

# ПРОСТОЙ ПРИМЕР ПАРАЛЛЕЛЬНОЙ ПРОГРАММЫ

# MPI – MESSAGE PASSING INTERFACE

**Это спецификация для разработчиков и пользователей библиотекой передачи сообщений.**

**1992 – 1994 - первое появление стандарта MPI-1**

**1996 – появление стандарта MPI-2**

**Наиболее известные реализации стандарта MPI:**

**MVAPICH, MVAPICH2, OpenMPI - [smh11.cc.dvo.ru](http://smh11.cc.dvo.ru)**

**MPICH2 - [mvs17.cc.dvo.ru](http://mvs17.cc.dvo.ru)**

# КОМПИЛЯЦИЯ MPI-ПРОГРАММ

|         | MPICH2          | MVAPICH2        | OpenMPI         |
|---------|-----------------|-----------------|-----------------|
| C       | mpicc           | mpicc           | mpicc           |
| C++     | mpic++   mpicxx | mpic++   mpicxx | mpic++   mpicxx |
| Fortran | mpif77   mpif90 | mpif77   mpif90 | mpif77   mpif90 |

**Компилятор MPI** – обертка для компиляторов GNU (gcc, g++), которая линкует соответствующую библиотеку MPI.

*Если в системе присутствует несколько реализаций MPI, то можно выбрать одну из них командой mpi-selector*

```
mpi-selector --list
  mvarich-1.2.0-3635
  mvarich2-1.6
  openmpi-1.4.3
```

# АТРИБУТЫ MPI-ПРОГРАММЫ

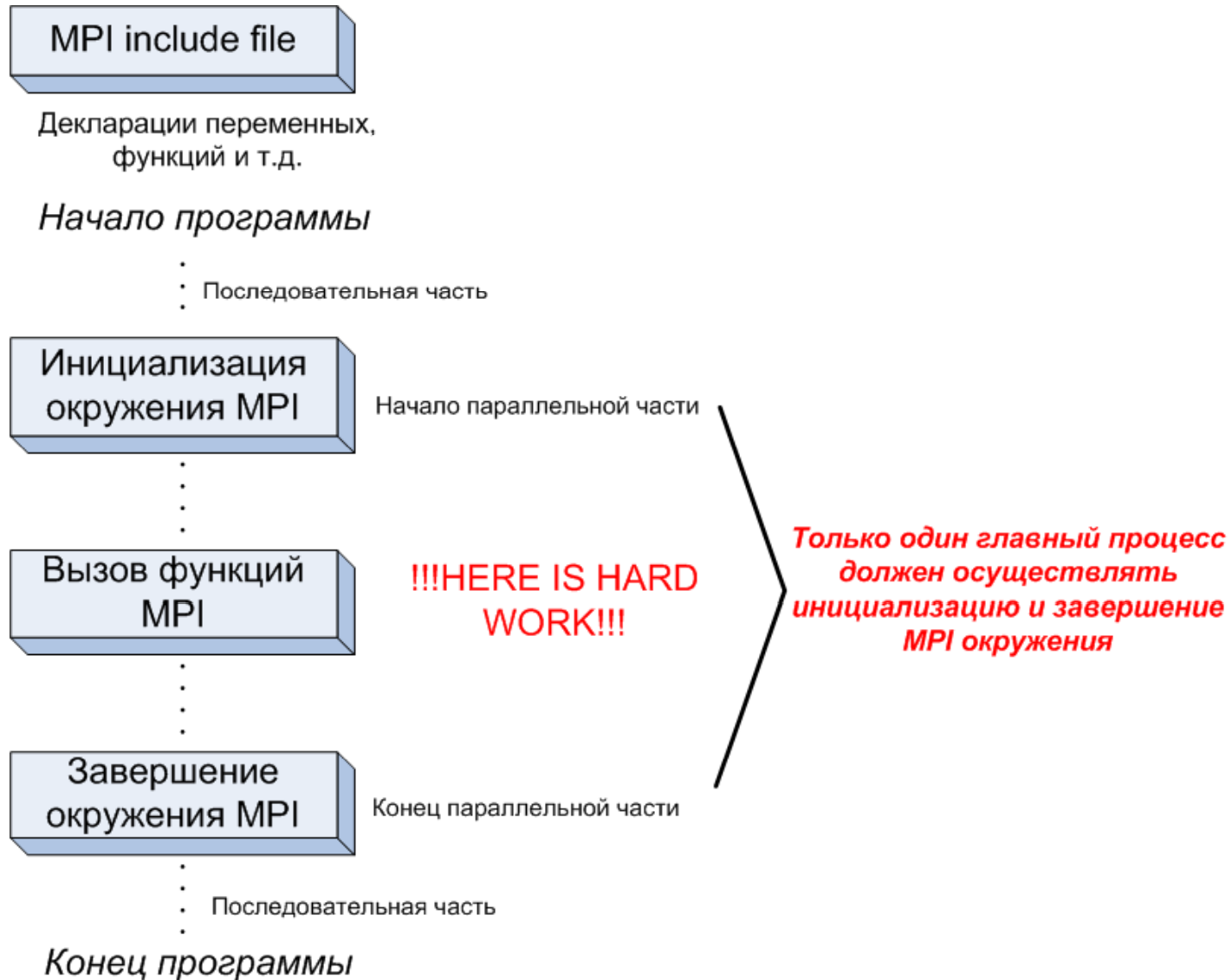
➤ Включение заголовочного файла:

|                    |                      |
|--------------------|----------------------|
| C/C++ include file | Fortran include file |
| #include "mpi.h"   | include 'mpif.h'     |

➤ Формат вызова MPI-функций:

| <b>C/C++</b>       |  |
|--------------------|--|
| <b>Формат:</b>     | rc = MPI_Xxxxx(parameter, ... );   |
| <b>Пример:</b>     | rc = MPI_Bsend(&buf, count, type, dest, tag, comm);                          |
| <b>Код ошибки:</b> | Передается в переменную rc. (MPI_SUCCESS)                                    |
| <b>Fortran</b>     |  |
| <b>Формат:</b>     | CALL MPI_XXXXX(parameter, ..., ierr)<br>call mpi_xxxxx(parameter, ..., ierr) |
| <b>Пример:</b>     | CALL MPI_BSEND(buf, count, type, dest, tag, comm, ierr)                      |
| <b>Код ошибки:</b> | Передается в переменную ierr. (MPI_SUCCESS)                                  |

# ОБЩАЯ СТРУКТУРА MPI-ПРОГРАММЫ



# ПРИМИТИВНЫЕ ТИПЫ ДАННЫХ В MPI

| C data type               |                    | Fortran data type           |                  |
|---------------------------|--------------------|-----------------------------|------------------|
| <b>MPI_CHAR</b>           | signed char        | <b>MPI_CHARACTER</b>        | character(1)     |
| <b>MPI_SHORT</b>          | signed short int   |                             |                  |
| <b>MPI_INT</b>            | signed int         | <b>MPI_INTEGER</b>          | integer          |
| <b>MPI_LONG</b>           | signed long int    |                             |                  |
| <b>MPI_UNSIGNED_CHAR</b>  | unsigned char      |                             |                  |
| <b>MPI_UNSIGNED_SHORT</b> | unsigned short int |                             |                  |
| <b>MPI_UNSIGNED</b>       | unsigned int       |                             |                  |
| <b>MPI_FLOAT</b>          | float              | <b>MPI_REAL</b>             | real             |
| <b>MPI_DOUBLE</b>         | double             | <b>MPI_DOUBLE_PRECISION</b> | double precision |
| <b>MPI_LONG_DOUBLE</b>    | long double        |                             |                  |
|                           |                    | <b>MPI_COMPLEX</b>          | complex          |
|                           |                    | <b>MPI_LOGICAL</b>          | Logical          |
| <b>MPI_BYTE</b>           | 8 binary digits    | <b>MPI_BYTE</b>             | 8 binary digits  |

# НЕКОТОРЫЕ ОСНОВНЫЕ ФУНКЦИИ MPI

|                      |   |   |
|----------------------|---|---|
| <b>MPI_Init</b>      | <b>MPI_Init (&amp;argc,&amp;argv)</b>   | Инициализация среды исполнения MPI-программы. Должна вызываться во <b>ВСЕХ</b> программах MPI и перед всеми другими функциями MPI. Вызывается <b>ТОЛЬКО ОДИН</b> раз. |
|                      | <b>MPI_INIT (ierr)</b>                  |   |
| <b>MPI_Finalize</b>  | <b>MPI_Finalize ()</b>                  | Завершение окружения MPI. Является <b>ПОСЛЕДНЕЙ</b> функцией MPI во <b>ВСЕХ</b> MPI-программах – ни одна MPI-функция не может быть вызвана после нее.                 |
|                      | <b>MPI_FINALIZE (ierr)</b>              |   |
| <b>MPI_Wtime</b>     | <b>MPI_Wtime ()</b>                     | Возвращает текущее время в секундах   |
|                      | <b>MPI_WTIME ()</b>                     |   |
| <b>MPI_Comm_rank</b> | <b>MPI_Comm_rank (comm,&amp;rank)</b>   | Возвращает номер вызывающего процесса (каждый вычислительный процесс внутри своего коммутатора MPI_COMM_WORLD получает номер из диапазона [0..Ncount-1])              |
|                      | <b>MPI_COMM_RANK (comm, rank, ierr)</b> |   |
| <b>MPI_Barrier</b>   | <b>MPI_Barrier (comm)</b>               | Создает барьер синхронизации для группы процессов. Каждый процесс при достижении барьера блокируется пока все процессы из группы не достигнут барьера                 |
|                      | <b>MPI_BARRIER (comm, ierr)</b>         |   |

# ПРИМЕР №1

```
#include <stdio.h>
#include "mpi.h"
int main(int argc, char *argv[]){
    int rank, num_proc, rc, tag = 1;
    double t0, t1, time;
    /*MPI SECTION BEGIN*/
    rc = MPI_Init(&argc, &argv);
    MPI_Status Stat;
    if(rc != MPI_SUCCESS){
        printf("\nERROR initializing MPI. Exit.\n");
        MPI_Abort(MPI_COMM_WORLD, rc);
        return 0; }
    t0 = MPI_Wtime();
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &num_proc);
    if( rank == 0) printf("\nBegin example #1\n");
    if( num_proc == 1){
        printf("\nProcess number equals to 1. Exit.\n");
        MPI_Finalize();
        return 0; }
    if( rank != 0) {
        printf("\nHello from process #%d", rank);
        t1 = MPI_Wtime(); time = 1.e6 * (t1 - t0);
        rc = MPI_Send(&time, 1, MPI_DOUBLE, 0, tag, MPI_COMM_WORLD);}
    else{
        t1 = MPI_Wtime();
        time = 1.e6 * (t1 - t0);
        FILE *time_file = fopen("./times.txt", "wt");
        for( int i = 1; i < num_proc; i++ ){
            double recv_time;
            rc = MPI_Recv(&recv_time, 1, MPI_DOUBLE, i, tag, MPI_COMM_WORLD, &Stat);
            fprintf(time_file, "#%d process: %9.0f msec\n", i, recv_time);}
        fprintf(time_file, "#%d process: %9.0f msec\n", rank, time);
        fclose(time_file); }
    MPI_Barrier(MPI_COMM_WORLD);
    if( rank == 0 ) printf("\nEnd example\n");
    MPI_Finalize();
    /*MPI SECTION END*/
    return 0;}

```

*MPI section*



# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ

- Посылка сообщений, которые содержат значения различных типов (например, целое число с последующим набором вещественных чисел)
- Посылка несмежных данных (например, подблоки матрицы)

**Производный тип MPI** – скрытый объект, который специфицирует две вещи: последовательность базовых типов и последовательность смещений. Последовательность таких пар определяется как *отображение (карта) типа*:  
$$\text{Typemap} = \{(\text{type}_0, \text{disp}_0), \dots, (\text{type}_{n-1}, \text{disp}_{n-1})\}$$

Стандартный сценарий определения и использования производных типов:

- Производный тип строится из predetermined типов MPI и ранее определенных производных типов с помощью специальных функций-конструкторов.
- Новый производный тип регистрируется вызовом функции `MPI_Type_commit`.
- Уничтожается функцией `MPI_Type_free`.

Любой тип данных в MPI имеет две характеристики: протяженность и размер, выраженные в байтах:


*Протяженность типа* определяет, сколько байт переменная данного типа занимает в памяти (опрашивается подпрограммой `MPI_Type_extent`).

*Размер типа* определяет количество реально передаваемых байт в коммуникационных операциях (опрашивается подпрограммой `MPI_Type_size`).


**Для простых типов протяженность и размер совпадают.**

# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ -1

|   |   |
|---|---|
| C   | <code>int MPI_Type_contiguous(int count, MPI_Datatype oldtype, MPI_Datatype *newtype)</code>                                |
| Fortran   | <code>MPI_TYPE_CONTIGUOUS(COUNT, OLDDTYPE, NEWTYPE, IERROR)</code><br><code>INTEGER COUNT, OLDDTYPE, NEWTYPE, IERROR</code> |
| IN  | <code>count</code> - число элементов базового типа  |
| IN  | <code>oldtype</code> - базовый тип данных   |
| OUT   | <code>newtype</code> - новый производный тип данных   |
| создает новый тип, элементы которого состоят из указанного числа элементов базового типа, занимающих смежные области памяти |   |

`oldtype = MPI_REAL` 

`count = 4`

`newtype` 

# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ -2

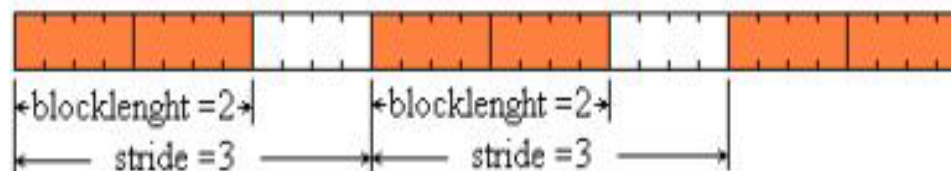
|   |  |
|---|--|
| C   | <code>int MPI_Type_vector(int count, int blocklength, int stride, MPI_Datatype oldtype, MPI_Datatype *newtype)</code>                              |
| Fortran   | <code>MPI_TYPE_VECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR)</code><br>INTEGER COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR |
| IN  | count - число блоков   |
| IN  | blocklength - число элементов базового типа в каждом блоке   |
| IN  | stride - шаг между началами соседних блоков, измеренный числом элементов базового типа   |
| IN  | oldtype - базовый тип данных   |
| OUT   | newtype - новый производный тип данных   |
| создает тип, элемент которого представляет собой несколько равноудаленных друг от друга блоков из одинакового числа смежных элементов базового типа |  |

oldtype = MPI\_REAL



count = 3, blocklength = 2, stride = 3

newtype



# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ -3

|  |  |
|--|--|
| C  | int MPI_Type_indexed(int count, int *array_of_blocklengths,int *array_of_displacements, MPI_Datatype oldtype, MPI_Datatype *newtype)   |
| Fortran  | MPI_TYPE_INDEXED(COUNT, ARRAY_OF_BLOCKLENGTHS,ARRAY_OF_DISPLACEMENTS, OLDTYPE, NEWTYPE, IERROR)INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*),OLDTYPE, NEWTYPE, IERROR |
| IN   | count - число блоков   |
| IN   | array_of_blocklengths - массив, содержащий число элементов базового типа в каждом блоке  |
| IN   | array_of_displacements - массив смещений каждого блока от начала размещения элемента нового типа, смещения измеряются числом элементов базового типа                                       |
| IN   | oldtype - базовый тип данных   |
| OUT  | newtype - новый производный тип данных   |
| создает тип, элемент которого состоит из произвольных по длине блоков с произвольным смещением блоков от начала размещения элемента. Смещения измеряются в элементах старого типа. |  |

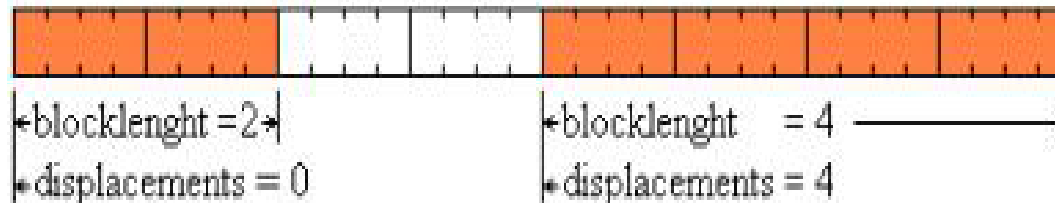
# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ -3

Эта функция создает тип `newtype`, каждый элемент которого состоит из `count` блоков, где  $i$ -ый блок содержит `array_of_blocklengths[i]` элементов базового типа и смещен от начала размещения элемента нового типа на `array_of_displacements[i]` элементов базового типа.

`oldtype = MPI_REAL` 

`count = 2` `blocklength = (2, 4)` `displacements = (0, 4)`

`newtype`



# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ -4

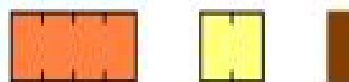
|  |   |
|--|---|
| C  | int MPI_Type_create_struct(int count, int array_of_blocklengths, MPI_Aint *array_of_displacements, MPI_Datatype *array_of_types, MPI_Datatype *newtype)   |
| Fortran  | MPI_TYPE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS, ARRAY_OF_TYPES, NEWTYPE, IERROR) INTEGER COUNT, ARRAY_OF_BLOCKLENGTHS(*), ARRAY_OF_DISPLACEMENTS(*), ARRAY_OF_TYPES(*), NEWTYPE, IERROR |
| IN   | count - число блоков  |
| IN   | array_of_blocklengths - массив, содержащий число элементов одного из базовых типов в каждом блоке   |
| IN   | array_of_displacements - массив смещений каждого блока от начала размещения структуры, смещения измеряются в байтах   |
| IN   | array_of_type - массив, содержащий тип элементов в каждом блоке   |
| OUT  | newtype - новый производный тип данных  |
| создает тип, являющийся структурой, состоящей из произвольного числа блоков, каждый из которых может содержать произвольное число элементов одного из базовых типов и может быть смещен на произвольное число байтов от начала размещения структуры. |   |

# ПРОИЗВОДНЫЕ ТИПЫ ДАННЫХ -4

Функция создает тип `newtype`, элемент которого состоит из `count` блоков, где  $i$ -ый блок содержит `array_of_blocklengths[i]` элементов типа `array_of_types[i]`.

Смещение  $i$ -ого блока от начала размещения элемента нового типа измеряется в байтах и задается в `array_of_displacements[i]`.

`oldtypes = (MPI_INT, MPI_SHORT, MPI_CHAR)`



`count = 3` `blocklength = (1, 6, 4)` `displacements = (0, 12, 26)`

`newtype`



`*blocklength = 1`  
`*displacements = 0`

`*blocklength = 6`  
`*displacements = 12`

`*blocklength = 4`  
`*displacements = 26`

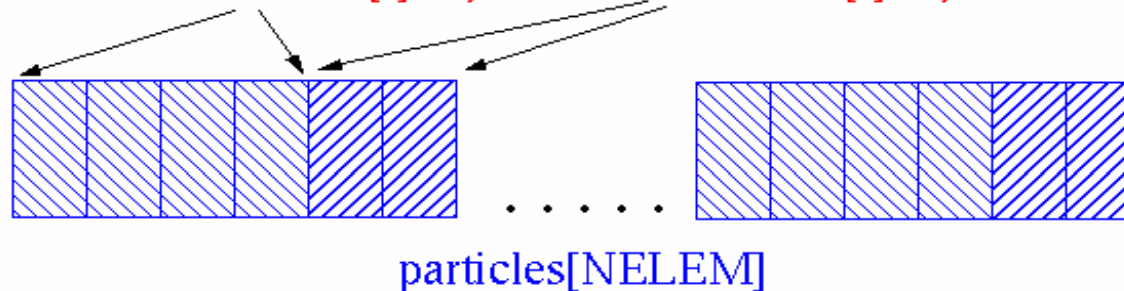
# ПРИМЕР ПРОГРАММЫ С ИСПОЛЬЗОВАНИЕМ ПРОИЗВОДНЫХ ТИПОВ ДАННЫХ

## MPI\_Type\_struct

```
typedef struct { float x, y, z, velocity; int n, type; } Particle;  
Particle particles[NELEM];
```

```
MPI_Type_extent(MPI_FLOAT, &extent);
```

```
count = 2; oldtypes[0] = MPI_FLOAT; oldtypes[1] = MPI_INT  
offsets[0] = 0; offsets[1] = 4 * extent;  
blockcounts[0] = 4; blockcounts[1] = 2;
```



```
MPI_Type_struct(count, blockcounts, offsets, oldtypes, &particletype);
```

```
MPI_Send(particles, NELEM, particletype, dest, tag, comm);
```

Sends entire (NELEM) array of particles, each particle being comprised four floats and two integers.



# ПОСТАНОВКА ЗАДАЧИ В ОЧЕРЕДЬ MVS17.CC.DVO.RU

**TORQUE OpenPBS + MAUI**

## *OpenPBS*

|                                     |  |
|-------------------------------------|--|
| <b>qsub {файл паспорта задания}</b> | Запуск задания                             |
| <b>qstat</b>                        | Проверка статуса задания                   |
| <b>qdel {номер задачи}</b>          | Остановка задачи                           |
| <b>qstat -f -Q</b>                  | Перечень доступных очередей с их описанием |
| <b>tracejob &lt;JOBID&gt;</b>       | Выводит служебную информацию по задаче     |

## *MAUI*

|                               |  |
|-------------------------------|--|
| <b>showq</b>                  | Отражает состояние очереди задач         |
| <b>checkjob &lt;JOBID&gt;</b> | Выводит информацию о задаче              |
| <b>showbf</b>                 | Показывает количество доступных ресурсов |

# ПРИМЕР ПАСПОРТА ЗАДАЧИ ДЛЯ OpenBPS

```
#!/bin/bash
```

```
#PBS -V
```

```
#PBS -S /bin/bash
```

```
#PBS -l walltime=00:15:00
```

```
#PBS -l nodes=1:ppn=4
```

```
#PBS -q simple
```

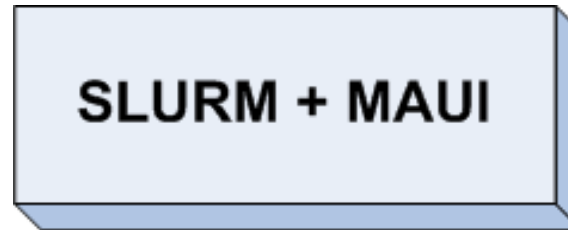
```
#PBS -e /home/daria/Lecture/ex1.err
```

```
#PBS -o /home/daria/Lecture/ex1.out
```

```
cd $PBS_O_WORKDIR
```

```
mpirun -np 8 -machinefile $PBS_NODEFILE ./ex1
```

# ПОСТАНОВКА ЗАДАЧИ В ОЧЕРЕДЬ SMH11.CC.DVO.RU



## *SLURM*

|   |                  |
|---|------------------|
| <b>srun -n X ExecutableFile [optional params]</b>   | Запуск задания   |
| <b>submit -n X ExecutableFile [optional params]</b> | Запуск задания   |
| <b>squeue</b>                                       | Проверка очереди |
| <b>scancel &lt;JOBID&gt;</b>                        | Остановка задачи |

## *MAUI*

|                               |  |
|-------------------------------|--|
| <b>showq</b>                  | Отражает состояние очереди задач         |
| <b>checkjob &lt;JOBID&gt;</b> | Выводит информацию о задаче              |
| <b>showbf</b>                 | Показывает количество доступных ресурсов |