

Лекция 5

Многопоточное программирование Стандарт OpenMP

Курносов Михаил Георгиевич

E-mail: mkurnosov@gmail.com

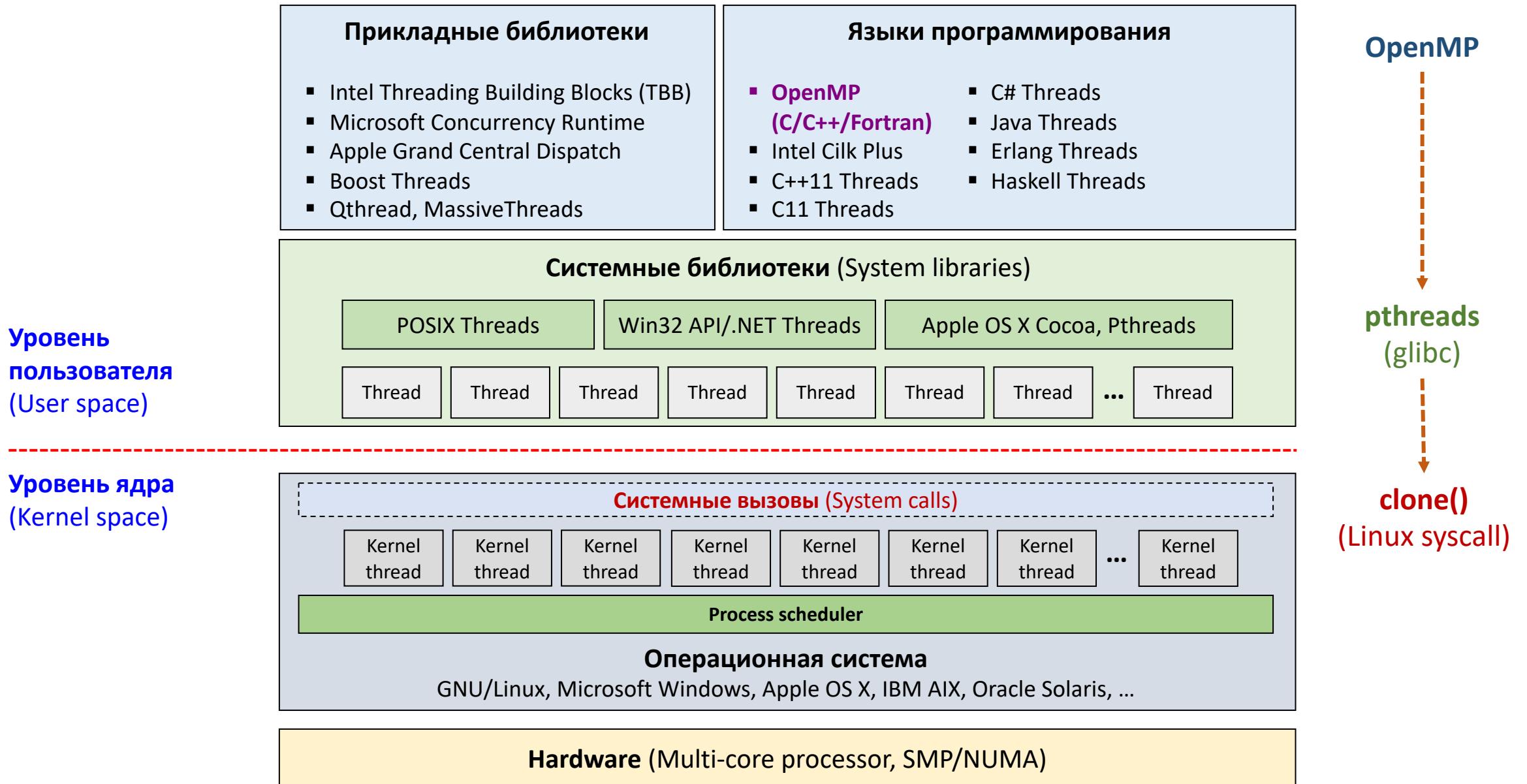
WWW: www.mkurnosov.net

Курс «Параллельные вычислительные технологии»

Сибирский государственный университет телекоммуникаций и информатики (г. Новосибирск)

Весенний семестр, 2020

Средства многопоточного программирования



Стандарт OpenMP

- **OpenMP (Open Multi-Processing)** – стандарт, определяющий набор директив компилятора, библиотечных процедур и переменных среды окружения для создания многопоточных программ
- Разрабатывается в рамках OpenMP Architecture Review Board с 1997 года
 - ❑ OpenMP 2.5 (2005), OpenMP 3.0 (2008), OpenMP 4.5 (2015), **OpenMP 5.0 (2018)**
 - ❑ <http://www.openmp.org>
 - ❑ <https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5.0.pdf>
- Требуется поддержка со стороны компилятора



Стандарт OpenMP

<http://www.openmp.org/wp-content/uploads/openmp-4.5.pdf>



OpenMP Application Programming Interface

Version 4.5 November 2015

Copyright © 1997-2015 OpenMP Architecture Review Board.
Permission to copy without fee all or part of this material is granted, provided the OpenMP Architecture Review Board copyright notice and the title of this document appear. Notice is given that copying is by permission of OpenMP Architecture Review Board.

Contents

1	Introduction	1
1.1	Scope	1
1.2	Glossary	2
1.2.1	Threading Concepts	2
1.2.2	OpenMP Language Terminology	2
1.2.3	Loop Terminology	8
1.2.4	Synchronization Terminology	9
1.2.5	Tasking Terminology	9
1.2.6	Data Terminology	11
1.2.7	Implementation Terminology	13
1.3	Execution Model	14
1.4	Memory Model	17
1.4.1	Structure of the OpenMP Memory Model	17
1.4.2	Device Data Environments	18
1.4.3	The Flush Operation	19
1.4.4	OpenMP Memory Consistency	20
1.5	OpenMP Compliance	21
1.6	Normative References	21
1.7	Organization of this Document	23
2	Directives	25
2.1	Directive Format	26
2.1.1	Fixed Source Form Directives	28
2.1.2	Free Source Form Directives	29
2.1.3	Stand-Alone Directives	32
2.2	Conditional Compilation	33
2.2.1	Fixed Source Form Conditional Compilation Sentinels	34

2.2.2	Free Source Form Conditional Compilation Sentinel	34
2.3	Internal Control Variables	36
2.3.1	ICV Descriptions	36
2.3.2	ICV Initialization	37
2.3.3	Modifying and Retrieving ICV Values	39
2.3.4	How ICVs are Scoped	41
2.3.4.1	How the Per-Data Environment ICVs Work	42
2.3.5	ICV Override Relationships	43
2.4	Array Sections	44
2.5	parallel Construct	46
2.5.1	Determining the Number of Threads for a parallel Region	50
2.5.2	Controlling OpenMP Thread Affinity	52
2.6	Canonical Loop Form	53
2.7	Worksharing Constructs	56
2.7.1	Loop Construct	56
2.7.1.1	Determining the Schedule of a Worksharing Loop	64
2.7.2	sections Construct	65
2.7.3	single Construct	67
2.7.4	workshare Construct	69
2.8	SIMD Constructs	72
2.8.1	simd Construct	72
2.8.2	declare simd Construct	76
2.8.3	Loop SIMD Construct	81
2.9	Tasking Constructs	83
2.9.1	task Construct	83
2.9.2	taskloop Construct	87
2.9.3	taskloop simd Construct	91
2.9.4	taskyield Construct	93
2.9.5	Task Scheduling	94
2.10	Device Constructs	95
2.10.1	target data Construct	95
2.10.2	target enter data Construct	97
2.10.3	target exit data Construct	100

Поддержка компиляторами

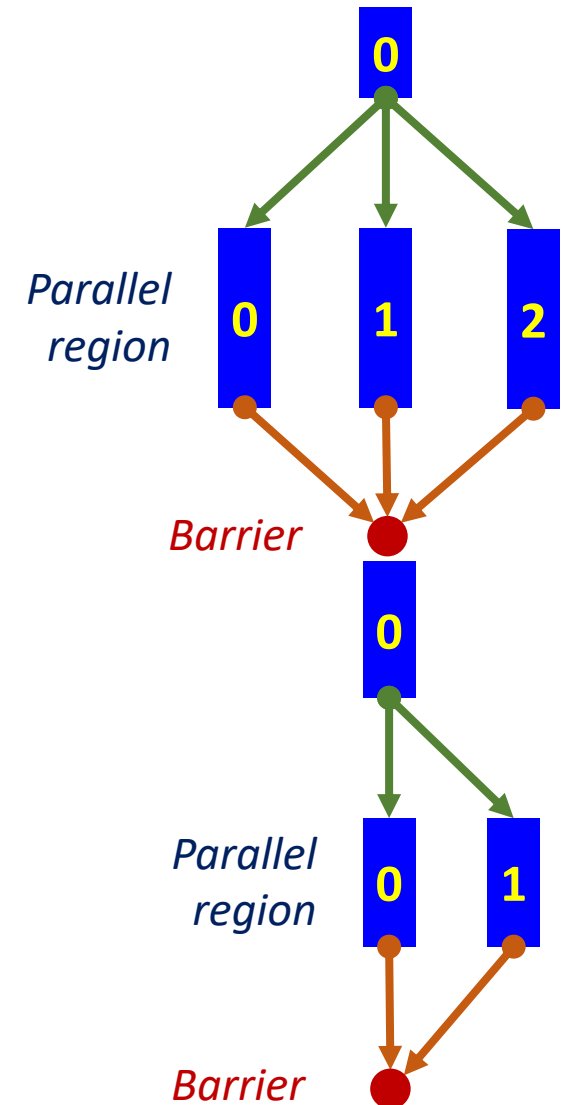
Compiler	Information
GNU GCC https://gcc.gnu.org/wiki/openmp	Option: -fopenmp gcc 4.7 – OpenMP 3.1 gcc 4.9 – OpenMP 4.0 gcc 5.x – Offloading
Clang (LLVM) http://openmp.llvm.org/	Option: -fopenmp Clang 3.8.0 – OpenMP 3.1 Clang + Intel OpenMP RTL (http://clang-omp.github.io/)
Intel C/C++, Fortran https://software.intel.com/en-us/c-compilers/	OpenMP 4.x Option: -Qopenmp, -openmp
Oracle Solaris Studio C/C++/Fortran	OpenMP 4.0 Option: -xopenmp
Microsoft Visual Studio C++	Option: /openmp OpenMP 2.0 only
http://www.openmp.org/resources/openmp-compilers/	

Модель выполнения OpenMP-программы

- **Динамическое управление потоками в модели Fork-Join:**

- ✓ **Fork** – порождение нового потока
- ✓ **Join** – ожидание завершения потока (объединение потоков управления)

- OpenMP-программа – совокупность последовательных участков кода (serial code) и параллельных регионов (parallel region)
- Каждый поток имеет логический номер: 0, 1, 2, ...
- Главный поток (master) имеет номер 0
- Параллельные регионы могут быть вложенными



Hello, World

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

int main(int argc, char **argv)
{
    #pragma omp parallel    /* <-- Fork */
    {
        printf("Hello, multithreaded world: thread %d of %d\n",
               omp_get_thread_num(), omp_get_num_threads());

    }                          /* <-- Barrier & join */

    return 0;
}
```

Компиляция и запуск OpenMP-программы

```
$ gcc -fopenmp -o hello ./hello.c
```

```
$ ./hello
```

```
Hello, multithreaded world: thread 0 of 4
```

```
Hello, multithreaded world: thread 1 of 4
```

```
Hello, multithreaded world: thread 3 of 4
```

```
Hello, multithreaded world: thread 2 of 4
```

- По умолчанию количество потоков в параллельном регионе равно числу логических процессоров в системе
- Порядок выполнения потоков заранее неизвестен – определяется планировщиком операционной системы

Указание числа потоков в параллельных регионах

```
$ export OMP_NUM_THREADS=8
```

```
$ ./hello
```

```
Hello, multithreaded world: thread 1 of 8
```

```
Hello, multithreaded world: thread 2 of 8
```

```
Hello, multithreaded world: thread 3 of 8
```

```
Hello, multithreaded world: thread 0 of 8
```

```
Hello, multithreaded world: thread 4 of 8
```

```
Hello, multithreaded world: thread 5 of 8
```

```
Hello, multithreaded world: thread 6 of 8
```

```
Hello, multithreaded world: thread 7 of 8
```

Задание числа потоков в параллельном регионе

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

int main(int argc, char **argv)
{
    #pragma omp parallel num_threads(6)
    {
        printf("Hello, multithreaded world: thread %d of %d\n",
               omp_get_thread_num(), omp_get_num_threads());
    }

    return 0;
}
```

Задание числа потоков в параллельном регионе

```
$ export OMP_NUM_THREADS=8
```

```
$ ./hello
```

```
Hello, multithreaded world: thread 2 of 6
```

```
Hello, multithreaded world: thread 3 of 6
```

```
Hello, multithreaded world: thread 1 of 6
```

```
Hello, multithreaded world: thread 0 of 6
```

```
Hello, multithreaded world: thread 4 of 6
```

```
Hello, multithreaded world: thread 5 of 6
```

- Директива num_threads имеет приоритет над значением переменной среды окружения OMP_NUM_THREADS

Список потоков процесса

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

int main(int argc, char **argv)
{
    #pragma omp parallel num_threads(6)
    {
        printf("Hello, multithreaded world: thread %d of %d\n",
            omp_get_thread_num(), omp_get_num_threads());

        /* Sleep for 30 seconds */
        nanosleep(&(struct timespec){.tv_sec = 30}, NULL);
    }
    return 0;
}
```

Список потоков процесса

```
$ ./hello &  
$ ps -eLo pid,tid,psr,args | grep hello  
6157 6157 0 ./hello  
6157 6158 1 ./hello  
6157 6159 0 ./hello  
6157 6160 1 ./hello  
6157 6161 0 ./hello  
6157 6162 1 ./hello  
6165 6165 2 grep hello
```

- Номер процесса (PID)
- Номер потока (TID)
- Логический процессор (PSR)
- Название исполняемого файла

- **Информация о логических процессорах системы:**

- ☐ /proc/cpuinfo
- ☐ /sys/devices/system/cpu

#pragma omp parallel

<http://www.openmp.org/wp-content/uploads/openmp-4.5.pdf>

1 Syntax

C / C++

2 The syntax of the **parallel** construct is as follows:

```
#pragma omp parallel [clause[ [, ] clause] ... ] new-line
    structured-block
```

3 where *clause* is one of the following:

4 **if**(**[parallel :]** *scalar-expression*)

5 **num_threads**(*integer-expression*)

6 **default**(**shared** | **none**)

7 **private**(*list*)

8 **firstprivate**(*list*)

9 **shared**(*list*)

10 **copyin**(*list*)

11 **reduction**(*reduction-identifier* : *list*)

12 **proc_bind**(**master** | **close** | **spread**)

C / C++

#pragma omp parallel

12

Binding

13

The binding thread set for a **parallel** region is the encountering thread. The encountering thread becomes the master thread of the new team.

14

15

Description

16

When a thread encounters a **parallel** construct, a team of threads is created to execute the **parallel** region (see Section 2.5.1 on page 50 for more information about how the number of threads in the team is determined, including the evaluation of the **if** and **num_threads** clauses). The thread that encountered the **parallel** construct becomes the master thread of the new team, with a thread number of zero for the duration of the new **parallel** region. All threads in the new team, including the master thread, execute the region. Once the team is created, the number of threads in the team remains constant for the duration of that **parallel** region.

17

18

19

20

21

22

23

The optional **proc_bind** clause, described in Section 2.5.2 on page 52, specifies the mapping of OpenMP threads to places within the current place partition, that is, within the places listed in the *place-partition-var* ICV for the implicit task of the encountering thread.

24

25

26

Within a **parallel** region, thread numbers uniquely identify each thread. Thread numbers are consecutive whole numbers ranging from zero for the master thread up to one less than the number of threads in the team. A thread may obtain its own thread number by a call to the **omp_get_thread_num** library routine.

27

28

29

30

A set of implicit tasks, equal in number to the number of threads in the team, is generated by the encountering thread. The structured block of the **parallel** construct determines the code that

31

#pragma omp parallel

1 will be executed in each implicit task. Each task is assigned to a different thread in the team and
2 becomes tied. The task region of the task being executed by the encountering thread is suspended
3 and each thread in the team executes its implicit task. Each thread can execute a path of statements
4 that is different from that of the other threads

5 The implementation may cause any thread to suspend execution of its implicit task at a task
6 scheduling point, and switch to execute any explicit task generated by any of the threads in the
7 team, before eventually resuming execution of the implicit task (for more details see Section 2.9 on
8 page 83).

9 There is an implied barrier at the end of a **parallel** region. After the end of a **parallel**
10 region, only the master thread of the team resumes execution of the enclosing task region.

11 If a thread in a team executing a **parallel** region encounters another **parallel** directive, it
12 creates a new team, according to the rules in Section 2.5.1 on page 50, and it becomes the master of
13 that new team.

14 If execution of a thread terminates while inside a **parallel** region, execution of all threads in all
15 teams terminates. The order of termination of threads is unspecified. All work done by a team prior
16 to any barrier that the team has passed in the program is guaranteed to be complete. The amount of
17 work done by each thread after the last barrier that it passed and before it terminates is unspecified.

18 Restrictions

19 Restrictions to the **parallel** construct are as follows:

- 20 • A program that branches into or out of a **parallel** region is non-conforming.
- 21 • A program must not depend on any ordering of the evaluations of the clauses of the **parallel**
22 directive, or on any side effects of the evaluations of the clauses.
- 23 • At most one **if** clause can appear on the directive.
- 24 • At most one **proc_bind** clause can appear on the directive.
- 25 • At most one **num_threads** clause can appear on the directive. The **num_threads**
26 expression must evaluate to a positive integer value.

Умножение матрицы на вектор

DGEMV – BLAS Level 2

(BLAS – Basic Linear Algebra Subroutines)

Умножение матрицы на вектор (DGEMV)

- Требуется вычислить произведение прямоугольной матрицы **A** размера $m \times n$ на вектор-столбец **B** размера $n \times 1$ (BLAS Level 2, DGEMV)

$$\mathbf{C}_{m \times 1} = \mathbf{A}_{m \times n} \cdot \mathbf{B}_{n \times 1}$$

$$C = \begin{pmatrix} c_1 \\ c_2 \\ \dots \\ c_m \end{pmatrix} \quad A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \quad B = \begin{pmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{pmatrix}$$

$$c_i = \sum_{j=1}^n a_{ij} \cdot b_j, \quad i = 1, 2, \dots, m.$$

DGEMV: последовательная версия

```
/*  
 * matrix_vector_product: Compute matrix-vector product  $c[m] = a[m][n] * b[n]$   
 */  
void matrix_vector_product(double *a, double *b, double *c, int m, int n)  
{  
    for (int i = 0; i < m; i++) {  
        c[i] = 0.0;  
        for (int j = 0; j < n; j++)  
            c[i] += a[i * n + j] * b[j];  
    }  
}
```

$$c_i = \sum_{j=1}^n a_{ij} \cdot b_j, \quad i = 1, 2, \dots, m.$$

DGEMV: последовательная версия

```
void run_serial()
{
    double *a, *b, *c;

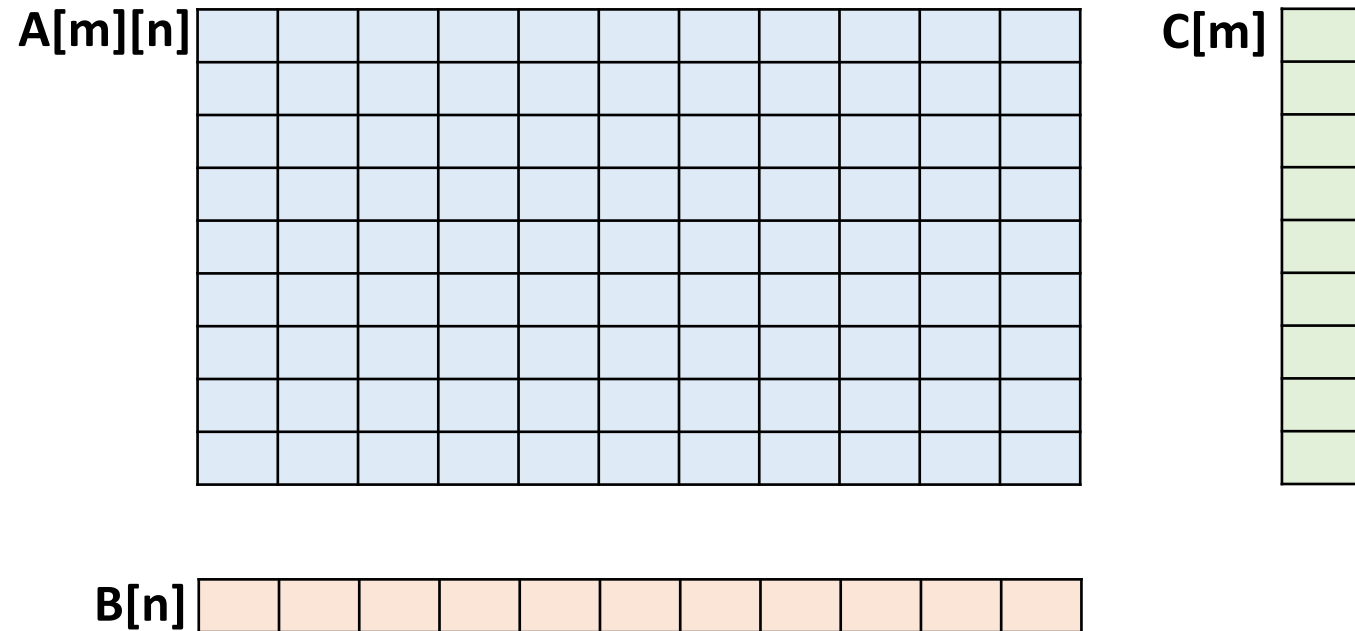
    a = xmalloc(sizeof(*a) * m * n);
    b = xmalloc(sizeof(*b) * n);
    c = xmalloc(sizeof(*c) * m);

    for (int i = 0; i < m; i++) {
        for (int j = 0; j < n; j++)
            a[i * n + j] = i + j;
    }
    for (int j = 0; j < n; j++)
        b[j] = j;

    double t = wtime();
    matrix_vector_product(a, b, c, m, n);
    t = wtime() - t;

    printf("Elapsed time (serial): %.6f sec.\n", t);
    free(a);
    free(b);
    free(c);
}
```

DGEMV: параллельная версия

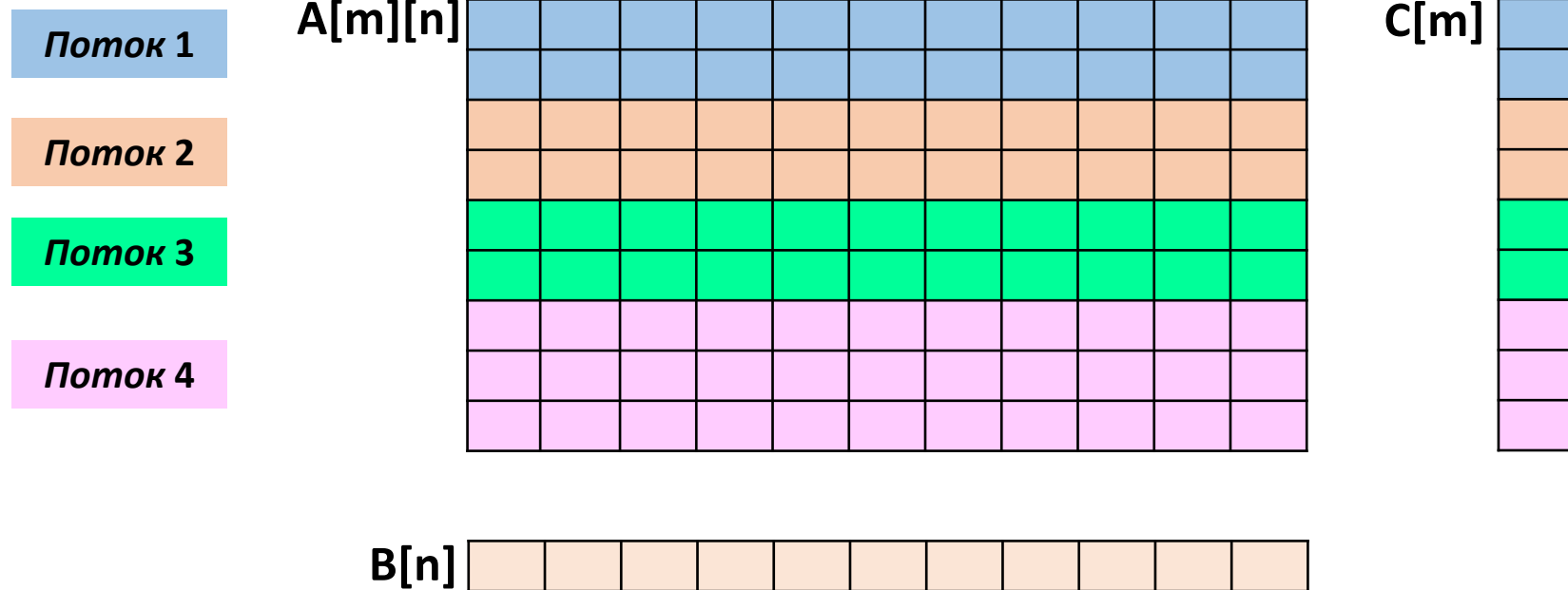


```
for (int i = 0; i < m; i++) {  
    c[i] = 0.0;  
    for (int j = 0; j < n; j++)  
        c[i] += a[i * n + j] * b[j];  
}
```

Требования к параллельному алгоритму

- Максимальная загрузка потоков вычислениями
- Минимум совместно используемых ячеек памяти – независимые области данных

DGEMV: параллельная версия



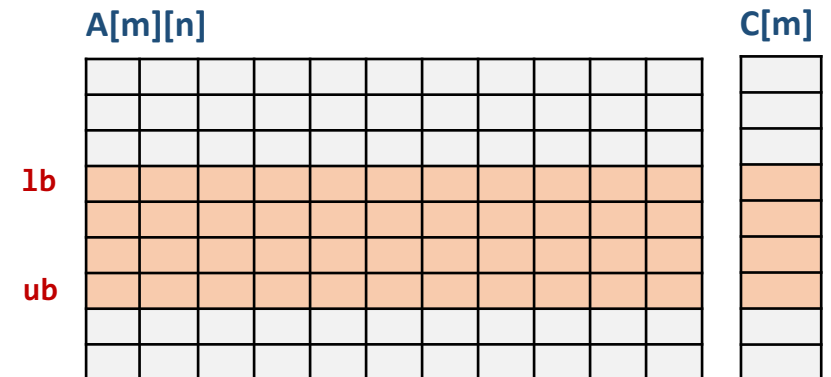
```
for (int i = 0; i < m; i++) {  
    c[i] = 0.0;  
    for (int j = 0; j < n; j++)  
        c[i] += a[i * n + j] * b[j];  
}
```

Распараллеливание внешнего цикла

- Каждому потоку выделяется часть строк матрицы A

DGEMV: параллельная версия

```
/* matrix_vector_product_omp: Compute matrix-vector product  $c[m] = a[m][n] * b[n]$  */  
void matrix_vector_product_omp(double *a, double *b, double *c, int m, int n)  
{  
    #pragma omp parallel  
    {  
        int nthreads = omp_get_num_threads();  
        int threadid = omp_get_thread_num();  
        int items_per_thread = m / nthreads;  
        int lb = threadid * items_per_thread;  
        int ub = (threadid == nthreads - 1) ? (m - 1) : (lb + items_per_thread - 1);  
  
        for (int i = lb; i <= ub; i++) {  
            c[i] = 0.0;  
            for (int j = 0; j < n; j++)  
                c[i] += a[i * n + j] * b[j];  
        }  
    }  
}
```



DGEMV: параллельная версия

```
void run_parallel()
{
    double *a, *b, *c;

    // Allocate memory for 2-d array a[m, n]
    a = xmalloc(sizeof(*a) * m * n);
    b = xmalloc(sizeof(*b) * n);
    c = xmalloc(sizeof(*c) * m);

    for (int i = 0; i < m; i++) {
        for (int j = 0; j < n; j++)
            a[i * n + j] = i + j;
    }
    for (int j = 0; j < n; j++)
        b[j] = j;

    double t = wtime();
    matrix_vector_product_omp(a, b, c, m, n);
    t = wtime() - t;

    printf("Elapsed time (parallel): %.6f sec.\n", t);
    free(a);
    free(b);
    free(c);
}
```


DGEMV: параллельная версия

```
int main(int argc, char **argv)
{
    printf("Matrix-vector product (c[m] = a[m, n] * b[n]; m = %d, n = %d)\n", m, n);
    printf("Memory used: %" PRIu64 " MiB\n", ((m * n + m + n) * sizeof(double)) >> 20);

    run_serial();
    run_parallel();

    return 0;
}
```

Анализ эффективности OpenMP-версии

- Введем обозначения:

- $T(n)$ – время выполнения последовательной программы (serial program) при заданном размере n входных данных
- $T_p(n, p)$ – время выполнения параллельной программы (parallel program) на p процессорах при заданном размере n входных данных

- Коэффициент $S_p(n)$ ускорения параллельной программ (Speedup):

$$S_p(n) = \frac{T(n)}{T_p(n)}$$

Во сколько раз параллельная программа выполняется на p процессорах быстрее последовательной программы при обработке одних и тех же данных размера n

- Как правило

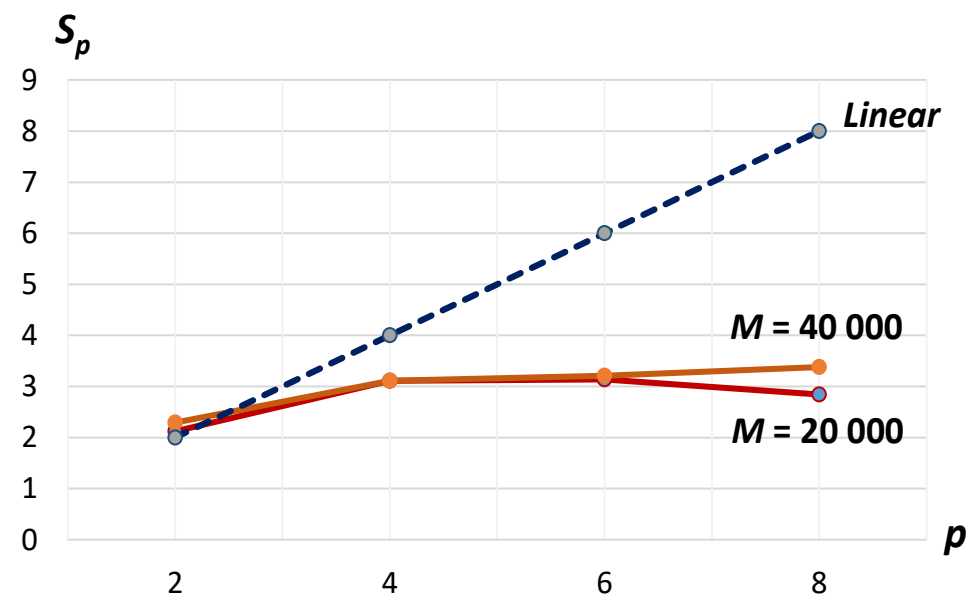
$$S_p(n) \leq p$$

- Цель распараллеливания –

достичь линейного ускорения на наибольшем числе процессоров: $S_p(n) \geq c \cdot p$, при $p \rightarrow \infty$ и $c > 0$

Анализ эффективности OpenMP-версии

M = N	Количество потоков								
	2			4		6		8	
	T_1	T_2	S_2	T_4	S_4	T_6	S_6	T_8	S_8
20 000 (~ 3 GiB)	0.73	0.34	2.12	0.24	3.11	0.23	3.14	0.25	2.84
40 000 (~ 12 GiB)	2.98	1.30	2.29	0.95	3.11	0.91	3.21	0.87	3.38
49 000 (~ 18 GiB)								1.23	3.69



Вычислительный узел кластера Oak (oak.crpt.sibsutis.ru):

- System board: Intel 5520UR
- 8 ядер – два Intel Quad Xeon E5620 (2.4 GHz)
- 24 GiB RAM – 6 x 4GB DDR3 1067 MHz
- CentOS 6.5 x86_64, GCC 4.4.7
- Ключи компиляции: -std=c99 -Wall -O2 -fopenmp

Низкая масштабируемость!

Причины ?

DGEMV: конкуренция за доступ к памяти

```
/* matrix_vector_product_omp: Compute matrix-vector product  $c[m] = a[m][n] * b[n]$  */  
void matrix_vector_product_omp(double *a, double *b, double *c, int m, int n)  
{  
    #pragma omp parallel  
    {  
        int nthreads = omp_get_num_threads();  
        int threadid = omp_get_thread_num();  
        int items_per_thread = m / nthreads;  
        int lb = threadid * items_per_thread;  
        int ub = (threadid == nthreads - 1) ? (m - 1) : (lb + items_per_thread - 1);  
  
        for (int i = lb; i <= ub; i++) {  
            c[i] = 0.0; // Store – запись в память  
            for (int j = 0; j < n; j++)  
                // 4 обращения к памяти: Load c[i], Load a[i][j], Load b[j], Store c[i]  
                // 2 арифметические операции + и *  
                c[i] = c[i] + a[i * n + j] * b[j];  
        }  
    }  
}
```

- DGEMV – *data intensive application*
- Конкуренция за доступ к контролеру памяти
- ALU ядер загружены незначительно

Конфигурация узла кластера Oak

Вычислительный узел кластера Oak (oak.crpt.sibsutis.ru):

- System board: Intel 5520UR (NUMA-система)
- Процессоры связаны шиной **QPI Link**: 5.86 GT/s
- 24 GiB RAM – 6 x 4GB DDR3 1067 MHz

```
$ numactl --hardware
```

```
available: 2 nodes (0-1)
```

```
node 0 cpus: 0 2 4 6
```

```
node 0 size: 12224 MB
```

```
node 0 free: 11443 MB
```

```
node 1 cpus: 1 3 5 7
```

```
node 1 size: 12288 MB
```

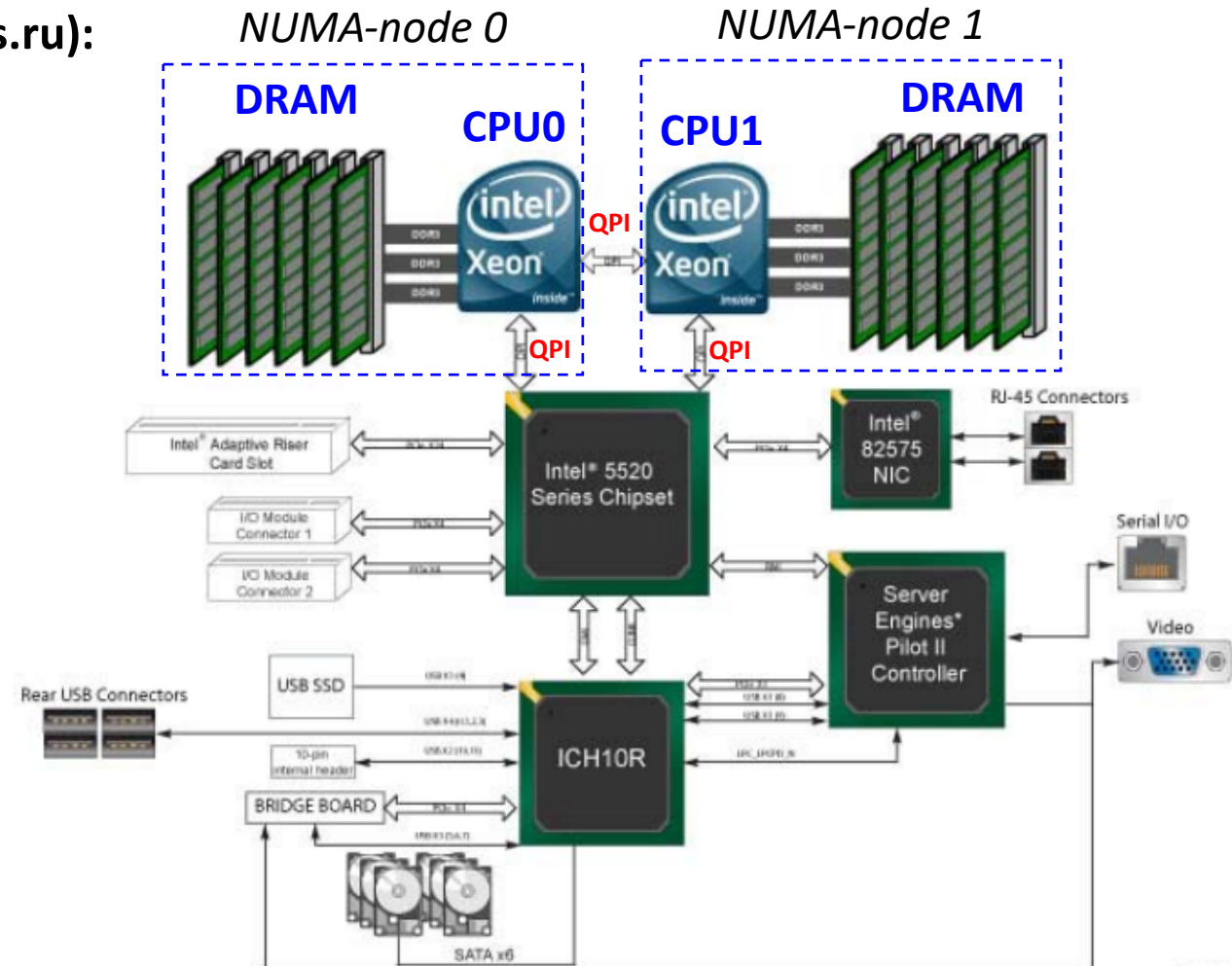
```
node 1 free: 11837 MB
```

```
node distances:
```

```
node  0  1
```

```
 0:  10  21
```

```
 1:  21  10
```



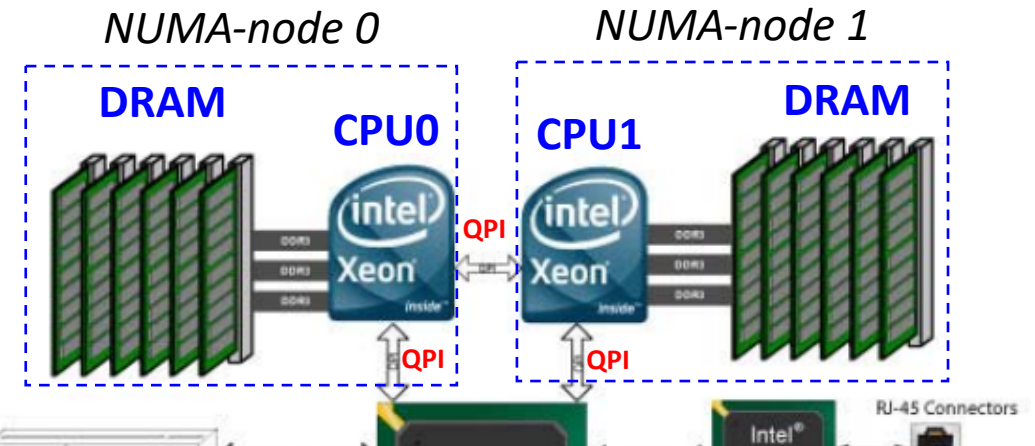
Конфигурация узла кластера Oak

Вычислительный узел кластера Oak (oak.crpt.sibsutis.ru):

- System board: Intel 5520UR (NUMA-система)
- Процессоры связаны шиной **QPI Link**: 5.86 GT/s
- 24 GiB RAM – 6 x 4GB DDR3 1067 MHz

```
$ numactl --hardware
```

```
available: 2 nodes (0-1)
node 0 cpus: 0 2 4 6
node 0 size: 12224 MB
node 0 free: 11443 MB
node 1 cpus: 1 3 5 7
node 1 size: 12288 MB
node 1 free: 11837 MB
node distances:
node  0  1
  0:  10  21
  1:  21  10
```

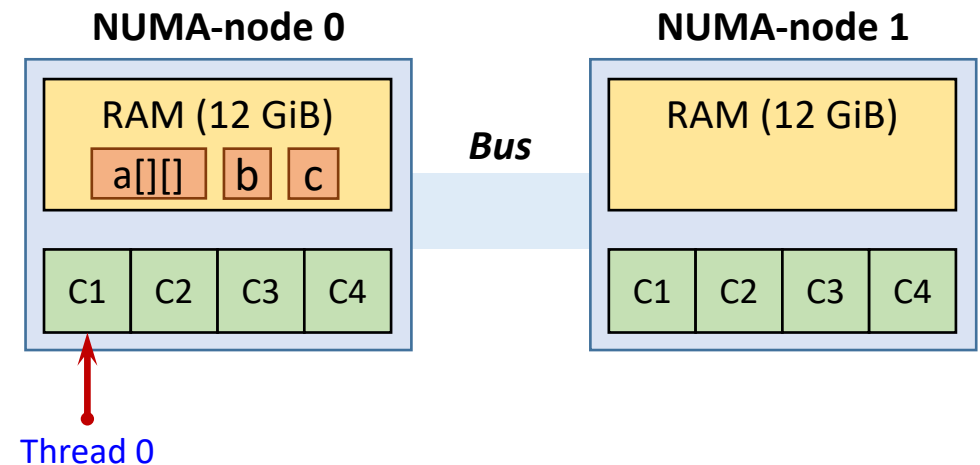


А на каком NUMA-узле (узлах) размещена матрица A и векторы B, C?

Выделение памяти потокам в GNU/Linux

- Страница памяти выделяется с NUMA-узла того потока, который первый к ней обратился (*first-touch policy*)
- Данные желательно инициализировать теми потоками, которые будут с ними работать

```
void run_parallel()  
{  
    double *a, *b, *c;  
  
    // Allocate memory for 2-d array a[m, n]  
    a = xmalloc(sizeof(*a) * m * n);  
    b = xmalloc(sizeof(*b) * n);  
    c = xmalloc(sizeof(*c) * m);  
  
    for (int i = 0; i < m; i++) {  
        for (int j = 0; j < n; j++)  
            a[i * n + j] = i + j;  
    }  
    for (int j = 0; j < n; j++)  
        b[j] = j;  
}
```

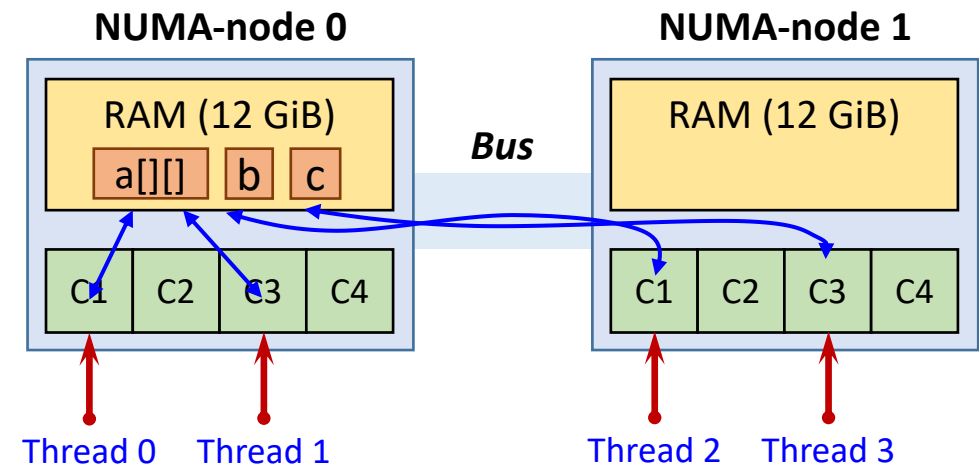


- Поток 0 запрашивает выделение памяти под массивы
- Пока хватает памяти, ядро выделяет страницы с NUMA-узла 0, затем с NUMA-узла 1

Выделение памяти потокам в GNU/Linux

- Страница памяти выделяется с NUMA-узла того потока, который первый к ней обратился (*first-touch policy*)
- Данные желательно инициализировать теми потоками, которые будут с ними работать

```
void run_parallel()  
{  
    double *a, *b, *c;  
  
    // Allocate memory for 2-d array a[m, n]  
    a = xmalloc(sizeof(*a) * m * n);  
    b = xmalloc(sizeof(*b) * n);  
    c = xmalloc(sizeof(*c) * m);  
  
    for (int i = 0; i < m; i++) {  
        for (int j = 0; j < n; j++)  
            a[i * n + j] = i + j;  
    }  
    for (int j = 0; j < n; j++)  
        b[j] = j;  
}
```



- Обращение к массивам из потоков NUMA-узла 1 будет идти через **межпроцессорную шину** в память узла 0

Параллельная инициализация массивов

```
void run_parallel()
{
    double *a, *b, *c;
    // Allocate memory for 2-d array a[m, n]
    a = xmalloc(sizeof(*a) * m * n);
    b = xmalloc(sizeof(*b) * n);
    c = xmalloc(sizeof(*c) * m);

    #pragma omp parallel
    {
        int nthreads = omp_get_num_threads();
        int threadid = omp_get_thread_num();
        int items_per_thread = m / nthreads;
        int lb = threadid * items_per_thread;
        int ub = (threadid == nthreads - 1) ? (m - 1) : (lb + items_per_thread - 1);

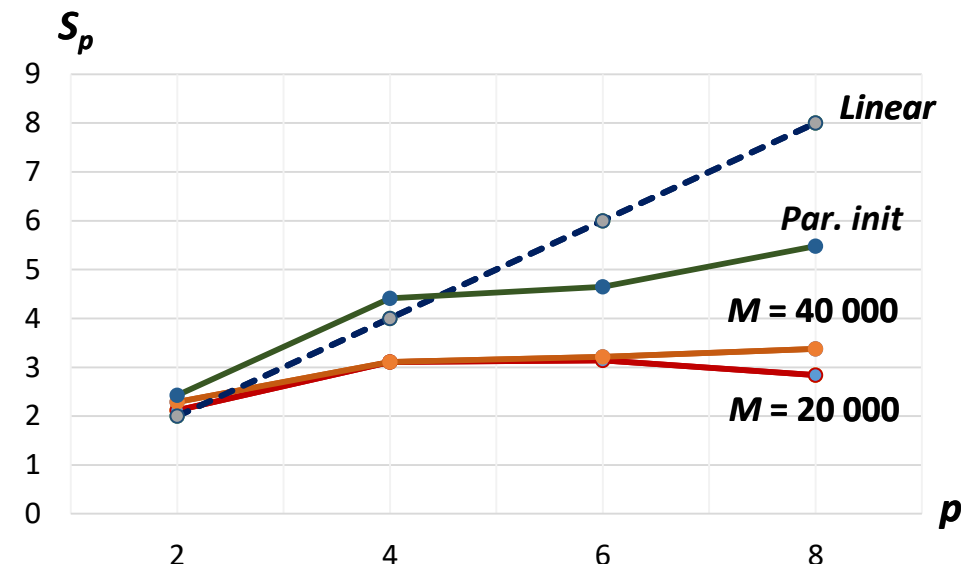
        for (int i = lb; i <= ub; i++) {
            for (int j = 0; j < n; j++)
                a[i * n + j] = i + j;
            c[i] = 0.0;
        }
    }
    for (int j = 0; j < n; j++)
        b[j] = j;

    /* ... */
}
```

Анализ эффективности OpenMP-версии (2)

M = N	Количество потоков									
	2			4		6		8		
	T_1	T_2	S_2	T_4	S_4	T_6	S_6	T_8	S_8	
20 000 (~ 3 GiB)	0.73	0.34	2.12	0.24	3.11	0.23	3.14	0.25	2.84	
40 000 (~ 12 GiB)	2.98	1.30	2.29	0.95	3.11	0.91	3.21	0.87	3.38	
49 000 (~ 18 GiB)								1.23	3.69	

Parallel initialization										
40 000 (~ 12 GiB)	2.98	1.22	2.43	0.67	4.41	0.65	4.65	0.54	5.48	
49 000 (~ 18 GiB)								0.83	5.41	



Улучшили масштабируемость

Дальнейшие оптимизации:

- Эффективный доступ к кеш-памяти
- Векторизация кода (SSE/AVX)
- ...

Суперлинейное ускорение (super-linear speedup): $S_p(n) > p$

Спасибо за внимание!

Время выполнения отдельных потоков

```
void matrix_vector_product_omp(double *a, double *b, double *c, int m, int n)
{
    #pragma omp parallel
    {
        double t = omp_get_wtime();
        int nthreads = omp_get_num_threads();
        int threadid = omp_get_thread_num();
        int items_per_thread = m / nthreads;
        int lb = threadid * items_per_thread;
        int ub = (threadid == nthreads - 1) ? (m - 1) : (lb + items_per_thread - 1);

        for (int i = lb; i <= ub; i++) {
            c[i] = 0.0;
            for (int j = 0; j < n; j++)
                c[i] += a[i * n + j] * b[j];
        }
        t = omp_get_wtime() - t;

        printf("Thread %d items %d [%d - %d], time: %.6f\n", threadid, ub - lb + 1, lb, ub, t);
    }
}
```