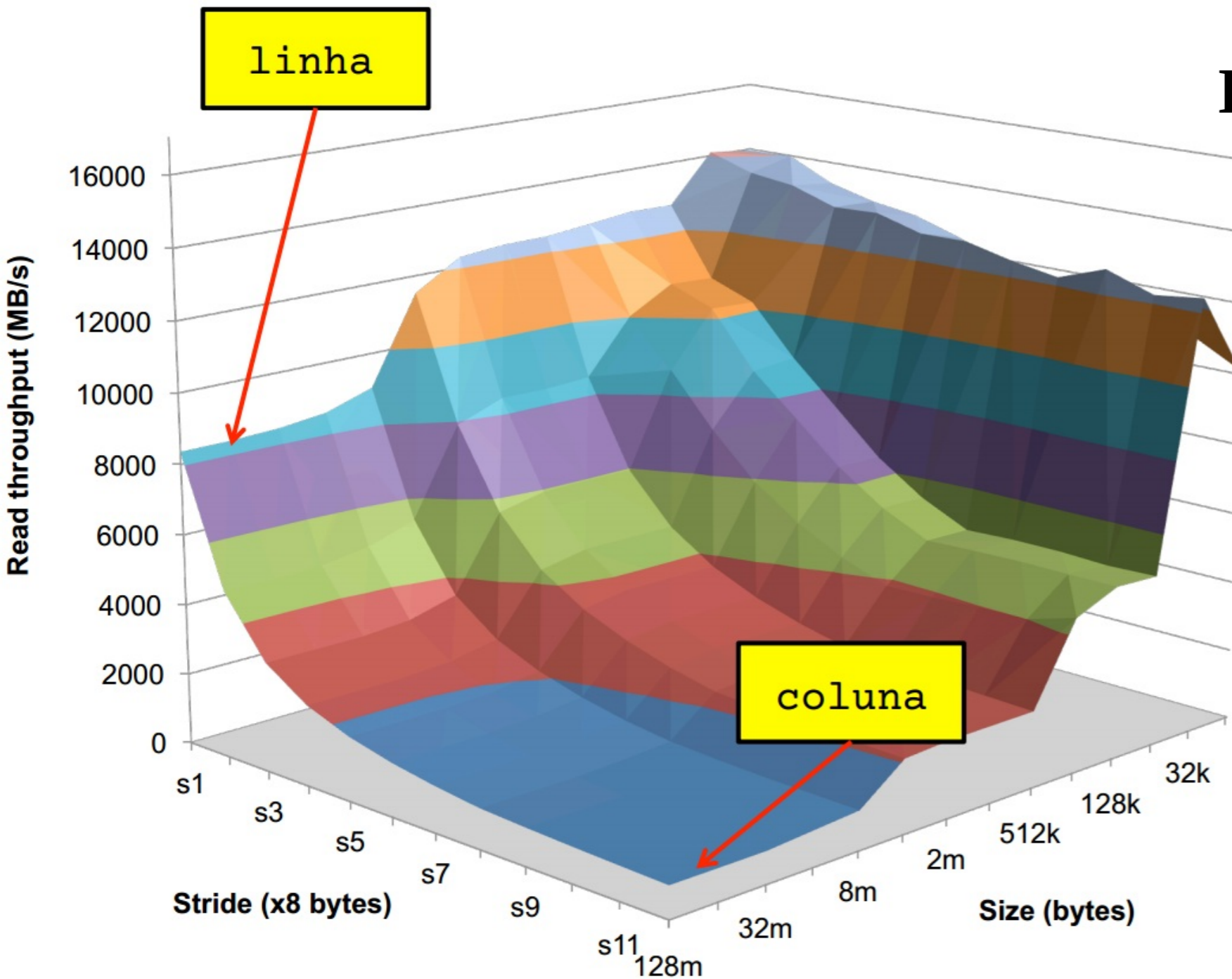
```
[abrantestasf@cosmos ~/cr6.190a/introducao]$ time ./copiar_matrizes1
```

```
real    0m0,133s
user    0m0,097s
sys     0m0,037s
```

Copiar por coluna foi 4 VEZES MAIS LENTO do que copiar por linha! Por que isso ocorreu?

A performance da memória não é uniforme!

$$\begin{pmatrix} b_{0,0} & b_{0,1} & b_{0,2} & \cdots & b_{0,4095} \\ b_{1,0} & b_{1,1} & b_{1,2} & \cdots & b_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{4095,0} & b_{4095,2} & b_{4095,3} & \cdots & b_{4095,4095} \end{pmatrix} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & \cdots & a_{0,4095} \\ a_{1,0} & a_{1,1} & a_{1,2} & \cdots & a_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{4095,0} & a_{4095,2} & a_{4095,3} & \cdots & a_{4095,4095} \end{pmatrix} \quad \begin{pmatrix} b_{0,0} & b_{0,1} & b_{0,2} & \cdots & b_{0,4095} \\ b_{1,0} & b_{1,1} & b_{1,2} & \cdots & b_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_{4095,0} & b_{4095,2} & b_{4095,3} & \cdots & b_{4095,4095} \end{pmatrix} = \begin{pmatrix} a_{0,0} & a_{0,1} & a_{0,2} & \cdots & a_{0,4095} \\ a_{1,0} & a_{1,1} & a_{1,2} & \cdots & a_{1,4095} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{4095,0} & a_{4095,2} & a_{4095,3} & \cdots & a_{4095,4095} \end{pmatrix}$$



Para entender porque a lentidão, é necessário:

- * Entender a hierarquia da memória
- * Entender como a performance depende dos padrões de acesso à memória (incluindo como percorrer um array multidimensional)

3ª verdade nua e crua: você precisa entender a arquitetura da memória!



Para ser um bom programador:

- Entender a hierarquia de memória**
- Entender como a arquitetura da memória e linguagens como C podem levar a bugs de referenciamento de memória que são complicados de debugar e que podem ser distantes do tempo e espaço**
- Entender que a performance da memória não é uniforme e que é necessário otimizar para o baixo nível também**

"Entender a representação em nível de máquina, na memória, das estruturas de dados e como elas funcionam faz uma grande diferença na sua habilidade de evitar e lidar com problemas de referenciamento de memória e vulnerabilidades no seu programa." Randal Bryant