

Introdução a OpenMP

MC-SD02-I

OpenMP Programação Avançada

MC-SD02-II

Carla Osthoff
LNCC/MCTI
osthoff@lncc.br

Material do curso: MC-SD02-I

<http://www.cenapad-rj.lncc.br/tutoriais/materiais-hpc/semana-sdumont/>

http://www.openmp.org/ Home - OpenMP OpenMP Books - OpenMP OpenMP Compilers & Tools - ...

OpenMP®

Enabling HPC since 1997

The OpenMP API specification for parallel programming

Home

Specifications

Blog

Community

Resources

News & Events

About



OpenMP ARB Members

The OpenMP API is jointly defined by a group of major computer hardware and software vendors and major parallel computing user facilities.

[READ MORE](#)

Latest News



An Interview with InsideHPC

Dec 14, 2017

View the InsideHPC video from SC17 where Michael Klemm discusses the OpenMP ARB, the latest Technical Report 6 (TR6) and asks for feedback via the OpenMP Forum. [more](#)

SC17 In-Booth Talks Video and Slides Available

Nov 27, 2017

Twelve in-booth talks from SC17 - Denver are now viewable from our [SC17 Presentations Page](#).

Release of OpenMP Technical Report 6 (TR6) Addresses Top User Requests

Nov 13, 2017

OpenMP ARB Technical Report 6 (TR6) extends TR4 adding a number of key features and is a preview of OpenMP 5.0, expected in November 2018. [more](#)

OpenMPCon 2017 Presentations Now Available for Download

Using OpenMP – The Next Step

Oct 01, 2017

OpenMP Accelerator Support for GPUs

@OpenMP_ARB

OpenMP ARB @OpenMP_ARB

Great work!

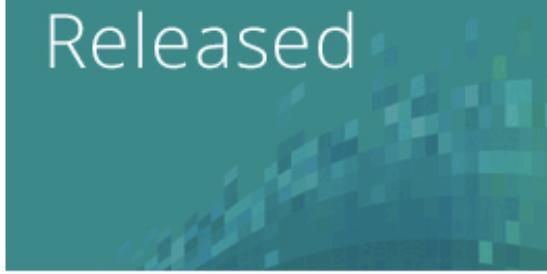
13 Feb 2018

OpenMP ARB @OpenMP_ARB

"As we prepare for the next gen of supercomputers and #GPUs, #OpenMP is growing to meet challenges of programming scientific apps in a world of accelerators, unified memory, and explicitly

OpenMP

OpenMP 5.1
Released



OpenMP ARB releases OpenMP 5.1 with vital usability enhancements

While the primary focus has been enhancements, clarifications and corrections to the 5.0 specification, several useful new features have been added such as support for interoperability with lower level APIs like CUDA and HIP.



SC20

Everywhere | more
we are | than hpc.

Join OpenMP
@ SC20
Virtual



OpenMP @ SC20

We're gearing up for SC'20 and delighted to have three OpenMP tutorials included in the program: Advanced OpenMP: Host Performance and 5.0 Features, Programming your GPU with OpenMP: A hands-on Introduction, and The OpenMP Common Core: A hands-on Introduction.

ECP OpenMP Hackathons 2020

ECP SOLLVE, in conjunction with ORNL and NERSC, are organizing two OpenMP Hackathons in July and August. We encourage participation of teams interested in using OpenMP to port and optimize their applications on GPUs, and in using OpenMP for energy-efficient processor architectures.



The OpenMP API specification for parallel programming

Home

Specifications

Blog

Community ▾

Resources ▾

News & Events ▾

OpenMP Books

- [OpenMP Common Core: Making OpenMP Simple Again](#) – by Tim Mattson, Helen He, Alice Koniges (2019)
- [Using OpenMP – The Next Step](#) – by Ruud van der Pas, Eric Stotzer and Christian Terboven (2017)
- [Using OpenMP – Portable Shared Memory Parallel Programming](#) – by Chapman, Jost, and Van Der Pas (2007)
- [Parallel Programming in OpenMP](#) – by Rohit Chandra et al.
- [Parallel Programming in C with MPI and OpenMP](#) – by Michael J. Quinn.
- [Parallel Programming Patterns: Working with Concurrency in OpenMP, MPI, Java, and OpenCL](#) – by Timothy C. Sanders
- [An Introduction to Parallel Programming with OpenMP, PThreads and MPI](#) – by Robert Cook
- [The International Journal of Parallel Programming](#) – Issues and articles devoted to OpenMP.

OpenMP Compilers & Tools

[Home](#) > [Resources](#) > [OpenMP Compilers & Tools](#)

GNU

GCC

The free and open-source GNU Compiler Collection (GCC) supports among others Linux, Solaris, AIX, MacOSX, Windows, FreeBSD, NetBSD, OpenBSD, DragonFly BSD, HPUX, RTEMS, for architectures such as x86_64, PowerPC, ARM, and many more.

C/C++/Fortran

Code offloading to NVIDIA GPUs (nvptx) and the AMD Radeon (GCN) GPUs Fiji and Vega is supported on Linux.

OpenMP 4.0 is fully supported for C, C++ and Fortran since GCC 4.9; OpenMP 4.5 is fully supported for C and C++ since GCC 6 and partially for Fortran since GCC 7. OpenMP 5.0 is partially supported for C and C++ since GCC 9 and extended in GCC 10. The next release, GCC 11, will fully support OpenMP 4.5 for Fortran and extend the OpenMP 5.0 support for C, C++ and Fortran; the devel/omp/gcc-10 (og10) branch augments the GCC 10 branch with OpenMP and offloading features, mostly from GCC 11 development branch.

Compile with -fopenmp to enable OpenMP.

GCC binary builds are provided by Linux distributions, often with offloading support provided by additional packages, and by multiple entities for other platforms – and you can build it from source.

Releases and release notes: <https://gcc.gnu.org/>

OpenMP documentation: <https://gcc.gnu.org/onlinedocs/libgomp/>

Building and using GCC for offloading: <https://gcc.gnu.org/wiki/Offloading>

OpenMP Compilers & Tools

[Home](#) > [Resources](#) > [OpenMP Compilers & Tools](#)

Intel

C/C++/Fortran

Windows, Linux, and MacOSX.

- OpenMP 3.1 C/C++/Fortran fully supported in version 12.0, 13.0, 14.0 compilers
- OpenMP 4.0 C/C++/Fortran supported in version 15.0 and 16.0 compilers
- OpenMP 4.5 C/C++/Fortran supported in version 17.0, 18.0, and 19.0 compilers
- OpenMP 4.5 and subset of OpenMP 5.0 in C/C++/Fortran compiler classic 2021.1
- OpenMP 4.5 and subset of OpenMP 5.1 supported in oneAPI DPC++/C++ compiler 2021.1 under -fopenmp -fopenmp-targets=spir64
- OpenMP 4.5 and subset of OpenMP 5.1 supported in oneAPI Fortran compiler (Beta) under -fopenmp -fopenmp-targets=spir64.

Compile with -Qopenmp on Windows, or just -qopenmp or -fopenmp on Linux or Mac OSX

Compile with -fopenmp -fopenmp-targets=spir64 on Windows and Linux for offloading support



The OpenMP API specification for parallel programming

Home

Specifications

Blog

Community ▾

Resources ▾

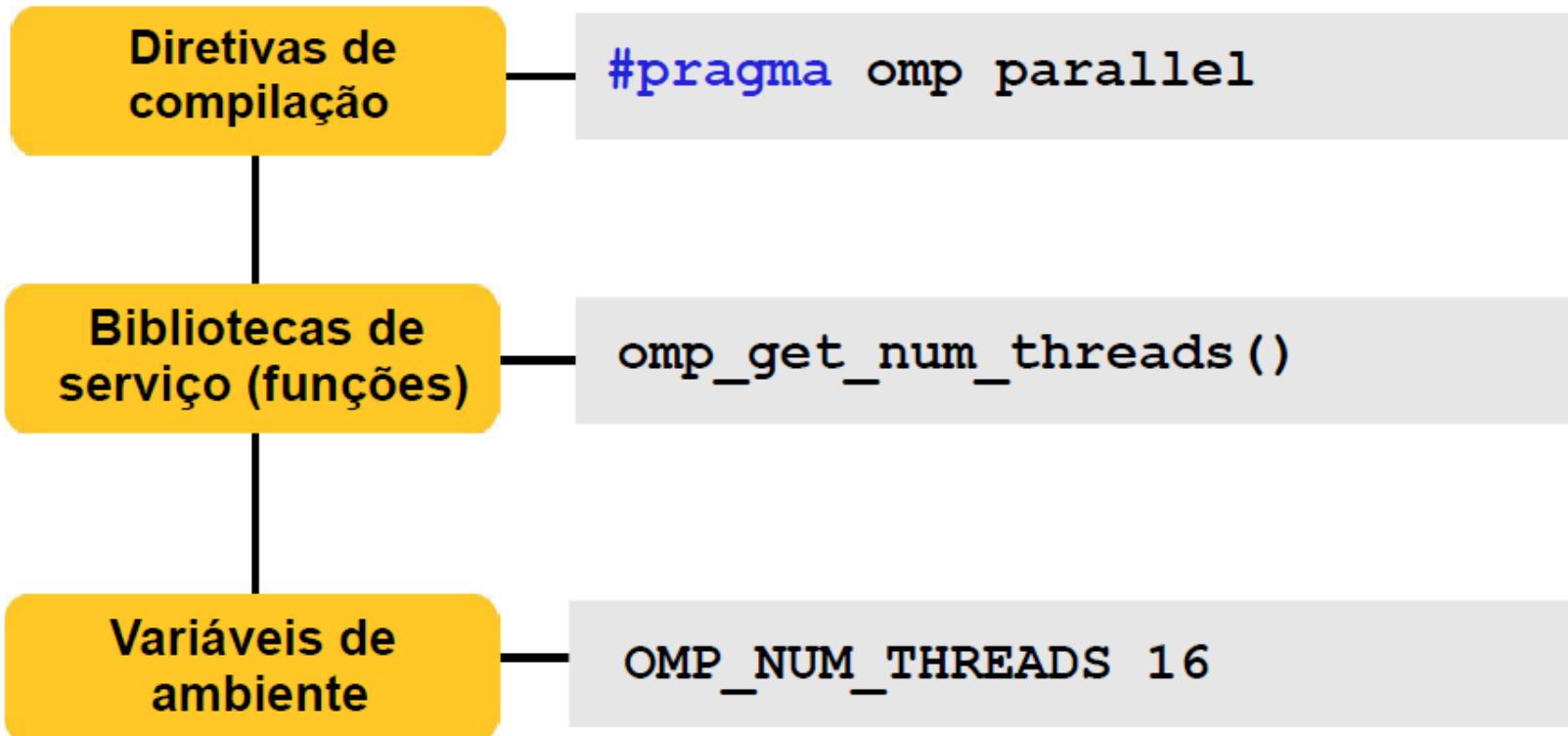
News & Events ▾

► OpenMP Is:

- An Application Program Interface (API) that may be used to explicitly direct ***multi-threaded, shared memory*** parallelism.
- Comprised of three primary API components:
 - Compiler Directives
 - Runtime Library Routines
 - Environment Variables
- An abbreviation for: **Open Multi-Processing**



Estrutura do OpenMP API



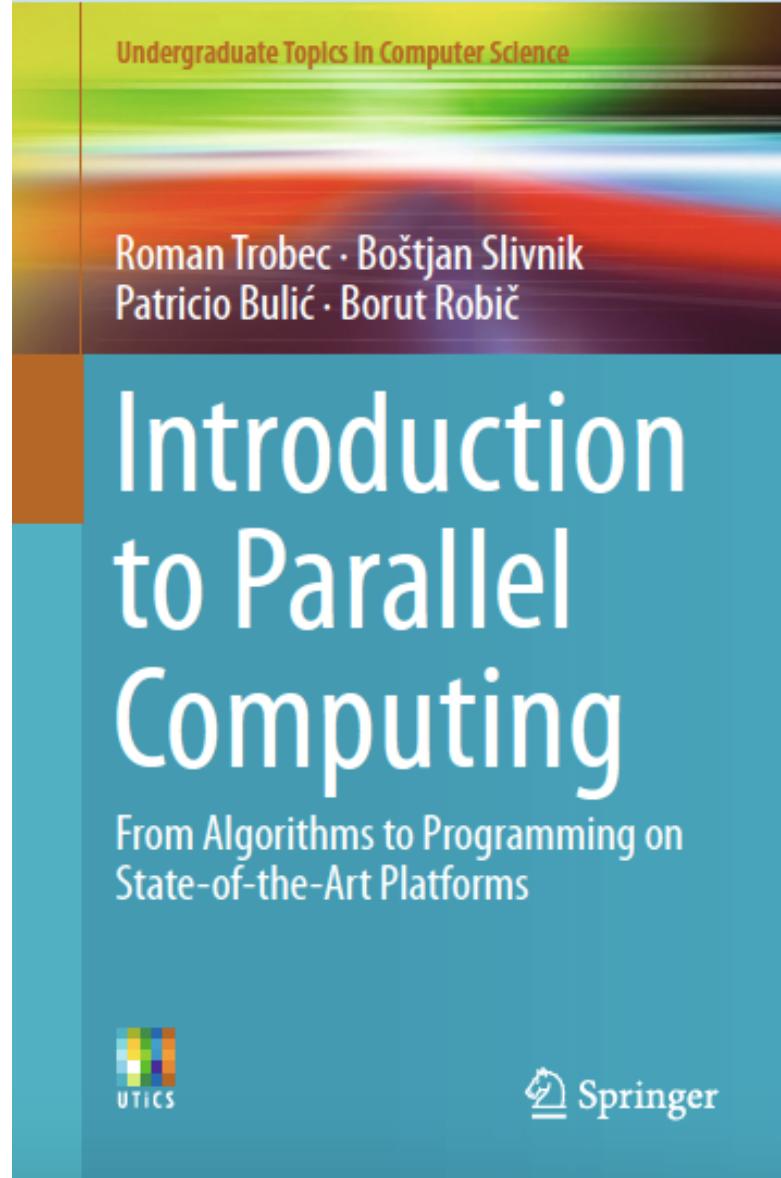
[Home](#)[Specifications](#)[Blog](#)[Community](#) ▾[Resources](#) ▾[News & Events](#) ▾[About](#) ▾

Tutorials & Articles

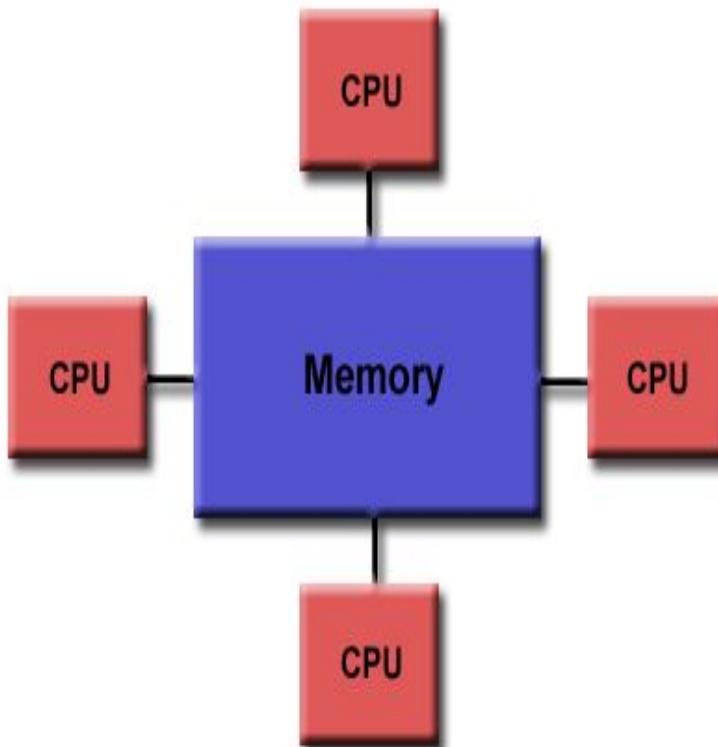
[Home](#) > [Resources](#) > [Tutorials & Articles](#)

Tutorial:<https://computing.llnl.gov/tutorials/openMP/>

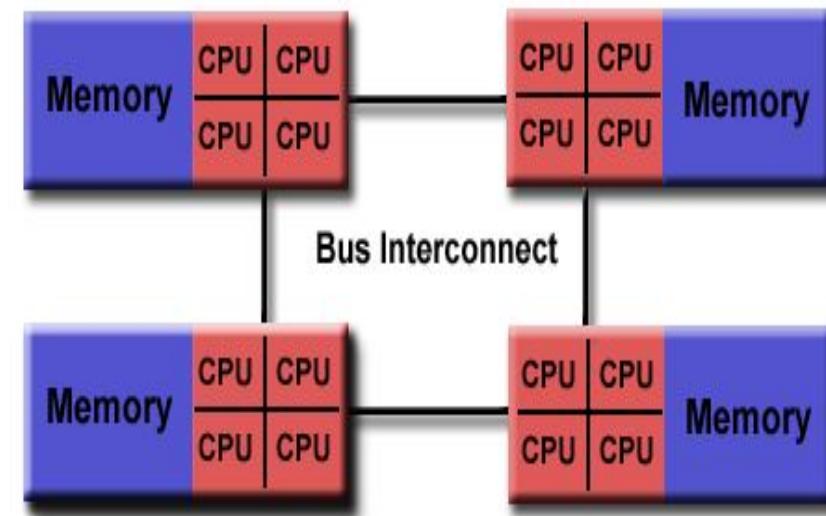
Capítulo 3 : Programming Multi-core and Shared Memory Multiprocessors Using OpenMP



Modelo de Programação do OpenMP : Modelo de Memória Compartilhada



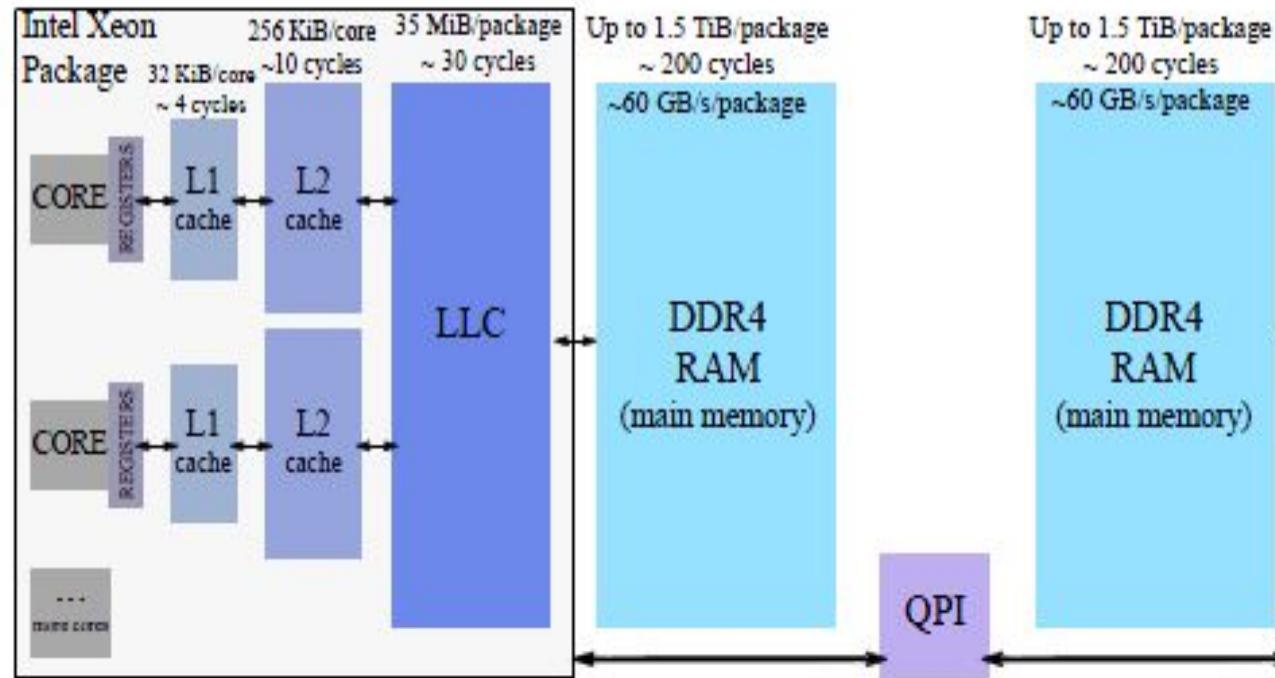
Uniform Memory Access



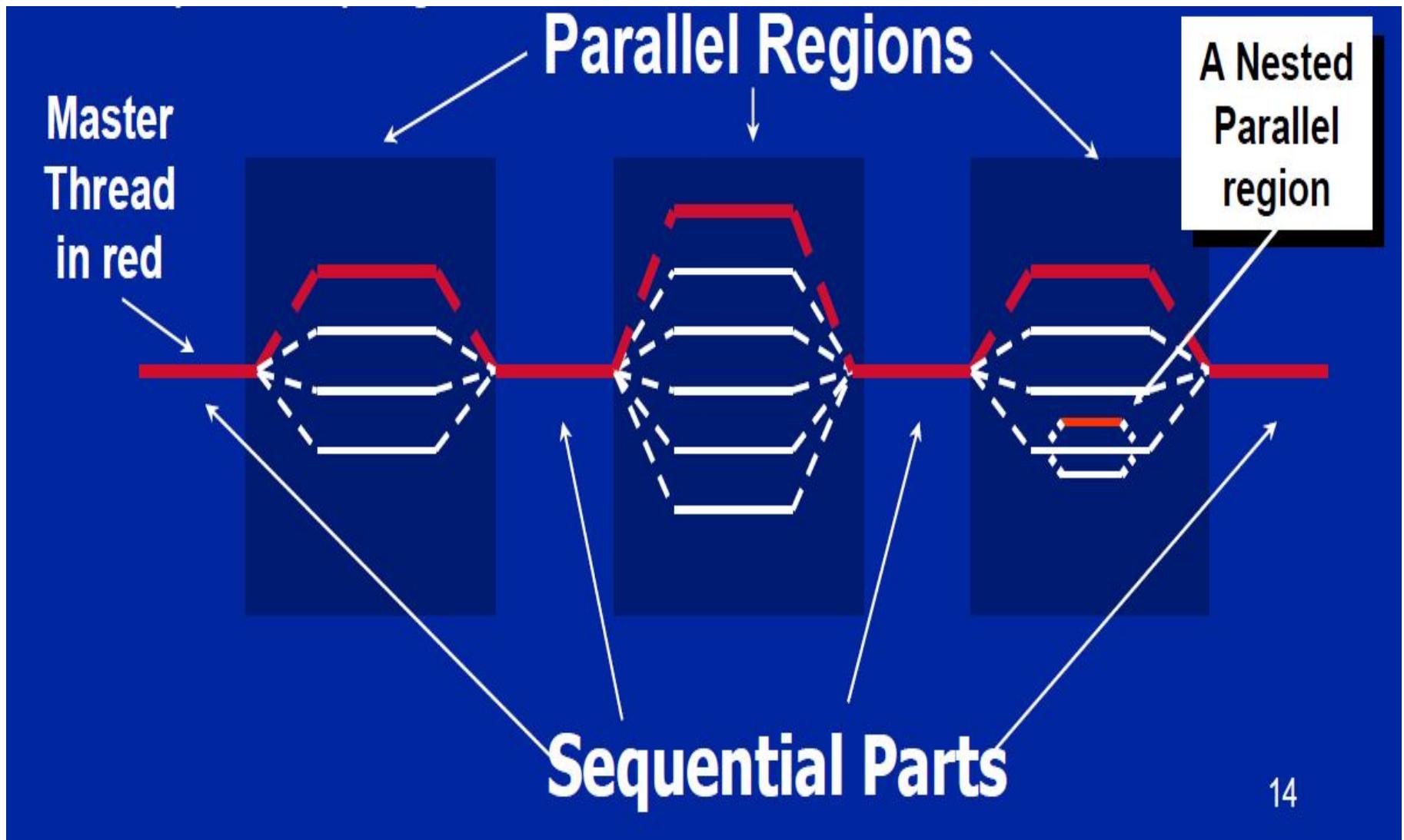
Non-Uniform Memory Access

INTEL XEON CPU: MEMORY ORGANIZATION

- ▷ Hierarchical cache structure
- ▷ Two-way processors have NUMA architecture



OpenMP utiliza Modelo de programação “Fork-Join”



Diretivas de Compilação

- **Diretivas** - Consiste em uma linha de código com significado “especial” para o compilador.

C/C++

```
#pragma omp parallel
```

FORTRAN

```
!$OMP OMP PARALLEL
```

Compiler / Platform	Compiler Commands	OpenMP Flag
Intel Linux	<code>icc</code> <code>icpc</code> <code>ifort</code>	<code>-qopenmp</code>
GNU Linux IBM Blue Gene Sierra, CORAL EA	<code>gcc</code> <code>g++</code> <code>g77</code> <code>gfortran</code>	<code>-fopenmp</code>
PGI Linux Sierra, CORAL EA	<code>pgcc</code> <code>pgCC</code> <code>pgf77</code> <code>pgf90</code>	<code>-mp</code>
Clang Linux Sierra, CORAL EA	<code>clang</code> <code>clang++</code>	<code>-fopenmp</code>
IBM XL Blue Gene	<code>bgxlcr</code> , <code>bgccr</code> <code>bgxlCr</code> , <code>bgxlC++r</code> <code>bgxlC89r</code> <code>bgxlC99r</code> <code>bgxlfr</code> <code>bgxlff90r</code> <code>bgxlff95r</code> <code>bgxlff2003r</code>	<code>-qsmp=omp</code>
IBM XL Sierra, CORAL EA	<code>xlc_r</code> <code>xlc_r</code> , <code>xlc++_r</code> <code>xlf_r</code> <code>xlf90_r</code> <code>xlf95_r</code> <code>xlf2003_r</code> <code>xlf2008_r</code>	<code>-qsmp=omp</code>

Instalação do material na sua conta no SDUMONT:

```
$ sftp user@login.sdumont.lncc.br
$ put OpenMP_2021_exemplos.tar
$ exit
```

```
$ ssh user@login.sdumont.lncc.br
$ cp OpenMP_2021_exemplos.tar $SCRATCH/.
$ cd $SCRATCH
$ tar xvf OpenMP_2021_exemplos.tar
```

Tem que atualizar fila de submissão (de treinamento para cpu_dev) no arquivo run.sh e o compilador para gcc/8.3 (também no arquivo compila.sh)

```
#SBATCH --nodes=1          #Numero de Nós
#SBATCH --ntasks-per-node=1 #Numero de tarefas por Nó
#SBATCH --ntasks=1          #Numero total de tarefas MPI
#SBATCH --cpus-per-task=24   #Numero de threads
#SBATCH -p cpu_dev          #Fila (partition) a ser utilizada
#SBATCH -J semana-sdumont-openmp #Nome job
```

```
#####
#
# USAGE: sbatch run.sh NUM_THREADS EXECUTABLE
#
#####
```

```
module load gcc/8.3

NUM_THREAD=${1}
EXEC=${2}

export OMP_NUM_THREADS=${NUM_THREAD}

srun -N 1 -c ${NUM_THREAD} ${EXEC}
```

Instalação do material na maquina que você tem acesso:

```
$ sftp user@intelknl.lncc.br
```

```
$ put OpenMP_2021_exemplos.tar
```

```
$ exit
```

```
$ ssh user@intelknl..lncc.br
```

```
$tar xvf OpenMP_2021_exemplos.tar
```



```
1 #include <omp.h>
2
3 main () {
4
5     int var1, var2, var3;
6
7     Serial code
8     .
9     .
10    .
11
12     Beginning of parallel region. Fork a team of threads.
13     Specify variable scoping
14
15     #pragma omp parallel private(var1, var2) shared(var3)
16     {
17
18         Parallel region executed by all threads
19         .
20         Other OpenMP directives
21         .
22         Run-time Library calls
23         .
24         All threads join master thread and disband
25
26     }
27
28     Resume serial code
29     .
30     .
31     .
32
33 }
```

- When a thread reaches a PARALLEL directive, it creates a team of threads and becomes the master of the team. The master is a member of that team and has thread number 0 within that team.
- Starting from the beginning of this parallel region, the code is duplicated and all threads will execute that code.
- There is an implied barrier at the end of a parallel region. Only the master thread continues execution past this point.
- If any thread terminates within a parallel region, all threads in the team will terminate, and the work done up until that point is undefined.

I - CONSTRUTOR PARALELO

```
#pragma omp parallel
```

- Informa ao compilador a existência de uma região que deve ser executada em paralelo.

```
#pragma omp parallel
{
    for (i = 0; i < n; i++)
        c[i] = a[i]+b[i];
}
```

C/C++

```
#pragma omp parallel [clause ...] newline
    if (scalar_expression)
    private (list)
    shared (list)
    default (shared | none)
    firstprivate (list)
    reduction (operator: list)
    copyin (list)
    num_threads (integer-expression)

    structured_block
```

Primeiro Exemplo

```
$ cd exemplo_1
```

```
$ gcc -fopenmp exemplo_01.c -o exemplo_01
```

```
$ ./exemplo_01
```

OBS:

O número de threads a serem gerados corresponde ao número de núcleos lógicos do processador . Você pode descobrir o número de núcleos lógicos através do comando:

```
$ cat /proc/cpuinfo
```

O número de threads a serem gerados pode ser alterado através coomando;

```
$ export OMP_NUM_THREADS=2
```

Execução do exemplo_01 no SDumont

- OpenMP “script” para compilar
 - `./compilar.sh`
- OpenMP “script” para submeter na fila de execução:
 - `$sbatch sub.sh 10 exemplo_01`
 - `$squeue -u $USER` (para visualizar o job na fila de execução)
 - `$scancel JOB_ID` (para cancelar o job da fila de execução, caso seja necessário)
 - `$cat slurm-JOB_ID.out`

Funções de Ambiente de Execução biblioteca omp.h

```
void omp_set_num_threads(int num)
```

- Define o número de *threads* padrão a serem usadas na região paralela.

```
int omp_get_num_threads()
```

- Retorna o número de *threads* ativas na região paralela onde ela foi chamada.

```
int omp_get_max_threads()
```

- Retorna o número de *threads* disponíveis para executar a região paralela.

```
int omp_get_thread_num()
```

- Retorna o identificador da thread relativo ao grupo ao qual ela pertence.

```
int omp_get_num_procs()
```

- Retorna o número de processadores disponíveis no momento da chamada da função.

exemplo_01_num_threads.c

```
#include <stdio.h>
#include <omp.h>

int main()
{
    #pragma omp parallel num_threads(10)
    {
        printf("omp_get_thread_num = %d\n", omp_get_thread_num());
    #pragma omp single
        printf("omp_get_num_procs = %d\n", omp_get_num_procs());
    }
}
```

exemplo_01_omp_set_num_threads.c

```
void omp_set_num_threads(int num)
```

- Define o número de *threads* padrão a serem usadas na região paralela.

```
int omp_get_num_threads()
```

- Retorna o número de *threads* ativas na região paralela onde ela foi chamada.

```
int omp_get_max_threads()
```

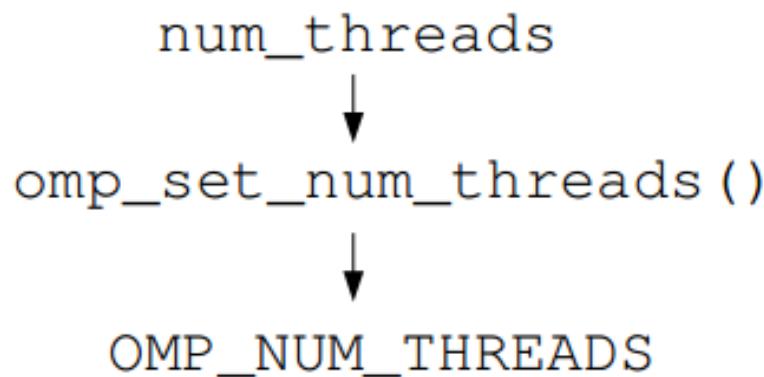
- Retorna o número de *threads* disponíveis para executar a região paralela.

Alteração do padrão de execução: Cláusula, Função da OMP.h , Variável de ambiente ,

- Define o número de *threads* que irá executar as regiões paralelas.

```
BASH : export OMP_NUM_THREADS = 8  
CSH  : setenv OMP_NUM_THREADS 8
```

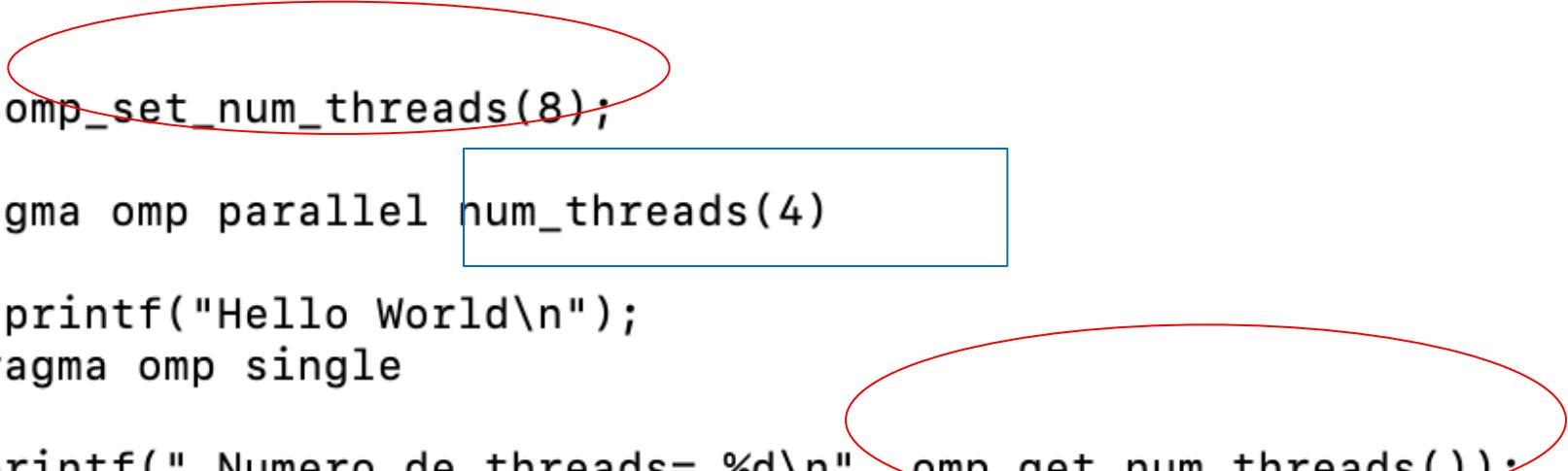
- Diagrama de Precedência:



exemplo_01_set_num_threads.c

```
#include <stdio.h>
#include <omp.h>

int main()
{
    omp_set_num_threads(8);
    #pragma omp parallel num_threads(4)
    {
        printf("Hello World\n");
    #pragma omp single
    {
        printf(" Numero de threads= %d\n", omp_get_num_threads());
        printf(" Numero max de threads= %d\n", omp_get_max_threads());
    }
}
}
```



II- CONSTRUTORES DE TRABALHO

- São responsáveis pela distribuição de trabalho entre as *threads* e determinam como o trabalho será dividido.

```
#pragma omp for
```

```
#pragma omp sections
```

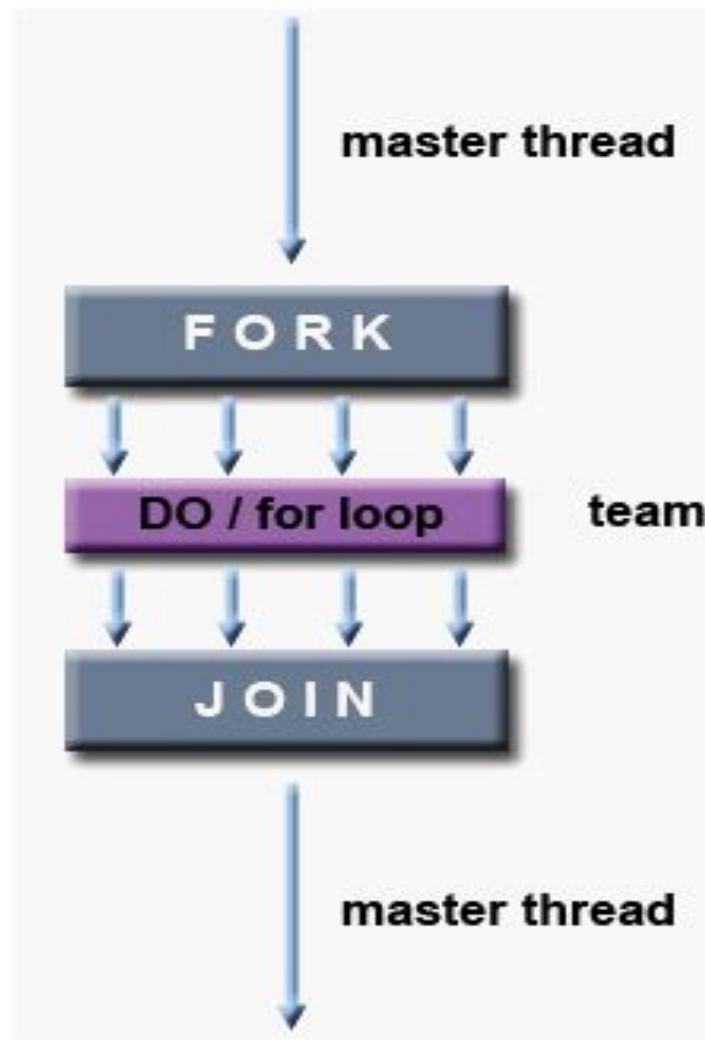
```
#pragma omp single
```

II.1- CONSTRUTOR DE TRABALHO FOR

```
#pragma omp for [cláusula, ...]
```

- Esse construtor é responsável pela divisão das iterações do laço a serem realizadas entre as *threads*.
- O número de iterações do laço deve ser previamente conhecido.

DO / for - shares iterations of a loop across the team. Represents a type of "data parallelism".



exemplo_04_1.c

```
#include <stdio.h>
#include <omp.h>

int main (int argc, char * argv[])
{
    int n,id,i;

    n=10;

    #pragma omp parallel
    {
        #pragma omp for
        for (i = 0; i < n; i++)
        {
            id= omp_get_thread_num();
            printf("Thread = %d, i= %d\n",id,i);
        }
    }
}
```

Observem que a thread 0 sempre pega as primeiras iterações
(quando o numero de threads é menor que o número de iterações) :
Teste o exemplo_04_1.c para. \$export OMP_NUM_THREADS=4

```
[professor@sdumont12 exemplo_4]$ ./exemplo_04_1
Thread = 0, i= 0
Thread = 0, i= 1
Thread = 1, i= 4
Thread = 1, i= 5
Thread = 1, i= 6
Thread = 2, i= 7
Thread = 2, i= 8
Thread = 2, i= 9
Thread = 0, i= 2
Thread = 0, i= 3
[professor@sdumont12 exemplo_4]$ ./exemplo_04_1
Thread = 0, i= 0
Thread = 0, i= 1
Thread = 0, i= 2
Thread = 0, i= 3
Thread = 1, i= 4
Thread = 1, i= 5
Thread = 1, i= 6
Thread = 2, i= 7
Thread = 2, i= 8
Thread = 2, i= 9
[professor@sdumont12 exemplo_4]$ ./exemplo_04_1
Thread = 0, i= 0
Thread = 0, i= 1
Thread = 1, i= 4
Thread = 1, i= 5
Thread = 1, i= 6
Thread = 2, i= 7
Thread = 2, i= 8
Thread = 2, i= 9
```

As cláusulas servem para alterar a execução padrão:

```
#pragma omp for [clause ...] newline
    schedule (type [,chunk])
    ordered
    private (list)
    firstprivate (list)
    lastprivate (list)
    shared (list)
    reduction (operator: list)
    collapse (n)
    nowait
```

for_loop

Cláusula IF : exemplo_03.c (n=10)

if (expressão lógica)

- Se a expressão lógica for verdadeira a região paralela será executada por mais de uma *thread*.

EXEMPLO:

```
#pragma omp parallel [if (n > 100*var)]  
{  
    #pragma omp for  
    for (i = 0; i < n; i++)  
    {  
        c[i] = a[i]*b[i];  
    }  
}
```



if(n>100*var)
 região paralela **ativa**
else
 região paralela **inativa**

Cláusula schedule do construtor de trabalho FOR

```
schedule(tipo, [tamanho])
```

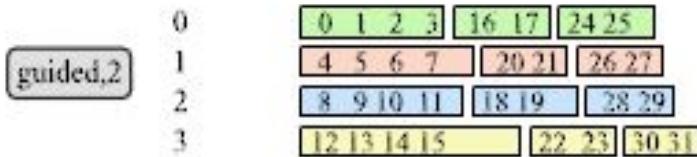
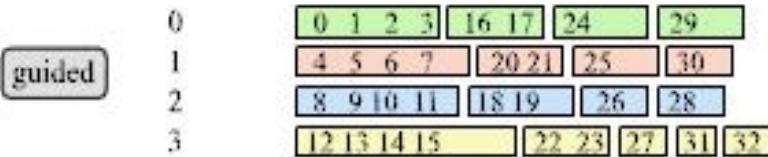
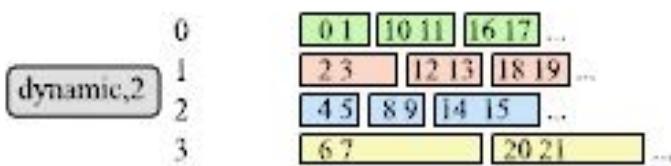
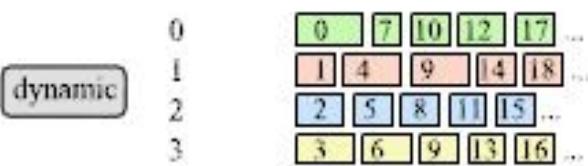
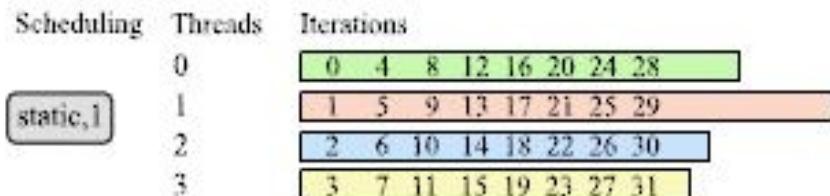
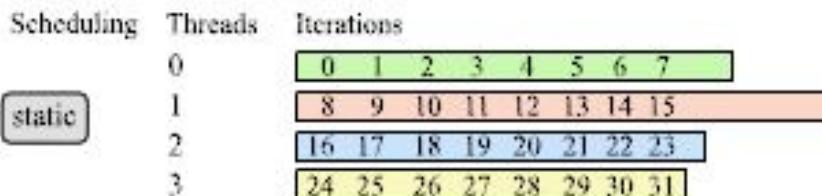
- É utilizada apenas no construtor for e controla a forma com as iterações são distribuídas entre as *threads*.
- `schedule` pode ser do tipo:
 - static
 - dynamic
 - guided
 - runtime

static é eficiente com carga do SO ou dados homogêneos

dynamic é eficiente com carga do SO ou dados heterogêneos

guided mistura dos dois

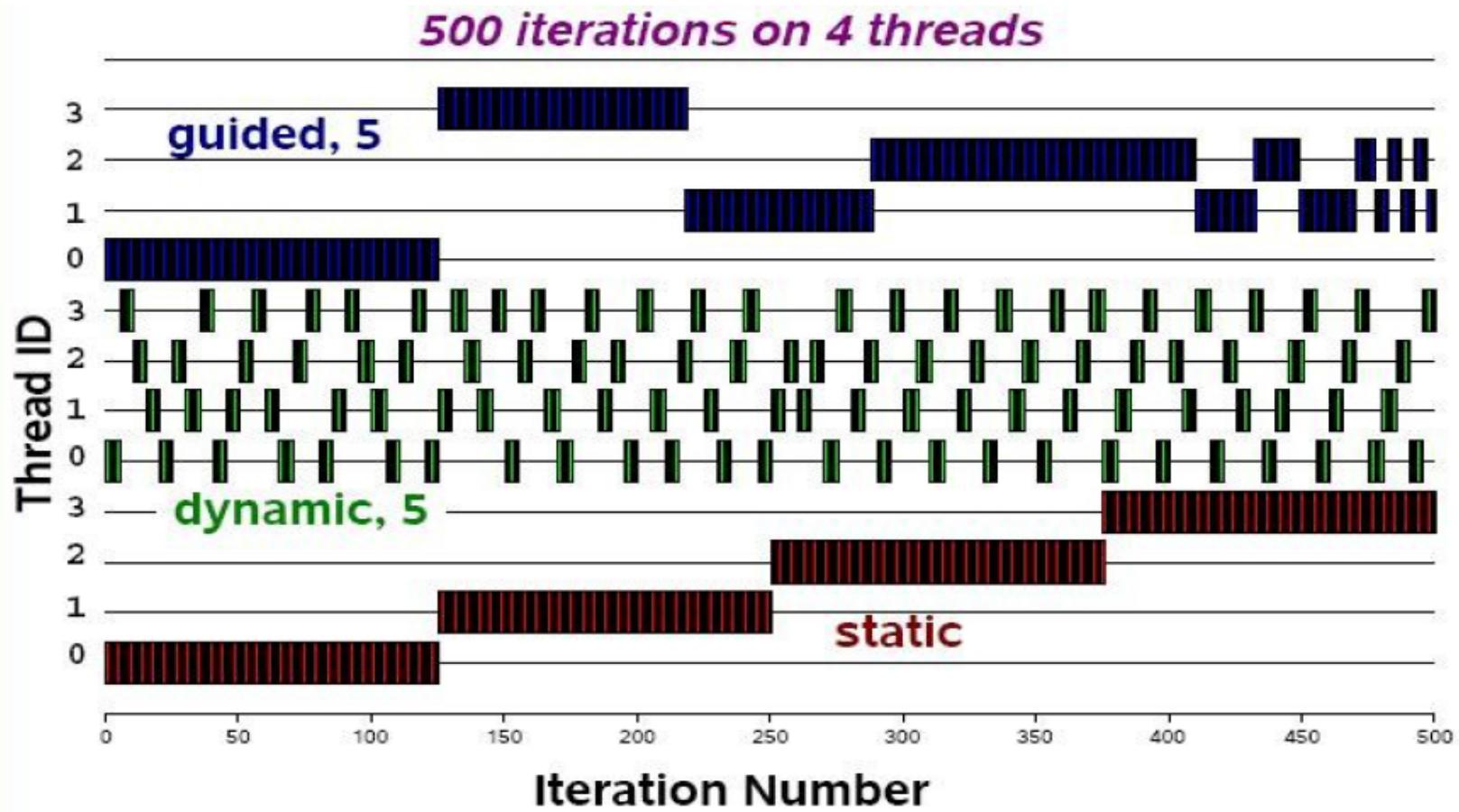
LOOP SCHEDULING MODES IN OPENMP



Time

Time

Escalonamento das cláusulas do divisor de trabalho FOR



Cláusula schedule

`$./exemplo_04_1` (executa configuração “default”)

`$./exemplo_04_2` (cláusula schedule dynamic, chunk=5)

`$./exemplo_04_4` (cláusula schedule dynamic, chunk=1)

Teste o exemplo_04_3 para diversos valores de n, schedules e chunks, conforme abaixo:

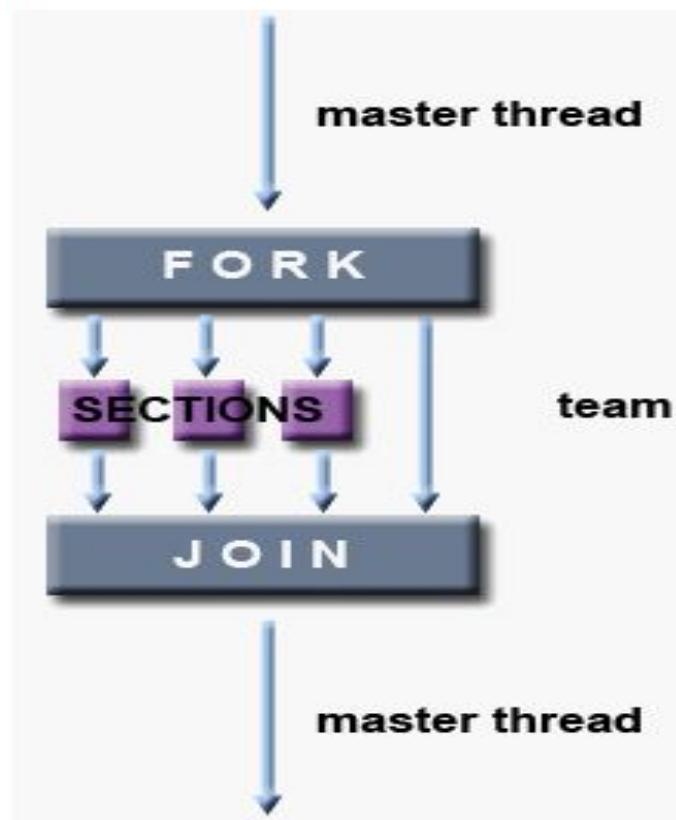
`$time ./exemplo_04_3` (runtime: usa variável de ambiente)

`$export OMP_NUM_THREADS=x`

`$export OMP_SCHEDULE="mode, chunk"`

II.2- CONSTRUTOR DE TRABALHO SECTIONS

SECTIONS - breaks work into separate, discrete sections. Each section is executed by a thread. Can be used to implement a type of "functional parallelism".



#pragma omp section

```
#pragma omp parallel
{
    #pragma omp sections
    {
        #pragma omp section
        soma_vetores(a,b,c,n);
        #pragma omp section
        subtrai_vetores(a,b,d,n);
    }
}
```

- O construtor **section** define o bloco que será executado por uma *thread* dentro do construtor **sections**

Cláusulas para alterar a execução padrão

```
#pragma omp sections [clause ...] newline
    private (list)
    firstprivate (list)
    lastprivate (list)
    reduction (operator: list)
    nowait

{

#pragma omp section    newline
    structured_block

#pragma omp section    newline
    structured_block

}
```

-exemplo_06.c :

(3 threads trabalham em paralelo)

-exemplo_06_1.c :

divisor de trabalho “single”, testar para:
export OMP_NUM_THREADS=3

(Apenas duas threads recebem trabalho no
“sections” e o terceiro thread tem que
esperar o trabalho das outras duas acabar
observe que neste caso um mesmo thread
pode executar uma das duas na sections e
depois no single). obs: Existe uma barreira
implícita no final de todo “construtor de
trabalho”

Exemplos para Construtor de Trabalho Sections

exemplo_05_corrida.c: (as variáveis "id" e "i" são globais e são alteradas simultaneamente pelas threads .)

Observem que cada vez que se executa aparece uma saída diferente. Isto se deve porque a variável id e i estão sendo alteradas simultaneamente pelas threads da section vetor B e section vetor A. Para corrigir, tem que transformar estas variáveis em "private".

exemplo_05.c (cada thread executa uma section e escreve em variáveis privativas)

exemplo_05_1.c (as variáveis dos enlaces são diferentes)

Observação:

- Apenas os índices do enlace dos construtores de trabalho “for” são “private”
- Os índices de enlaces dos demais construtores de trabalho necessitam ser declarados como “private” .
- Variáveis “private” precisam ser inicializadas dentro da região paralela.

VARIÁVEIS “PRIVATE”

Só existem dentro da região paralela e precisam ser inicializadas dentro da região paralela

- **exemplo_04_1_private.c** (id e m são variáveis private e precisam ser inicializadas dentro da região paralela)

```
3  #include <stdio.h>
4  #include <omp.h>
5
6  int main (int argc, char * argv[])
7  {
8      int n,m,id,i;
9
10     n=10;
11     m=10;
12     #pragma omp parallel private(id,m)
13     {
14         #pragma omp for
15         for (i = 0; i < n; i++)
16         {
17             id= omp_get_thread_num();
18             printf("Thread = %d, i= %d , m=%d\n",id,i,m);
19         }
20     }
21 }
22 }
```

- **exemplo_04_firstprivate.c** (inicializa as variáveis privativas m e id com valor global que tinham antes da região paralela)

```
3  #include <stdio.h>
4  #include <omp.h>
5
6  int main (int argc, char * argv[])
7  {
8      int n,m,id,i;
9
10     n=10;
11     m=10;
12     #pragma omp parallel firstprivate(id,m)
13     {
14         #pragma omp for
15         for (i = 0; i < n; i++)
16         {
17             id= omp_get_thread_num();
18             printf("Thread = %d, i= %d , m=%d\n",id,i,m);
19         }
20     }
21 }
```

- **exemplo_04_lastprivate.c** (“i” armazena na variável global o último valor dentro da região paralela) **Observe que ao retirar o lastprivate da variável i ela volta ter o valor da variável global i.**

```
2
3     #include <stdio.h>
4     #include <omp.h>
5
6     int main (int argc, char * argv[])
7     {
8         int n,m,id,i;
9
10        n=10;
11        #pragma omp parallel
12        {
13            #pragma omp for  lastprivate(i)
14            for (i = 0; i < n; i++)
15            {
16                id= omp_get_thread_num();
17            }
18        }
19        printf("Thread = %d, i= %d\n",id,i);
20    }
21
```

- exemplo_04_variavel_local.c (as variáveis id e i são inicializadas dentro da região paralela e precisam ser “inicializadas”.)

```
#include <stdio.h>
#include <omp.h>

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

    n=10;

    #pragma omp parallel
    {
        int id,i;
        printf("valor ao ser inicializada:Thread = %d, i= %d\n",id,i);
        #pragma omp for
        for (i = 0; i < n; i++)
        {
            id= omp_get_thread_num();
            printf(" valor dentro do divisor de trabalho for: Thread = %d, i= %d\n"
id,i);
        }
    }
}
```

Sincronizador barrier:

```
#pragma omp barrier
```

- É utilizado para sincronizar todas as *threads* em determinado ponto do código.
- Em alguns construtores existe uma barreira implícita: “Na saída” `parallel`, `for`, `sections`, `single`.



nowait

- Essa cláusula faz com que as *threads* ignorem as barreiras implícitas.

➢ **Que diretivas podem utilizar essa cláusula ?**

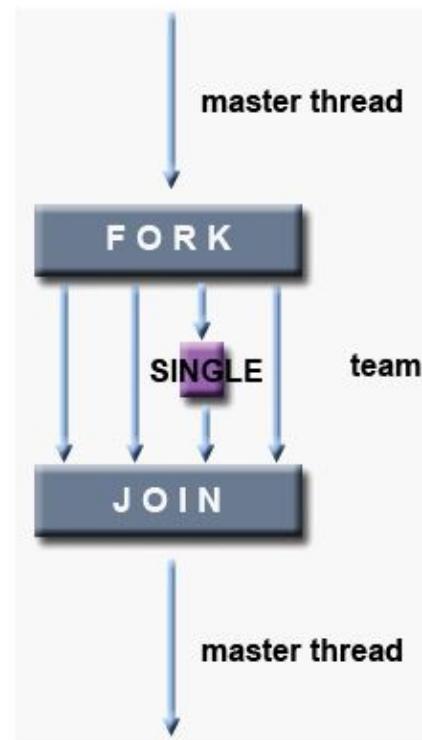
```
#pragma omp single  
#pragma omp for  
#pragma omp sections  
#pragma omp parallel for  
#pragma omp parallel sections
```

exemplo_5

- Exemplo_05.c: cada section é executada com uma thread distinta
- exemplo_05_nowait.c (a section single pode ser executada por um thread que havia executado outra section e já terminou.)
- exemplo_05_sem_nowait.c : a thread single executará após que todas as threads da section anterior tiverem terminado

II.3- CONSTRUTOR DE TRABALHO SINGLE

SINGLE - serializes a section of code



exemplo_05_sem_nowait.c

exemplo_01_set_num_threads.c

exemplo_01_num_threads.c

#pragma omp single [cláusula, ...]

- Esse construtor indica que o bloco sintático localizado logo abaixo do mesmo deve ser executado por apenas uma *thread*.
- Não necessariamente a *thread master*

```
#pragma omp parallel
{
    #pragma omp single
        a = function(t);
    ...
}
```

CLÁUSULAS

- private
- firstprivate
- copyprivate
- nowait

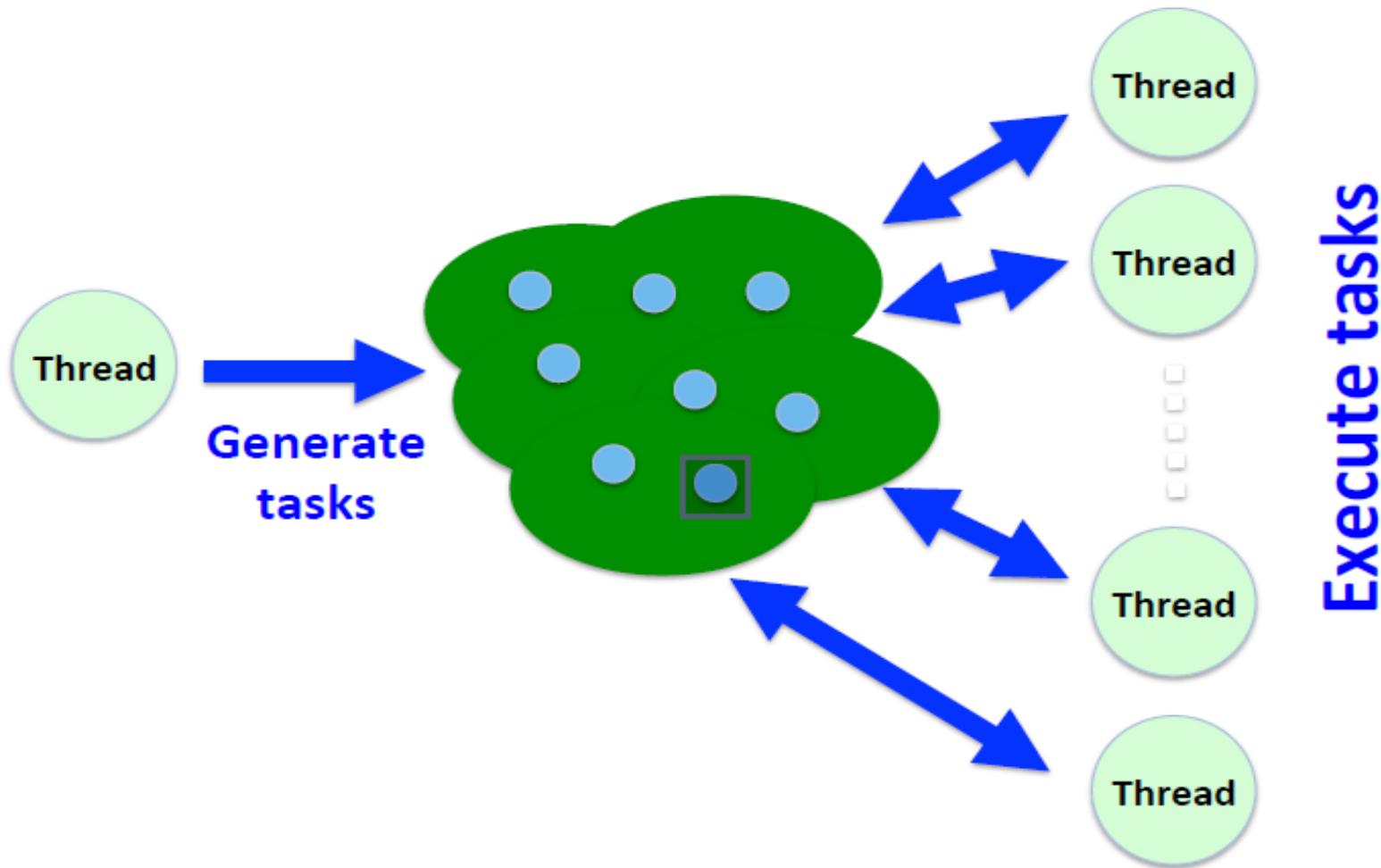
#pragma omp master

- Define um bloco de código que será executado apenas pela *thread mestre*.
- Não possui barreira implícita, na entrada e na saída.

exemplo_5_nowait_master.c

(só será executada pela “thread 0”, inclusive quando as outras estiverem livres

II.4- CONSTRUTOR DE TRABALHO TASK



Diretivas do Construtor de Trabalho

TASK

```
#pragma omp task [clause ...] newline
    if (scalar expression)
    final (scalar expression)
    untied
    default (shared | none)
    mergeable
    private (list)
    firstprivate (list)
    shared (list)
```

C/C++

structured_block

- Quando uma thread encontra um construtor de task, uma nova task é gerada.
- A execução da task cabe ao sistema de runtime
- O término de uma task pode ser forçado através de um “task”

Diretivas do Construtor de Trabalho

TASK

```
#pragma omp task  
!$omp task
```

Defines a task

#pragma omp barrier

```
#pragma omp barrier  
!$omp barrier
```

#pragma omp taskwait

```
#pragma omp taskwait  
!$omp taskwait
```

exemplos/racecar/racecar.c

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

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

    printf("A ");
    printf("race ");
    printf("car ");

    printf("\n");
    return(0);
}
```

```
$ cc -fast hello.c
$ ./a.out
A race car
$
```

What will this program print ?

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

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

    #pragma omp parallel
    {
        printf("A ");
        printf("race ");
        printf("car ");

    } // End of parallel region

    printf("\n");
    return(0);
}
```

***What will this program print
using 2 threads ?***

```
$ cc -xopenmp -fast hello.c
$ export OMP_NUM_THREADS=2
$ ./a.out
A race car A race car
```

*Note that this program could (for example) also print
“A A race race car car” or
“A race A car race car”, or
“A race A race car car”, or
.....*

But I have not observed this (yet)

```
#include <stdlib.h>
#include <stdio.h>
int main(int argc,
```

***What will this program print
using 2 threads ?***

```
#pragma omp parallel
{
    #pragma omp single
    {
        printf("A ");
        printf("race ");
        printf("car ");
    }
} // End of parallel region

printf("\n");
return(0);
}
```



```
$ cc -xopenmp -fast hello.c
$ export OMP_NUM_THREADS=2
$ ./a.out
A race car
```

```
int main(int argc, char *argv[]) {  
    #pragma omp parallel  
    {  
        #pragma omp single  
        {  
            printf("A ");  
            #pragma omp task  
            {printf("race ");}  
            #pragma omp task  
            {printf("car ");}  
        }  
    } // End of parallel region  
  
    printf("\n");  
    return(0);  
}
```

What will this program print using 2 threads ?

```
$ cc -xopenmp -fast hello.c  
$ export OMP_NUM_THREADS=2  
$ ./a.out  
A race car  
$ ./a.out  
A race car  
$ ./a.out  
A car race  
$
```

```
int main(int argc, char *argv[]) {  
    #pragma omp parallel  
    {  
        #pragma omp single  
        {  
            printf("A ");  
            #pragma omp task  
            {printf("race ");}  
            #pragma omp task  
            {printf("car ");}  
            printf("is fun to watch ");  
        }  
    } // End of parallel region  
  
    printf("\n");  
    return(0);  
}
```

What will this program print using 2 threads ?

```
$ cc -xopenmp -fast hello.c  
$ export OMP_NUM_THREADS=2  
$ ./a.out  
  
A is fun to watch race car  
$ ./a.out  
  
A is fun to watch race car  
$ ./a.out  
  
A is fun to watch car race  
$
```

```
int main(int argc, char *argv[]) {  
    #pragma omp parallel  
    {  
        #pragma omp single  
        {  
            printf("A ");  
            #pragma omp task  
            {printf("race ");}  
            #pragma omp task  
            {printf("car ");}  
            printf("is fun to watch ");  
        }  
    } // End of parallel region  
  
    printf("\n");  
    return(0);  
}
```

What will this program print using 2 threads ?

```
$ cc -xopenmp -fast hello.c  
$ export OMP_NUM_THREADS=2  
$ ./a.out  
  
A is fun to watch race car  
$ ./a.out  
  
A is fun to watch race car  
$ ./a.out  
  
A is fun to watch car race  
$
```

```

int main(int argc, char *argv[])
{
    #pragma omp parallel
    {
        #pragma omp single
        {
            printf("A ");
            #pragma omp task
            {printf("car ");}
            #pragma omp task
            {printf("race ");}
            #pragma omp taskwait
            printf("is fun to watch ");
        }
    } // End of parallel region

    printf("\n");
}

```

What will this program print using 2 threads ?



```

$ cc -xopenmp -fast hello.c
$ export OMP_NUM_THREADS=2
$ ./a.out
$ 
A car race is fun to watch
$ ./a.out
A car race is fun to watch
$ ./a.out
A race car is fun to watch
$ 

```

```
my_pointer = listhead;  
  
#pragma omp parallel  
{  
    #pragma omp single  
    {  
        while(my_pointer)  
            #pragma omp task firstprivate(my_pointer)  
            {  
                (void) do_independent_work (my_pointer);  
            }  
        my_pointer = my_pointer->next ;  
    }  
} // End of single  
} // End of parallel region
```

*OpenMP Task is specified here
(executed in parallel)*

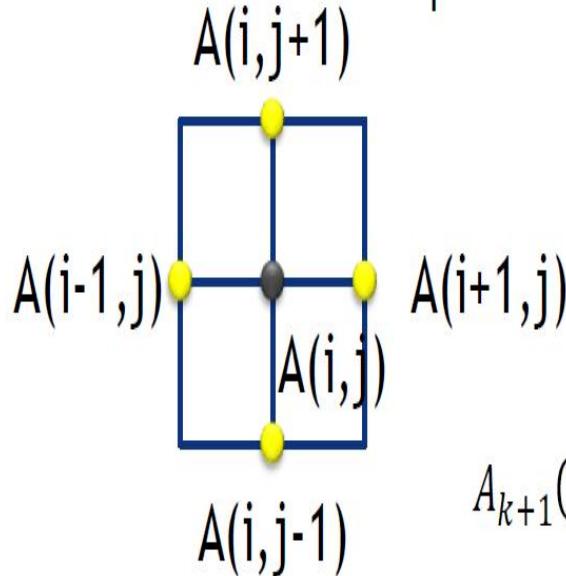


Example: Jacobi Iteration

Iteratively converges to correct value (e.g. Temperature), by computing new values at each point from the average of neighboring points.

Common, useful algorithm

Example: Solve Laplace equation in 2D: $\nabla^2 f(x, y) = 0$



$$A_{k+1}(i,j) = \frac{A_k(i-1,j) + A_k(i+1,j) + A_k(i,j-1) + A_k(i,j+1)}{4}$$

Exemplo: jacobi/laplace2d.c

Apresenta uma paralelização com o construtor "FOR"

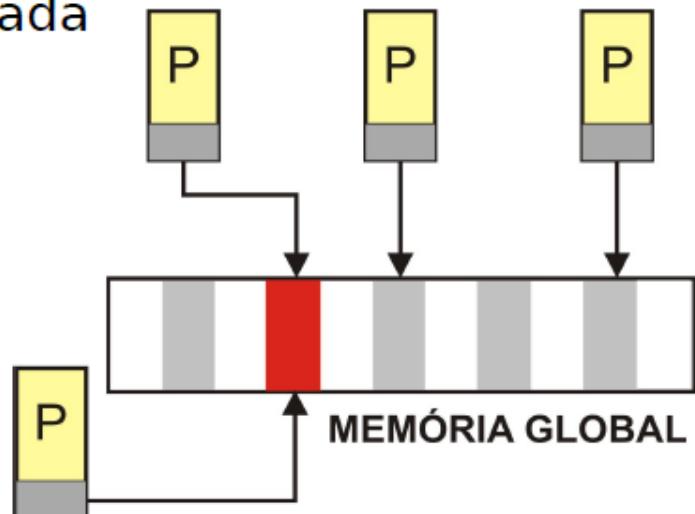
Exemplo laplace2d_task.c

Apresenta uma paralelização com o construtor "task"

Podemos observar que o construtor "task" introduz mais gastos de execução e deve ser implementado para tasks com "granularidade grossa"

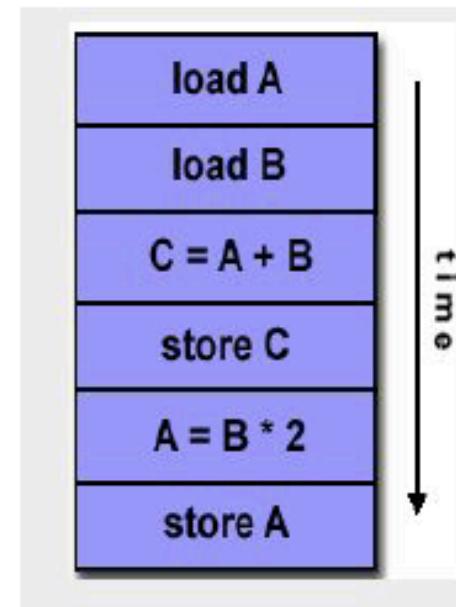
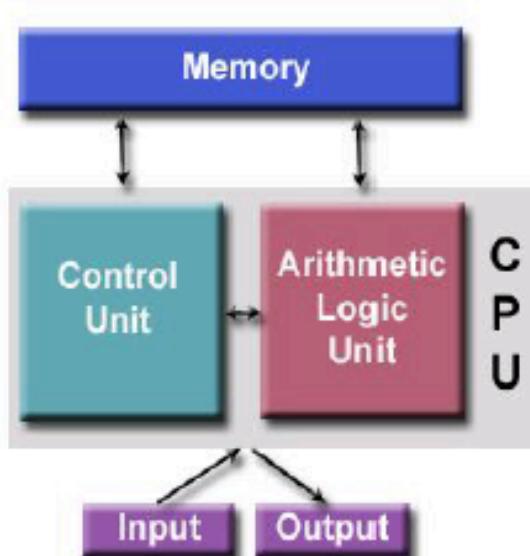
III- CONDIÇÃO DE CORRIDA

- › Acontece quando duas ou mais *threads* tentam atualizar, ao mesmo tempo, uma mesma variável ou quando uma *thread* atualiza uma variável e outra *thread* acessa o valor dessa variável ao mesmo tempo.
- › Dessa forma, as diretivas de sincronização garantem que o acesso ou atualização de uma determinada variável compartilhada aconteça no momento certo.

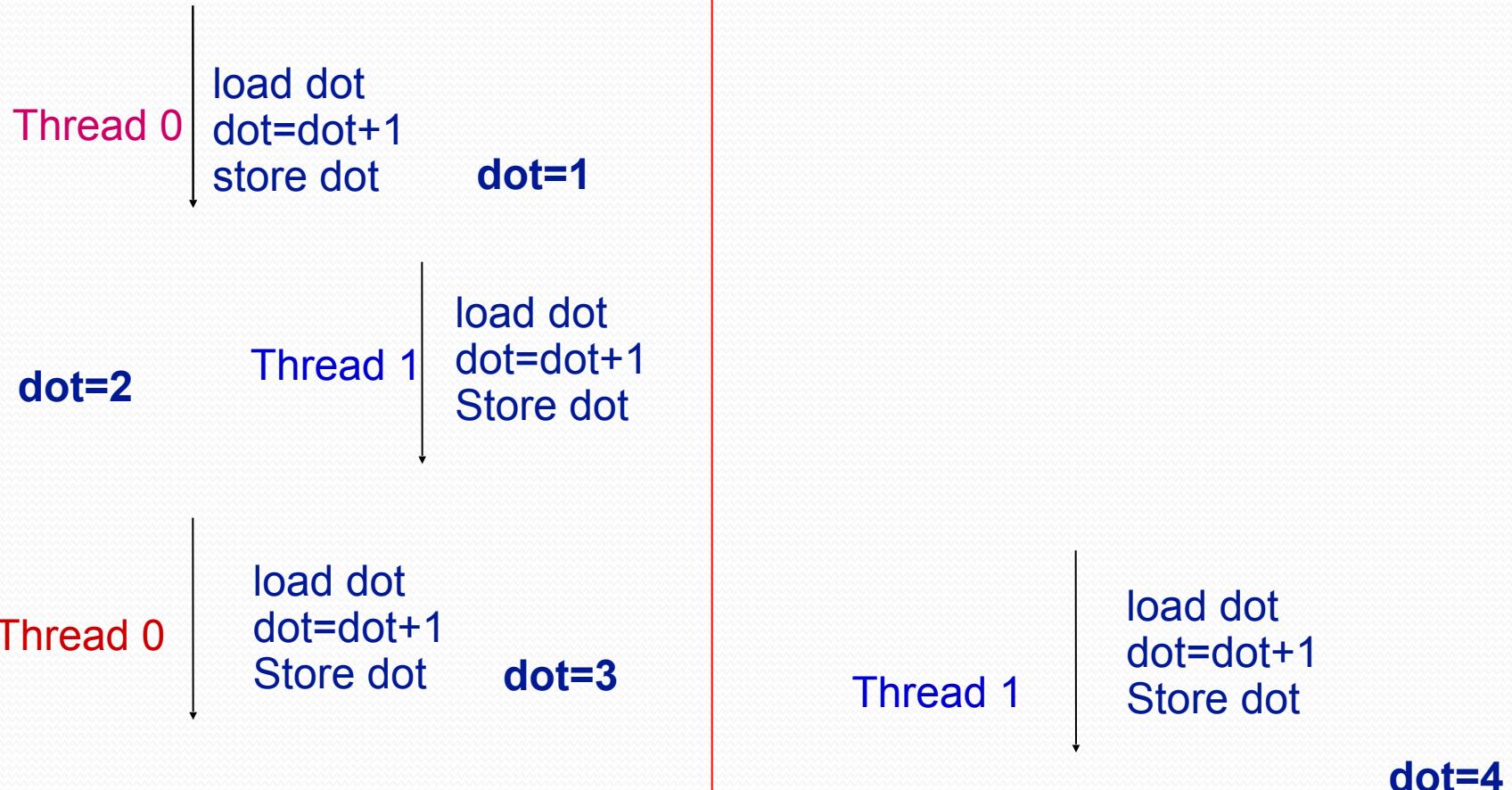


Arquitetura Von Neuman 1945 (Processador Sequencial)

Operações: $C = A + B$
 $A = B * 2$

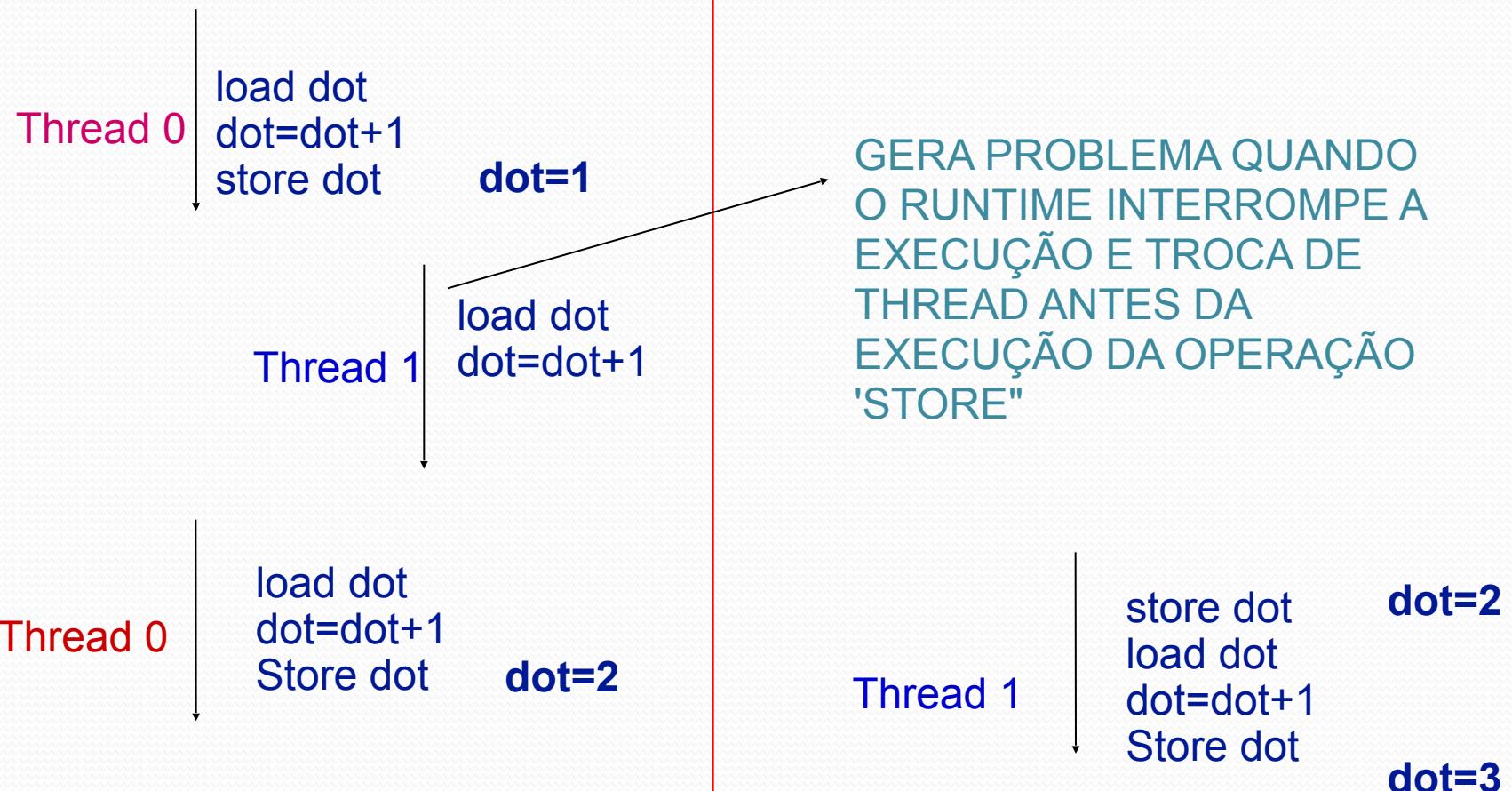


Operação “dot=dot+1”



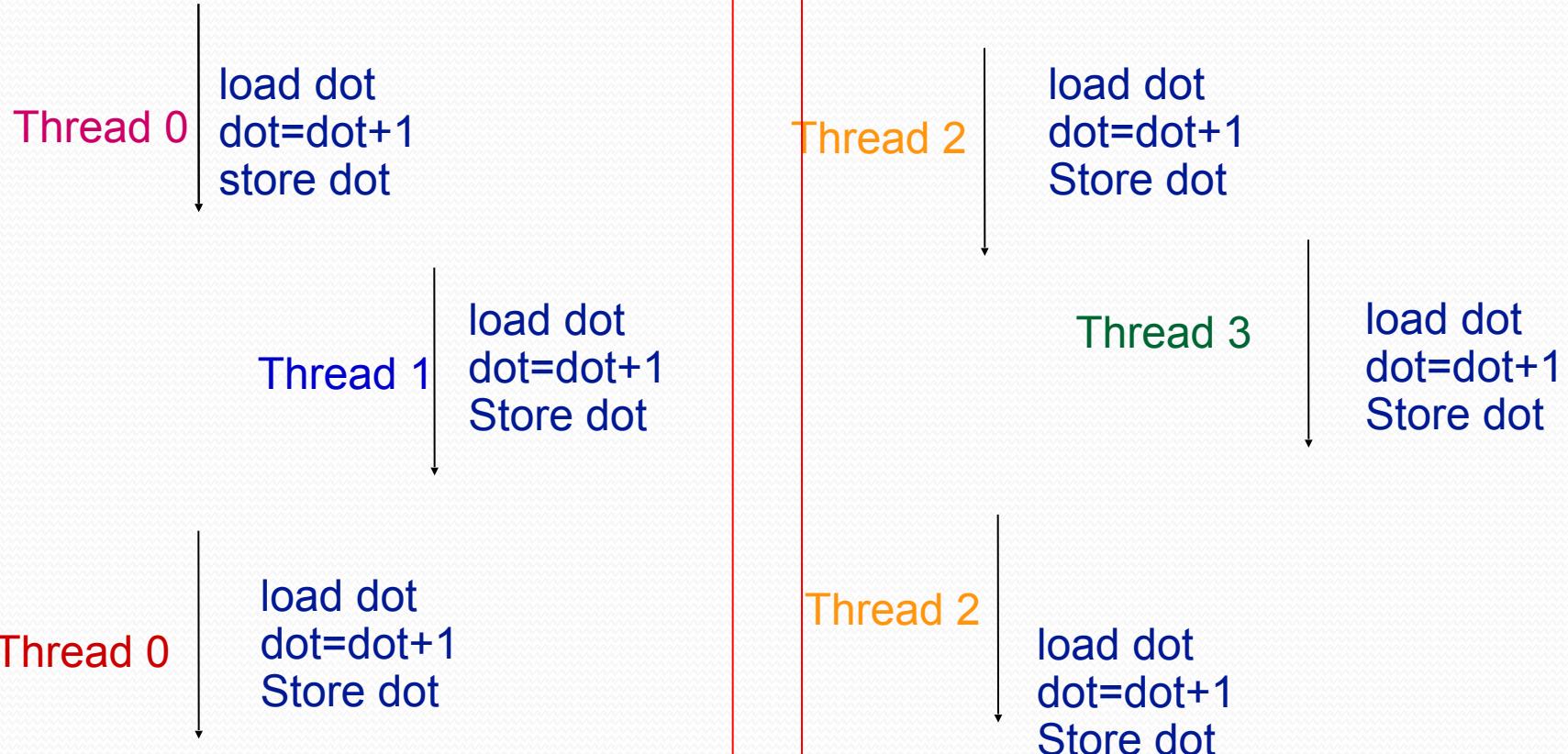
Para `dot=0` no inicio no final de 4 execuções temos `dot=4`

PROBLEMA 1: INTERRUPÇÃO DO RUNTIME ANTES DO THREAD TERMINAR A OPERAÇÃO DE ESCRITA



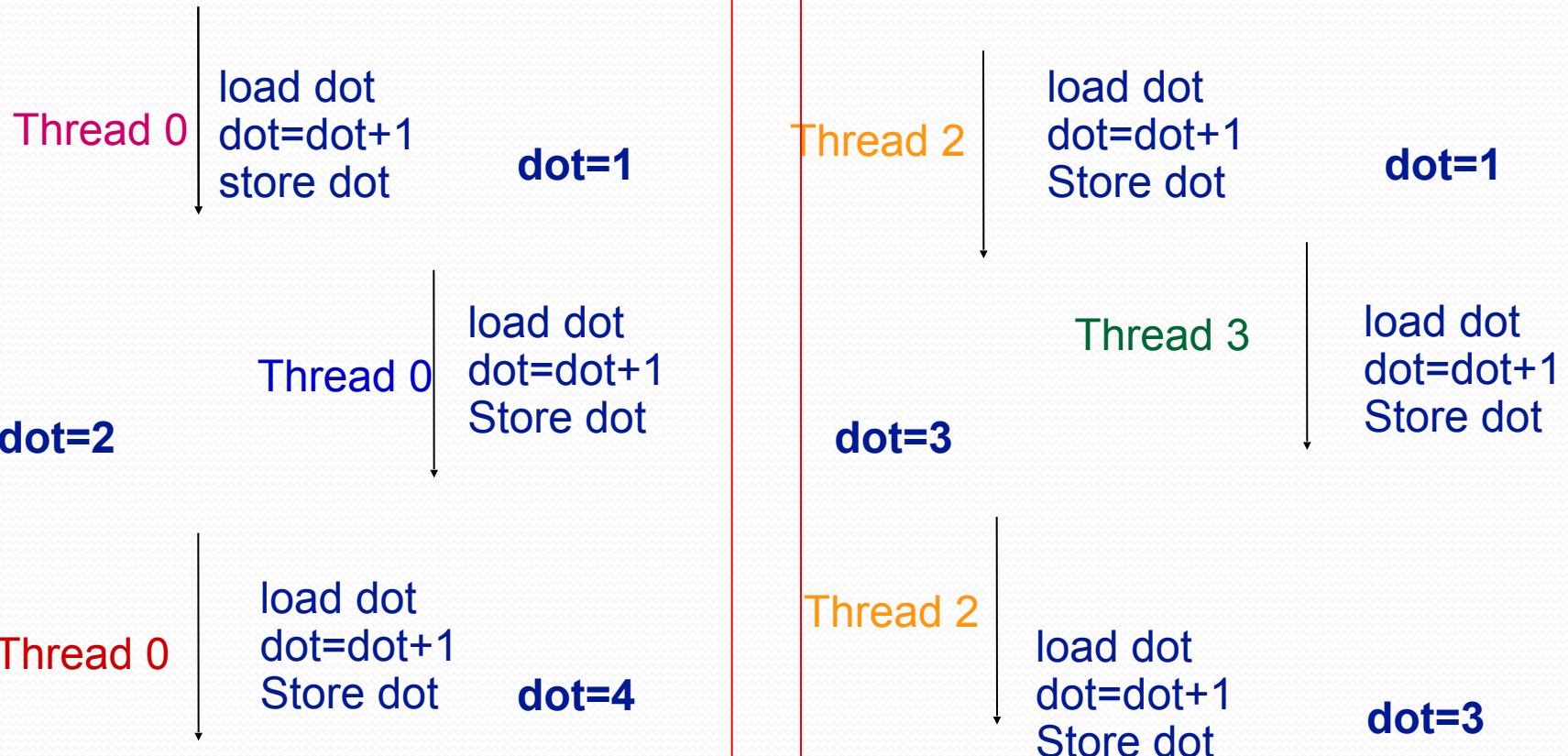
Para `dot=0` no inicio no final de 4 execuções temos `dot=3`

OPERAÇÃO DE ESCRITA POR DIVERSAS THREADS EM UMA MESMA VARIÁVEL GLOBAL



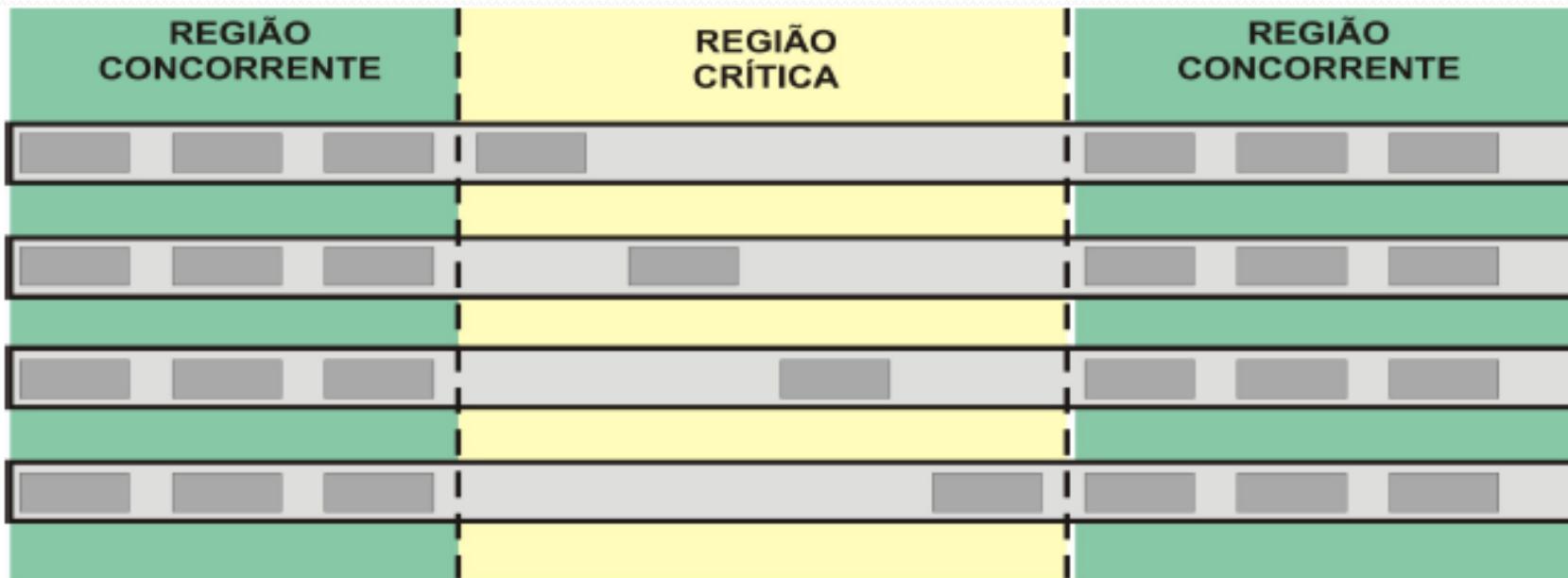
Para $\text{dot}=0$ no inicio no final de 6 execuções temos $\text{dot}=6$

PROBLEMA 2: ESCRITADAS DESORDENADAS NA VARIÁVEL GLOBAL



Para `dot=0` no inicio no final de 6 execuções temos `dot=3`

Operações em regiões com
“Condição de corrida” (escrita
simultânea em uma mesma variável
global por diversas threads) devem
ser executadas serialmente e todas
as operações devem ser
“finalizadas” dentro da região critica



Diretório exemplo_8

- \$./exemplo_08_corrida.c 100
- Observe que a cada execução do código acima o resultado de saída é diferente, ou seja, o código gera um resultado “inconsistente”, devido a condição de corrida gerada pela escrita simultânea de diversas thread na variável global dot”.

• exemplo_08_corrida.c

```
int main (int argc, char * argv[])
{
    int *V1,*V2;
    int i,j,id;
    int num_threads;
    int par,num,dot;

    printf("\n\n=====\\n");
    printf("\tPRODUTO ESCALAR\\n");
    printf("=====\\n\\n");

//Dimensao dos Vetores

    num = atoi(argv[1]);

//Alocacao dinamica dos vetores
    V1=(int*)malloc(num*sizeof(int));
    V2=(int*)malloc(num*sizeof(int));

//Inicializacao dos vetores
    InicializaVetores(num,V1,V2);

//Inicilializao da variavel dot
    dot = 0;
//Calculo do produto escalar
    #pragma omp parallel
    {
        #pragma omp for
        for(i=0;i<num;i++)
            dot += V1[i]*V2[i];
    }
    printf(" \\n\\n End Of Execution ::: dot = %d\\n\\n\\n",dot);
return 0;
}
```

Diretiva critical

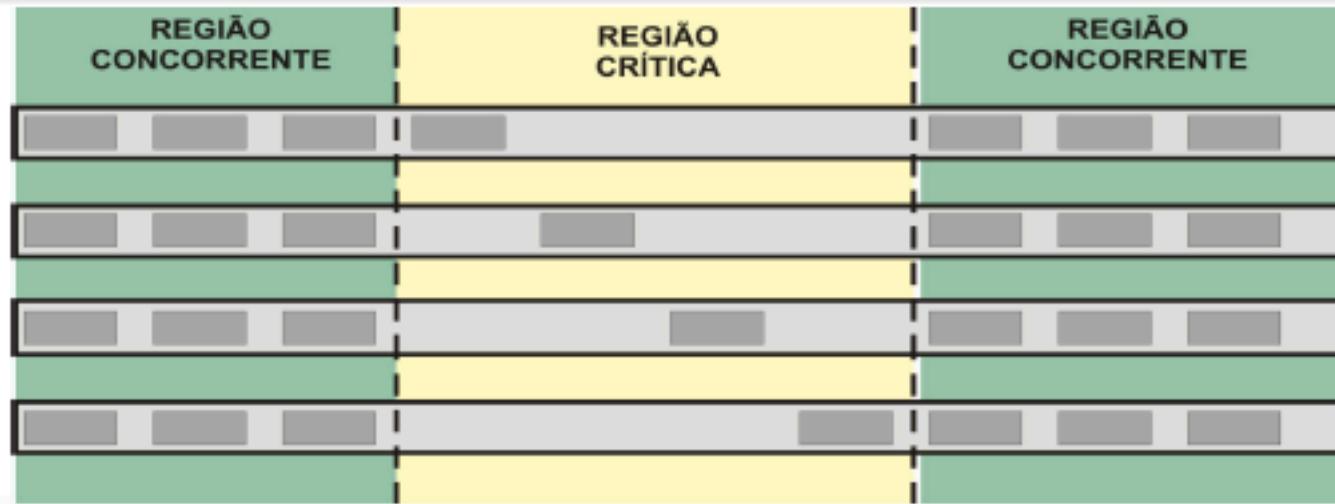
```
#pragma omp critical
```

- Restringe a execução de uma determinada tarefa a uma *thread* (*threads* do mesmo grupo) por vez.
- Atualização no sistema operacional

Diretiva critical

```
#pragma omp parallel
{
    int i, aux_dot, dot;
    #pragma omp for
    for (i = 0; i < n; i++){
        aux_dot += a[i]*b[i];
        #pragma omp critical
        dot += aux_dot;
    }
}
```

Garante que o acesso será feito por uma *thread* por vez



exemplo_08_sem_corrida_critical.c

```
//Alocacao dinamica dos vetores
    V1=(int*)malloc(num*sizeof(int));
    V2=(int*)malloc(num*sizeof(int));

//Inicializacao dos vetores
    InicializaVetores(num,V1,V2);

//Inicilializao da variavel dot
    dot = 0;
    int aux_dot=0;
//Calculo do produto escalar
    #pragma omp parallel
    {
        #pragma omp for
        for(i=0;i<num;i++)
            #pragma omp critical
            {
                dot += V1[i]*V2[i];
            }
    }
    printf(" \n\n End Of Execution :::: dot = %d\n\n",dot);
return 0;
```

Diretiva atomic

```
#pragma omp atomic
```

- Impedem que várias *threads* acessem essa variável ao mesmo tempo.
- Esse bloqueio é aplicado a todas as *threads* que executam o programa, não apenas as *threads* do mesmo grupo.
- Atualização a nível do processador

LIMITATIONS OF ATOMIC OPERATIONS

Read : operations in the form $v = x$

Write : operations in the form $x = v$

Update : operations in the form $x++, x--, --x, ++x, x \ binop= \ expr$
and $x = x \ binop \ expr$

Capture : operations in the form $v = x++, v = x-, v = -x, v = ++x,$
 $v = x \ binop \ expr$

- ▷ Here x and v are scalar variables
- ▷ $binop$ is one of $+, *, -, -, /, \&, ^, |, \ll, \gg$.
- ▷ No “trickery” is allowed for atomic operations:
 - no operator overload,
 - no non-scalar types,
 - no complex expressions.

Critical X atomic

ESCOPO

```
#pragma omp critical
{
    Instrução(1)
    ...
    Instrução(2)
}
```

```
#pragma omp atomic
Operação(1)
```

```
#pragma omp atomic
Operação(2)
```

```
#pragma omp critical
    x++
```

```
#pragma omp atomic
    x++
```

exemplo_08_sem_corrida_atomic.c

```
//Alocacao dinamica dos vetores
V1=(int*)malloc(num*sizeof(int));
V2=(int*)malloc(num*sizeof(int));

//Inicializacao dos vetores
    InicializaVetores(num,V1,V2);

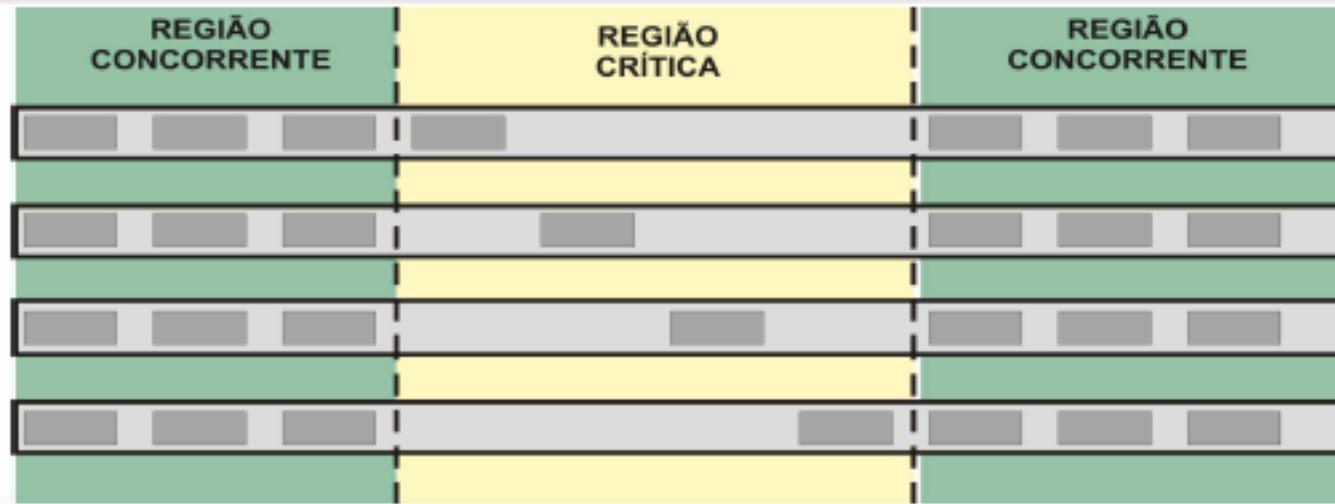
//Inicilializao da variavel dot
    dot = 0;
//Calculo do produto escalar
    #pragma omp parallel
    {
        #pragma omp for
        for(i=0;i<num;i++)
            #pragma omp atomic
            dot += V1[i]*V2[i];
    }
    printf(" \n\n End Of Execution :::: dot = %d\n\n\n",dot);
return 0;

}
```

Observem que podemos reduzir o número de operações na região critica utilizando variáveis privativas para as threads trabalharem na sua memória local e transferir a seguir os resultados parciais para a variável global utilizando uma região critica

```
#pragma omp parallel
{
    int i, aux_dot, dot;
    #pragma omp for
    for (i = 0; i < n; i++){
        aux_dot += a[i]*b[i];
        -----
        #pragma omp critical
        dot += aux_dot;
        -----
    }
}
```

Garante que o acesso será feito por uma *thread* por vez

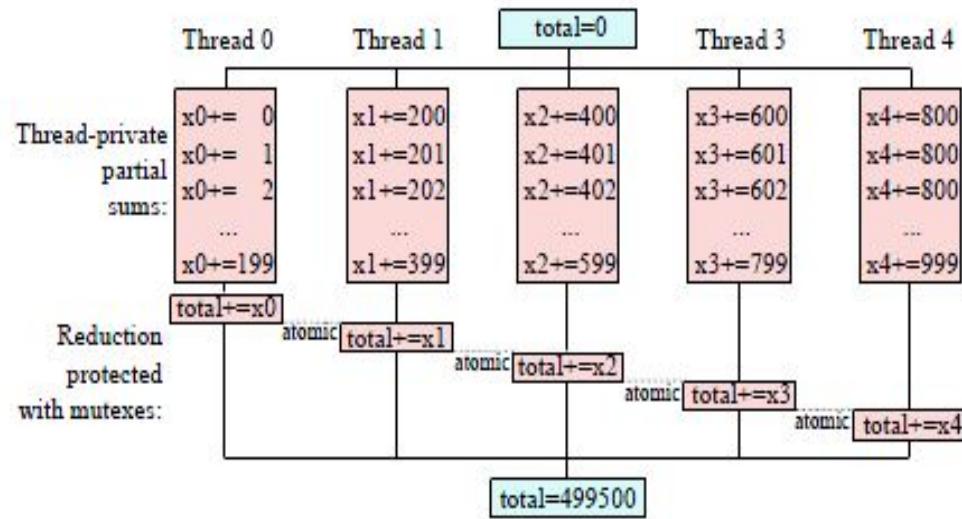


Para reduzir o número de operações dentro da região critica, cada thread deve trabalhar localmente com uma "variável auxiliar local" e a seguir, "passar" o resultado final para a variável global através de uma região critica:

AVOIDING RACES WITH THREAD-PRIVATE STORAGE

Correct and efficient code:

```
1 int total = 0;
2 #pragma omp parallel
3 {
4     int total_thr = 0;
5     #pragma omp for
6     for (int i=0; i<n; i++)
7         total_thr += i;
8
9     #pragma omp atomic
10    total += total_thr;
11 }
12 }
```



exemplo_08_critical_reduction.c

```
//Inicializao da variavel dot
    dot = 0;
//Calculo do produto escalar
#pragma omp parallel
{
    int aux_dot=0;
#pragma omp for
    for(i=0;i<num;i++)
        aux_dot += V1[i]*V2[i];

    #pragma omp critical
    {
        dot+=aux_dot;
    }
}
printf(" \n\n End Of Execution :::: dot = %d\n\n\n",dot);
return 0;
```

exemplo_08_atomic_reduction.c

```
//Inicializao da variavel dot
    dot = 0;
//Calculo do produto escalar
#pragma omp parallel
{
    int aux_dot=0;
#pragma omp for
    for(i=0;i<num;i++)
        aux_dot += V1[i]*V2[i];

    #pragma omp atomic
    dot+=aux_dot;
}
printf(" \n\n End Of Execution :::: dot = %d\n\n",dot);
return 0;
```

Cláusula Reduction do construtor de trabalho FOR

reduction (operador : lista de variáveis)

- Uma cópia de cada variável é criada para cada *thread*.
- Ao final da região paralela definida pelo construtor, a lista de variáveis original é atualizada com os valores da cópia privada de cada *thread* usando o operador especificado.

Operador	Valor Inicial
+	0
*	1
-	0
^	0

Operador	Valor Inicial
&	~0
	0
&&	1
	0

exemplo_08.c

```
//Alocacao dinamica dos vetores
V1=(int*)malloc(num*sizeof(int));
V2=(int*)malloc(num*sizeof(int));

//Inicializacao dos vetores
    InicializaVetores(num,V1,V2);

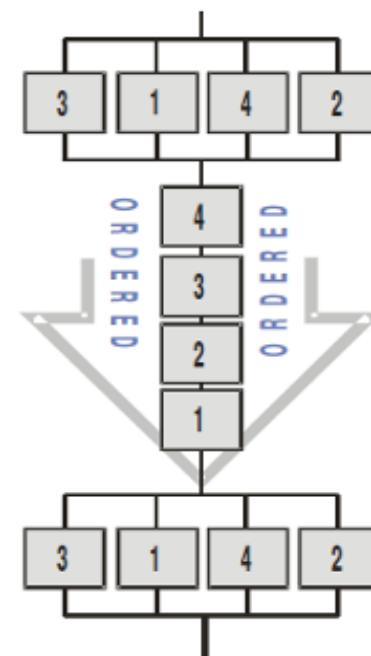
//Inicilializao da variavel dot
    dot = 0;
//Calculo do produto escalar
    #pragma omp parallel
    {
        #pragma omp single
            printf("numero de threads = %d      id = %d\n",omp_get_num_threads(),omp_get_thread_num());
        #pragma omp for reduction(+:dot)
            for(i=0;i<num;i++)
                dot += V1[i]*V2[i];
        #pragma omp single
            printf(" \n\n End Of Execution :::: dot = %d\n\n",dot);
    }
return 0;
}
```

Sincronizador: Ordered exemplo_04_ordered.c

```
#pragma omp ordered
```

- Garante que o loop será executado seqüencialmente
- Utilizado quando houver uma dependência dos dados atuais com as iterações anteriores.

```
#pragma omp parallel for ordered  
for(i = 1; i<= 4; i++)
```



Diretiva Threadprivate

```
#pragma omp threadprivate (lista de variáveis)
```

- Especifica quais variáveis serão privadas em todo o escopo do código.
- As variáveis serão privadas em todas as regiões paralelas.

• exemplo_9.c

```
{  
    int n,i;  
    static int var=0;  
    static int id=0;  
    #pragma omp threadprivate(var,id)  
    n = atoi(argv[1]);  
    omp_set_nested(1);  
  
    #pragma omp parallel if (n > 10)  
    {  
        #pragma omp for  
        for (i = 0; i < n; i++)  
        {  
            id= omp_get_thread_num();  
            var=var+id;  
            printf("Thread = %d, i= %d,var=%d \n",id,i,var);  
        }  
    }  
    printf("codigo serial Thread = %d, i= %d,var=%d \n",id,i,var);  
    #pragma omp parallel  
    {  
        #pragma omp for  
        for (i = 0; i < n; i++)  
        {  
            id= omp_get_thread_num();  
            var=var+id;  
            printf("loop 2 Thread = %d, i= %d,var=%d \n",id,i,var);  
        }  
    }  
}
```

Diretiva Threadprivate : cláusula copyin

- exemplo_09_copyin.c : permite que o valor da variável local (var) da thread master possa ser copiado para as variáveis privativas (var) das outras threads.
-

• exemplo_09_copyin.c

```
int n,i;
static int var=0;
static int id=0;
#pragma omp threadprivate(var,id)
n = atoi(argv[1]);
omp_set_nested(1);

#pragma omp parallel if (n > 10)
{
    #pragma omp for
    for (i = 0; i < n; i++)
    {
        id= omp_get_thread_num();
        var=var+id;
        printf("Thread = %d, i= %d,var=%d \n",id,i,var);
    }
    var=30000;
    printf("codigo serial Thread = %d, i= %d,var=%d \n",id,i,var);
#pragma omp parallel copyin(var)
{
    #pragma omp for
    for (i = 0; i < n; i++)
    {
        id= omp_get_thread_num();
        var=var+id;
        printf("loop 2 Thread = %d, i= %d,var=%d \n",id,i,var);
    }
}
```

Funções de tempo:

```
double omp_get_wtime(void)
```

- Retorna o valor decorrido em segundos relativo a um tempo passado.

```
double start, end;
start = omp_get_wtime();
...
end   = omp_get_wtime();

printf(Tempo decorrido%lf, end-start);
```

• exemplo_13.c

```
//Dimensao dos Vetores
    num = atoi(argv[1]);
//Alocacao dinamica dos vetores
    V1=(int*)malloc(num*sizeof(int));
    V2=(int*)malloc(num*sizeof(int));
//Inicializacao dos vetores
    InicializaVetores(num,V1,V2);
//Inicializao da variavel dot
    dot = 0;
    start = omp_get_wtime();
//Calculo do produto escalar
#pragma omp parallel
{
#pragma omp for reduction(+:dot)
    for(i=0;i<num;i++)
        dot += V1[i]*V2[i];
}
    printf(" \n\n End Of Execution :::: dot = %d\n\n\n",dot);
end = omp_get_wtime();
    printf(" \n\n Tempo decorrido :::: tempo = %lf\n\n\n",end-start);
return 0;
}
```

Funções de “lock”:

Procedimento para se utilizar as funções de Bloqueio:

- > Definir a variável “lock” (simples ou aninhada);
- > Inicializar a “lock” através da função `omp_init_lock()`;
- > Bloquear a “lock” através das funções `omp_set_lock()` e `omp_test_lock()`;
- > Desbloquear a “lock” através da função `omp_unset_lock()`;
- > Destruir a “lock” através da função `omp_destroy_lock()`.

Exemplo:

```
omp_lock_t * lock;
omp_init_lock(&lock);

#pragma omp parallel shared(lock,sum)
{
    #pragma omp for
    for(i=0;i<n;i++)
    {
        do_work(i);
        ...
        omp_set_lock(&lock);
        sum+=i;
        omp_unset_lock(&lock);
    }
}

omp_destroy_lock(&lock);
```

• exemplo_12.c

```
int main (int argc, char * argv[])
{
    int *V1,*V2;
    int i,j,id;
    int num_threads;
    int par,num,dot,aux_dot;
    omp_lock_t lock;
    omp_init_lock(&lock);
//Dimensao dos Vetores
    num = atoi(argv[1]);
//Alocacao dinamica dos vetores
    V1=(int*)malloc(num*sizeof(int));
    V2=(int*)malloc(num*sizeof(int));
//Inicializacao dos vetores
    InicializaVetores(num,V1,V2);
//Inicilizao da variavel dot
    dot = 0;
#pragma omp parallel shared(lock,dot)
    {
        int aux_dot = 0;
#pragma omp single
        printf("numeros de threads = %d      id = %d\n",omp_get_num_threads(),omp_get_thread_num());
#pragma omp for
        for(i=0;i<num;i++)
        {
            aux_dot += V1[i]*V2[i];
        }
        omp_set_lock(&lock);
        dot += aux_dot;
        omp_unset_lock(&lock);
    }
    omp_destroy_lock(&lock);
    printf(" \n\n End Of Execution :::: dot = %d\n\n",dot);
    return 0;
}
```

Geração de threads dentro de uma thread : “aninhado”:

```
void omp_set_nested(int num)
```

- Habilita ou desabilita o paralelismo aninhado.
 - > num = 0, desabilita o paralelismo aninhado
 - > num ≠ 0, habilita o paralelismo aninhado

```
int omp_get_nested(void)
```

- Retorna um valor que indica se o paralelismo aninhado está ativado ou não.
- Se o valor retornado for zero, o paralelismo aninhado está desativado. Se o valor retornado for diferente de zero, o paralelismo aninhado está ativado.

A thread que gera o paralelismo aninhado passa a ser a thread master(id=0) dentro do novo grupo de threads geradas por ela. Testar com o Exemplo_02.c

Exemplo:

```
omp_set_nested(0); // não seria necessário porque o default é zero.  
  
#pragma omp parallel num_threads(4)  
{  
    int id = omp_get_thread_num();  
  
    #pragma omp parallel num_threads(4)  
    {  
        printf(" thread %d do grupo da thread %d \n",omp_get_threads_num(),id);  
    }  
}
```

```
Thread 0 do grupo da thread 0  
Thread 0 do grupo da thread 1  
Thread 0 do grupo da thread 2  
Thread 0 do grupo da thread 3
```

VI-DIRETIVAS DE VETORIZAÇÃO

#pragma SIMD

#pragma vector aligned/unaligned

#pragma ivdep

#pragma vector always

#pragma no vector

Diretiva de vetorização do OpenMP

#pragma omp SIMD

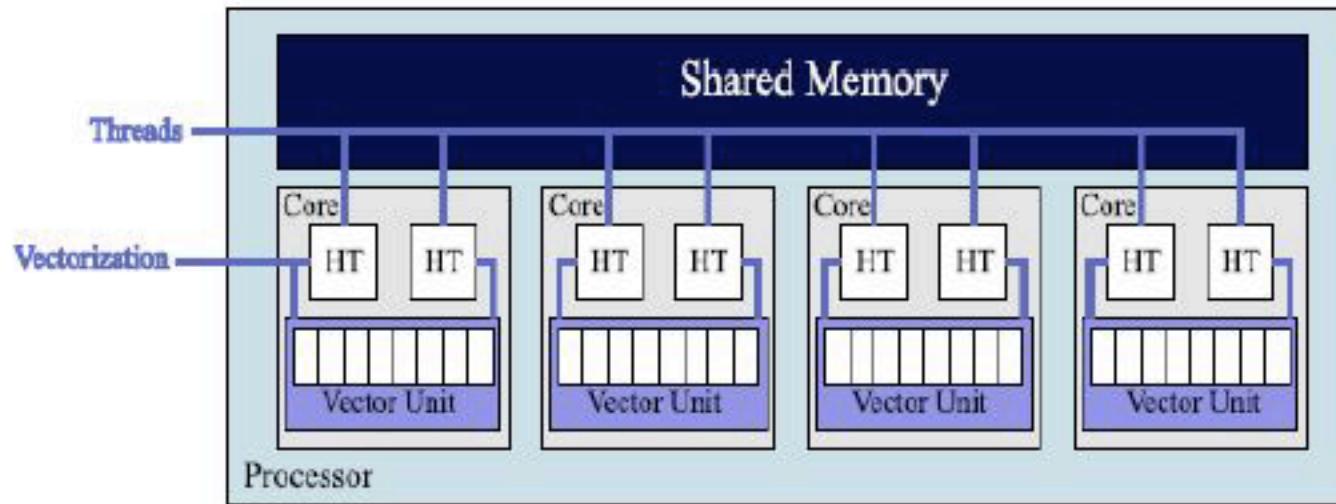
Existem diversos níveis de paralelismo em um processador multicore

PARALLELISM

21

CORES – multiple instructions on multiple data elements (MIMD)

VECTORS – single instruction on multiple data elements (SIMD)

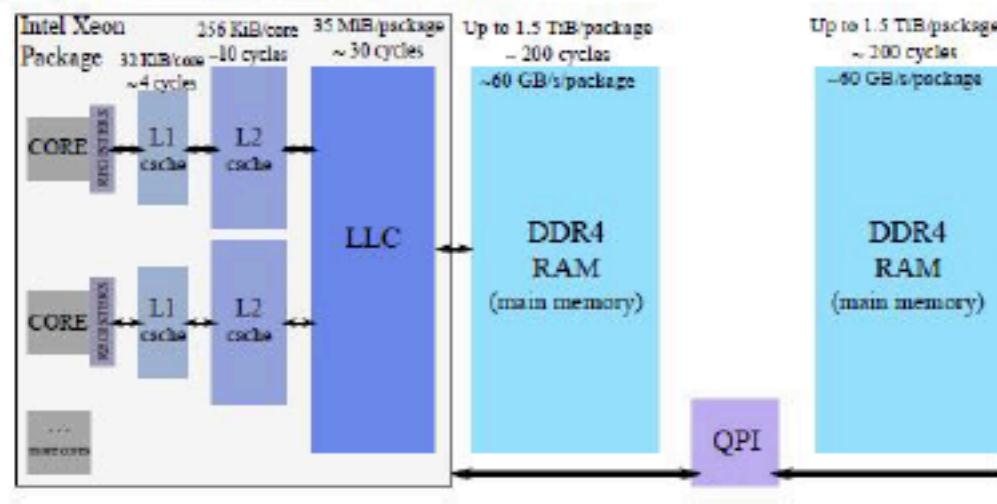


Unbounded growth opportunity, but **not automatic**

Cada core possui uma cache L1 e L2 e compartilha com os outros cores a cache LLC e a memória global RAM (no exemplo abaixo tem também a RAM de outro socket)

INTEL XEON CPU: MEMORY ORGANIZATION

- ▷ Hierarchical cache structure
- ▷ Two-way processors have NUMA architecture

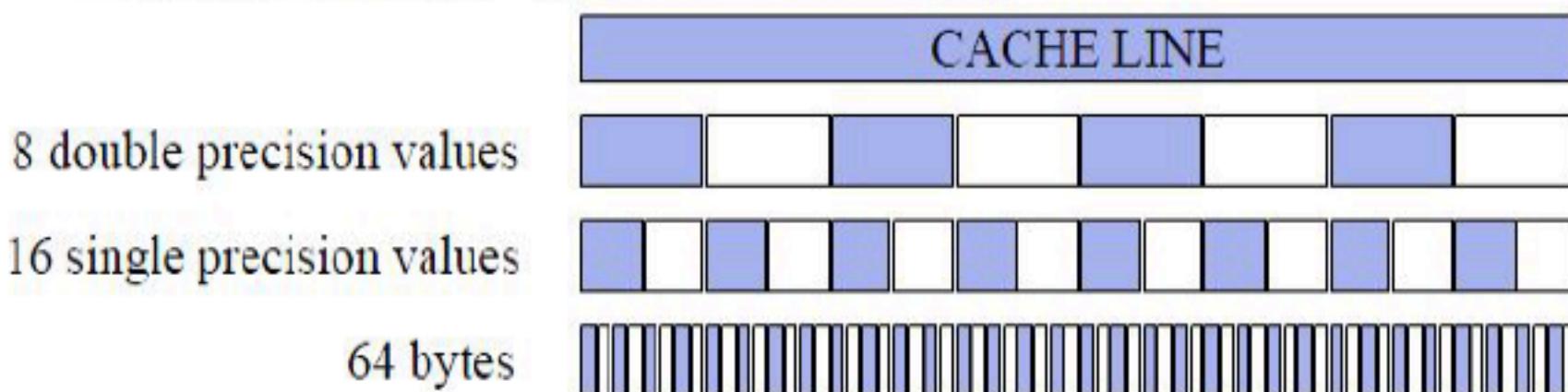


A transferência de um bloco de dados na arquitetura Intel é 64 bytes

CACHE LINES

17

- ▶ Minimal block of data transferred between memory and cache
- ▶ 64 bytes long in Intel Architecture
- ▶ Aligned on 64-byte boundaries in memory

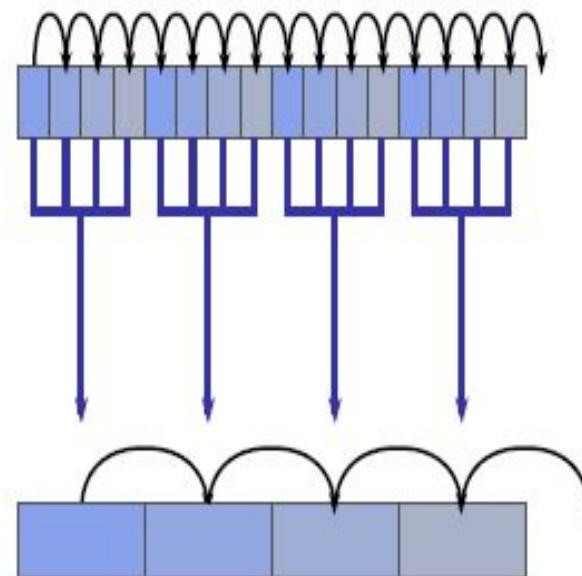


Abaixo são apresentada 4 limitações da implementação da vetorização automática pelo compilador:

LIMITATIONS ON AUTOMATIC VECTORIZATION

- ▶ Innermost loops*
- ▶ Known number of iterations
- ▶ No vector dependence
- ▶ Functions must be SIMD-enabled

* `#pragma omp simd` to override



A diretiva #pragma SIMD serve para ajudar na vetorização automática quando o compilador desconhece o numero de iterações

SIMULTANEOUS THREADING AND VECTORIZATION

Sometimes the compiler may need a little help:

```
1 const int STRIP_SIZE = 128; // A multiple of vector length
2 const int nTrunc = n - n%STRIP_SIZE; // A multiple of vector length
3
4 #pragma omp parallel for
5 for (int ii = 0; ii < nTrunc; ii += STRIP_SIZE) // Thread parallelism in outer
6 #pragma SIMD
7     for (int i = ii; i < ii + STRIP_SIZE; i++) // Vectorization in inner loop
8         DoSomeWork(A[i]);
9
10 // Remainder loop:
11 for (int i = nTrunc; i < n; i++)
12     DoSomeWork(A[i]);
```

#pragma omp simd

VECTORIZE MORE LOOPS: #pragma omp simd

Used to “enforce vectorization of loops”, which includes:

- ▷ Loops with SIMD-enabled functions
- ▷ Second innermost loops
- ▷ Failed vectorization due to compiler decision
- ▷ Where guidance is required (vector length, reduction, etc.)

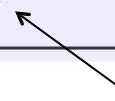
See OpenMP reference for syntax; #pragma simd

Vetorização no Second Innermost loop

Vetorização automática é ineficiente porque o “inner loop” (i) possui “stride access”.

EXAMPLE FOR #pragma omp simd

```
1 const int N=128, T=4;
2 float A[N*N], B[N*N], C[T*T];
3
4 for (int jj = 0; jj < N; jj+=T) // Tile in j
5     for (int ii = 0; ii < N; ii+=T) // and tile in i
6 #pragma omp simd    // Vectorize outer loop
7     for (int k = 0; k < N; ++k) // long loop, vectorize it
8         for (int i = 0; i < T; i++) { // Loop between ii and ii+T
9             // Instead of a loop between jj and jj+T, unrolling that loop:
10            C[0*T + i] += A[(jj+0)*N + k]*B[(ii+i)*N + k];
11            C[1*T + i] += A[(jj+1)*N + k]*B[(ii+i)*N + k];
12            C[2*T + i] += A[(jj+2)*N + k]*B[(ii+i)*N + k];
13            C[3*T + i] += A[(jj+3)*N + k]*B[(ii+i)*N + k];
14 }
```



Loop with SIMD-enable function

Diretivas para o compilador vetorizar funções é necessária quando a função está em arquivo separado do arquivo fonte (biblioteca).

SIMD-ENABLED FUNCTIONS

3

Define function in one file (e.g., library), use in another

```
1 // Compiler will produce 3 versions:  
2 #pragma omp declare simd  
3 float my_simple_add(float x1, float x2){  
4     return x1 + x2;  
5 }
```

```
1 // May be in a separate file  
2 #pragma omp simd  
3 for (int i = 0; i < N, ++i) {  
4     output[i] = my_simple_add(inputa[i], inputb[i]);  
5 }
```

OTIMIZAÇÕES PARA UMA VETORIZAÇÃO EFICIENTE

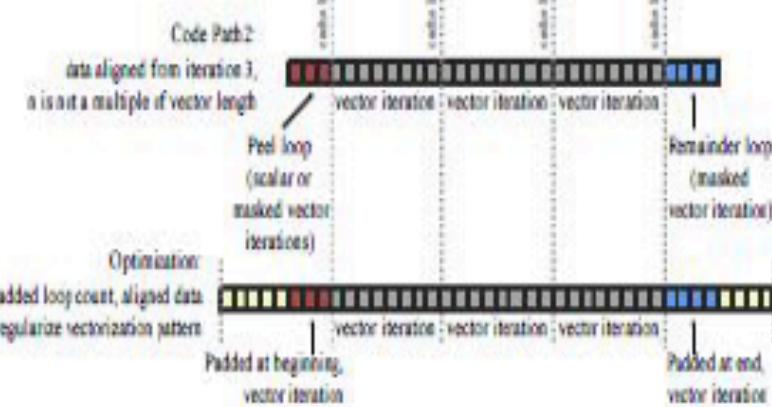
LOOP WAS VECTORIZED, NOW WHAT?

1. Unit-stride access
2. Data alignment
3. Container padding
4. Eliminate peel loops
5. Eliminate multiversoring
6. **Optimize data re-use in caches**

```
for (i = 0; i < n; i++) A[i] = ...
```

Code Path 1:
data aligned from iteration 0,
 n is multiple of vector length

Code Path 2:
data aligned from iteration 3,
 n is not a multiple of vector length

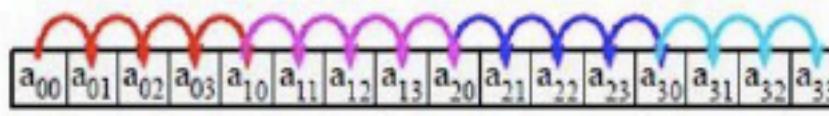


O loop deve percorrer a matriz evitando acessos de stride:

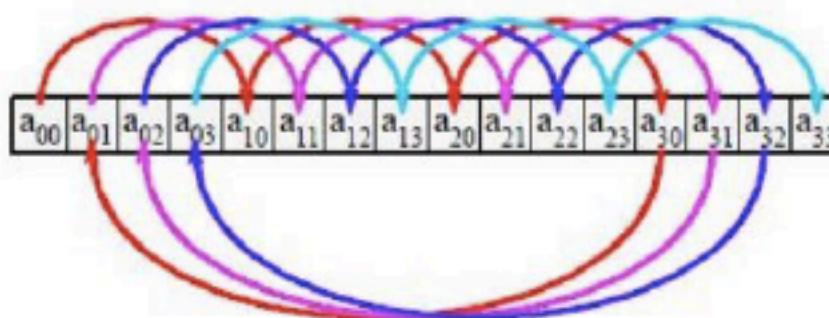
PRINCIPLE

Choose loop order to maintain unit-stride memory access

a ₀₀	a ₀₁	a ₀₂	a ₀₃
a ₁₀	a ₁₁	a ₁₂	a ₁₃
a ₂₀	a ₂₁	a ₂₂	a ₂₃
a ₃₀	a ₃₁	a ₃₂	a ₃₃



a ₀₀	a ₀₁	a ₀₂	a ₀₃
a ₁₀	a ₁₁	a ₁₂	a ₁₃
a ₂₀	a ₂₁	a ₂₂	a ₂₃
a ₃₀	a ₃₁	a ₃₂	a ₃₃



Compiler may or may not be able to automate loop permutation.

O “inner” loop passa para k, para percorrer a matriz sem “stride”

EXAMPLE: OVER-SIMPLIFIED MATRIX-MATRIX MULTIPLICATION

27

$$C = AB \Leftrightarrow C_{ij} = \sum_{k=0}^{n-1} A_{ik}B_{kj}$$

Before:

```
1 #pragma omp parallel for
2 for (int i = 0; i < n; i++)
3     for (int j = 0; j < n; j++)
4 #pragma vector aligned
5     for (int k = 0; k < n; k++)
6         C[i*n+j] += A[i*n+k]*B[k*n+j];
```

After:

```
1 #pragma omp parallel for
2 for (int i = 0; i < n; i++)
3     for (int k = 0; k < n; k++)
4 #pragma vector aligned
5     for (int j = 0; j < n; j++)
6         C[i*n+j] += A[i*n+k]*B[k*n+j];
```

sem_otimizacao.c

```
//Aloca a memória necessária para as matrizes
long int n = atol(argv[1]), i;
printf("Allocating memory...\\n");
int *A = (int*) malloc(n * n * sizeof(int));
int *B = (int*) malloc(n * n * sizeof(int));

//Inicializa A e B com os valores 2 e 3
printf("Initializing Matrices...\\n");
for(int i = 0;i < n;i++){
    for(int j = 0;i < n;j++){
        A[i*n + j] = 2;
        B[i*n + j] = 3;
    }
}

//Aloca memória para matriz C
int *C = (int*) malloc(n * n * sizeof(int));

//Multiplica A * B
#pragma omp parallel for
for(int i = 0;i < n;i++){
    for(int j = 0;j < n;j++){
#pragma vector aligned
        for(int k = 0;k < n;k++){
            C[i*n + j] += A[i*n + k] * B[k*n + j];
        }
    }
}
```

com_otimizacao.c

```
// Aloca a memória necessária para as matrizes
long int n = atol(argv[1]), i;
printf("Allocating memory...\n");
int *A = (int*) malloc(n * n * sizeof(int));
int *B = (int*) malloc(n * n * sizeof(int));

// Inicializa A e B com os valores 2 e 3
printf("Initializing Matrices...\n");
for(int i = 0; i < n; i++){
    for(int j = 0; j < n; j++){
        A[i*n + j] = 2;
        B[i*n + j] = 3;
    }
}

// Aloca memória para matriz C
int *C = (int*) malloc(n * n * sizeof(int));

#pragma omp parallel for
for(int i = 0; i < n; i++){
    for(int k = 0; k < n; k++){
#pragma vector aligned
        for(int j = 0; j < n; j++){
            C[i*n + j] += A[i*n + k] * B[k*n + j];
        }
    }
}
```

CONTAINER PADDING

Para não ter que align todas as variáveis pode escolher uma dimensão maior, múltiplo do vetor :

PADDING MULTI-DIMENSIONAL CONTAINERS FOR ALIGNMENT

To use aligned instructions, you may need to pad inner dimension of multi-dimensional arrays to a multiple of 16 (in SP) or 8 (DP) elements.

Incorrect:

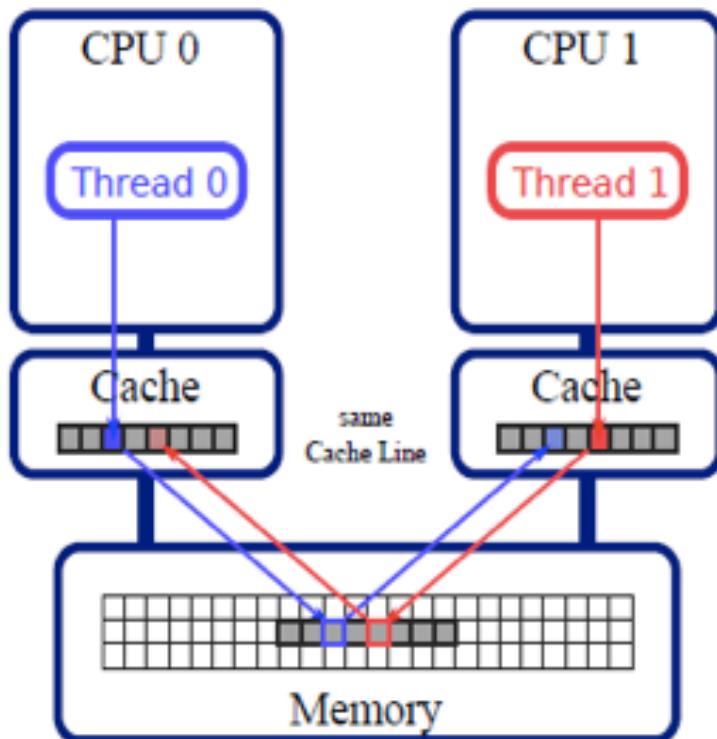
```
1 // A - matrix of size (n x n)
2 // n is not a multiple of 16
3 float* A =
4 _mm_malloc(sizeof(float)*n*n, 64);
5
6 for (int i = 0; i < n; i++)
7     // A[i*n + 0] may be unaligned
8     for (int j = 0; j < n; j++)
9         A[i*n + j] = ...
```

Correct:

```
1 // ... Padding inner dimension
2 int lda=n + (16-n%16); // lda%16==0
3 float* A =
4 _mm_malloc(sizeof(float)*n*lda, 64);
5
6 for (int i = 0; i < n; i++)
7     // A[i*lda + 0] aligned for any i
8     for (int j = 0; j < n; j++)
9         A[i*lda + j] = ...
```

FALSE SHARING. DATA PADDING AND PRIVATE VARIABLES

2

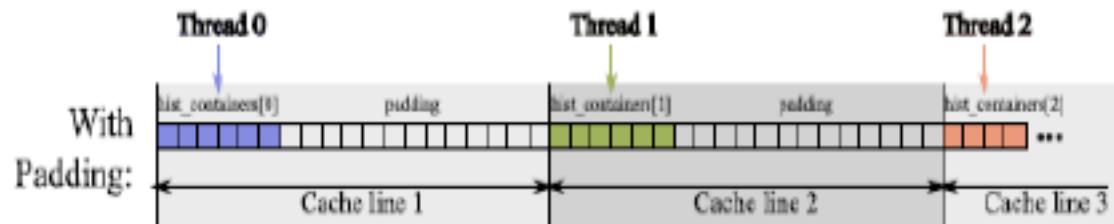
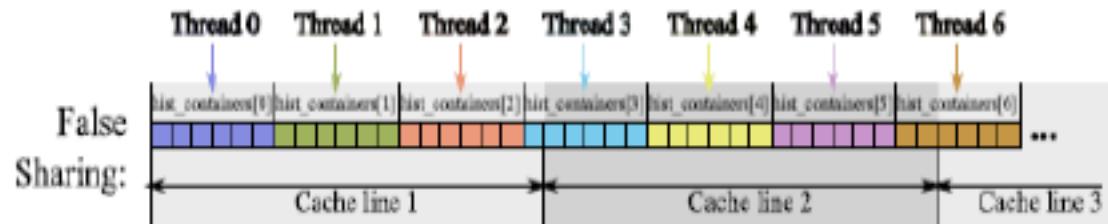


- ▶ Occurs when 2 or more threads access the same cache line, and at least one of the accesses is for writing
- ▶ Caused by *coherent caches*
- ▶ Cache line is 64-byte wide (in modern Intel architectures)

Solução para evitar false sharing: Padding + alignment

2

PADDING TO AVOID FALSE SHARING



```
1 // Padding to avoid sharing a cache line between threads
2 const int paddingBytes = 64;
3 const int paddingElements = paddingBytes / sizeof(int);
4 const int mPadded = m + (paddingElements-m%paddingElements);
5 int histContainers[nThreads][mPadded]; // New container
```

Ao alterar o tamanho é necessário usar o #pragma vector aligned

DATA ALIGNMENT HINTS

Programmer may promise to the compiler (under penalty of segmentation fault) that alignment has been taken care of:

```
1 // Promising that A[i*lda + 0] is aligned for every i
2 // and the same for every other array in this loop
3 #pragma vector aligned
4     for (int j = 0; j < n; j++)
5         A[i*lda + j] -= ...
```

This can lead to significant speedups, because compiler will not implement runtime checks for alignment situation and *peel loops*.

OPTIMIZE DATA REUSE IN CACHE

Strip-Mining

STRIP-MINING FOR VECTORIZATION

- ▶ Programming technique that turns one loop into two nested loops.
- ▶ Used to expose vectorization opportunities.

Original:

```
1 for (int i = 0; i < n; i++) {  
2     // ... do work  
3 }
```

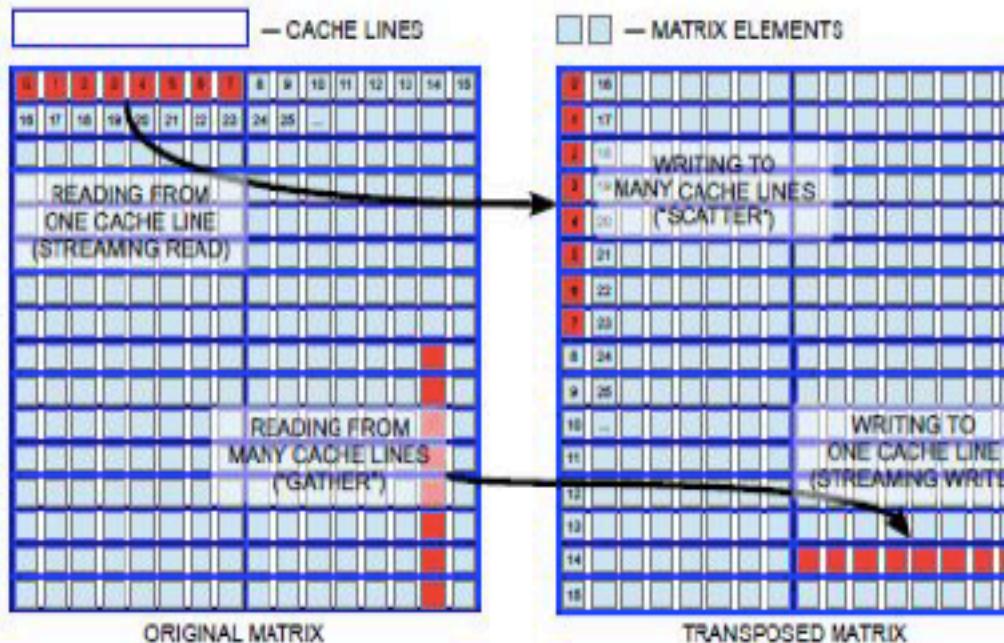
Strip-mined:

```
1 const int STRIP=1024;  
2 const int nPrime = n - n%STRIP;  
3 for (int ii=0; ii<nPrime; ii+=STRIP)  
4     for (int i=ii; i<ii+STRIP; i++)  
5         // ... do work  
6  
7 for (int i=nPrime; i<n; i++)  
8     // ... do work
```

Exemplo com código de transposição de matriz

LOOP TILING EXAMPLE: MATRIX TRANSPOSITION

$$B = A^T \Leftrightarrow B_{ij} = A_{ji}$$



Implementar os códigos abaixo com tile múltiplo de 64

MATRIX TRANSPOSITION

Before:

```
1 #pragma omp parallel for
2 for (int i = 0; i < n; i++)
3     for (int j = 0; j < n; j++)
4         B[i*n + j] = A[j*n + i];
```

After:

```
1 const int tile = 200;
2 if (n%tile != 0) exit(1);
3
4 #pragma omp parallel for
5 for (int ii=0; ii<n; ii+=tile)
6     for (int jj=0; jj<n; jj+=tile)
7         for (int i=ii; i<ii+tile; i++)
8             for (int j=jj; j<jj+tile; j++)
9                 B[i*n + j] = A[j*n + i];
```

EXAMPLE: MATRIX-VECTOR MULTIPLICATION

$$c_i = \sum_{j=0}^n A_{ij} b_j, \quad i = 0, 1, \dots, (m-1). \quad (1)$$

```
1 void Multiply(const double* const A, const double* const b,
2                 double* const c, const long n, const long m){
3     assert(n%64 == 0);
4     #pragma omp parallel for
5     for (long i = 0; i < m; i++)
6     #pragma vector aligned
7         for (long j = 0; j < n; j++) // Each value of A[i*n+j] is used only once
8             c[i] += A[i*n+j] * b[j]; // Each value of b[j] is used a total of m times
9 }
```

Non-optimal performance due to inefficient cache use

APPLYING TILING

```
1 const long jTile = 4096L; assert(n%jTile == 0);
2 #pragma omp parallel
3 {
4     double temp_c[m] __attribute__((aligned(64)));
5     temp_c[:] =0;
6 #pragma omp for
7     for (long jj =0; jj < n; jj+=jTile) // Loop Tiling in j
8         for (long i = 0; i < m; i++)
9 #pragma vector aligned
10        for (long j =jj; j < jj+jTile; j++)
11            temp_c[i] += A[i*n+j] * b[j];
12
13    for(long i = 0; i < m; i++) { // Reduction
14 #pragma omp atomic
15     c[i] += temp_c[i];
16 } } }
```



Home » Training

Training

Browse our webinars and on-demand training on parallel programming, performance optimization for parallel processors, new technology and development tools. Some of our training programs feature remote access to training servers, original programming exercises and certificates of accomplishment.

“HOW” Series: Deep Dive

The image is a promotional graphic for a training series. It features a blue-toned background with two scuba divers swimming in an underwater environment. One diver is in the foreground, and another is slightly behind and to the left. The water has a bright, sunlit appearance at the top. In the upper right corner, there is a red rectangular overlay containing white text. The text reads "THE 'HOW' SERIES TRAINING" in a bold, sans-serif font, and below it, "DEEP DIVE" in a large, stylized, blocky font. Underneath "DEEP DIVE", smaller text reads "WITH CODE MODERNIZATION EXPERTS". At the bottom of the graphic, there is a blue button-like shape containing the text "It's free". In the bottom right corner of the main image area, there is a dark, rounded rectangular button with the text "Learn More" in white. At the very bottom of the image, there is a small video player interface showing a timestamp of "0:45*", a progress bar, and some control icons.

Learn More