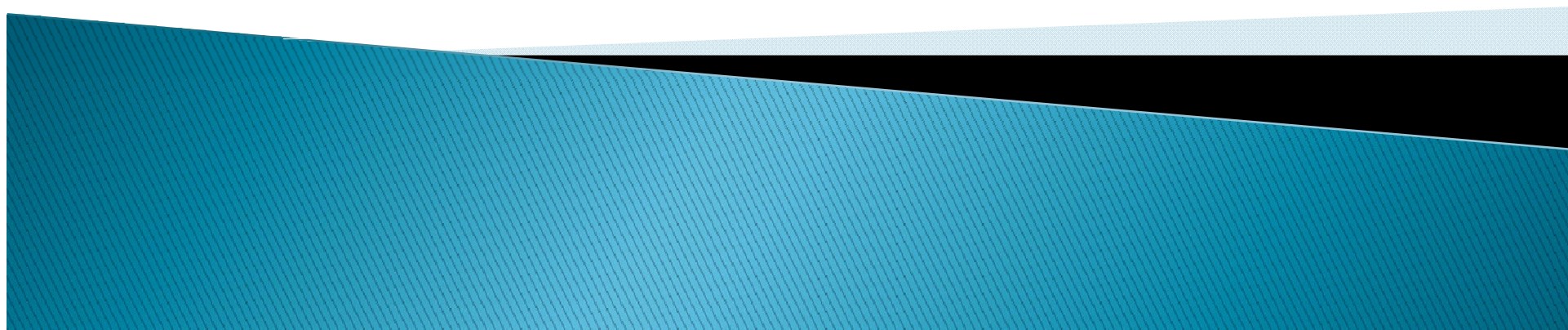
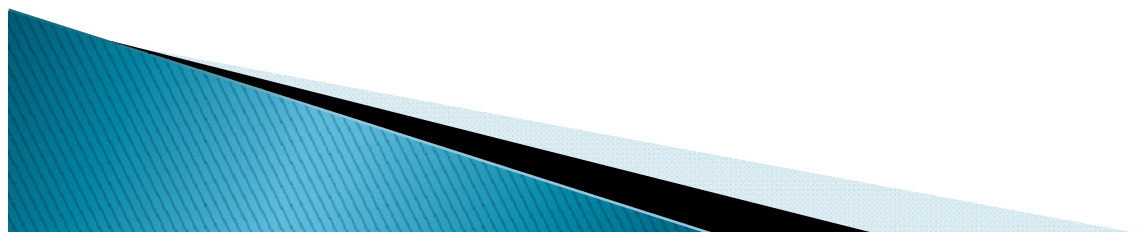


Обзор технологии параллельного программирования MPI



План лекции

- ▶ Стандарт MPI
- ▶ Основные понятия
- ▶ Блокирующие двухточечные обмены
- ▶ Двухточечные обмены с буферизацией, другие типы двухточечных обменов
- ▶ Коллективные операции



TOP-500



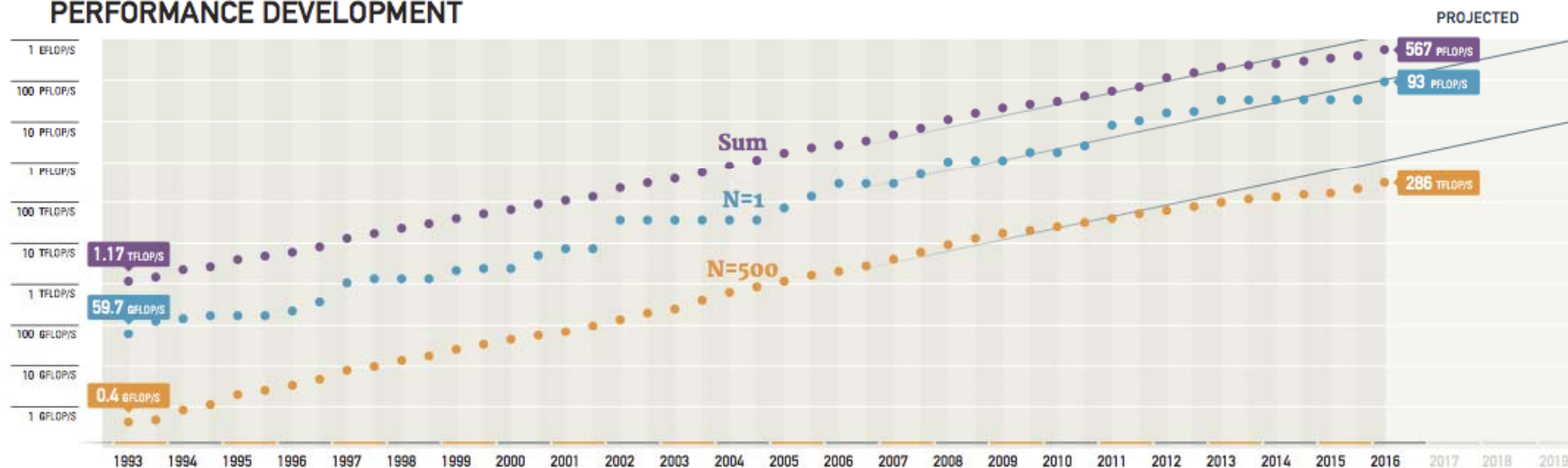
PRESENTED BY



FIND OUT MORE AT
top500.org

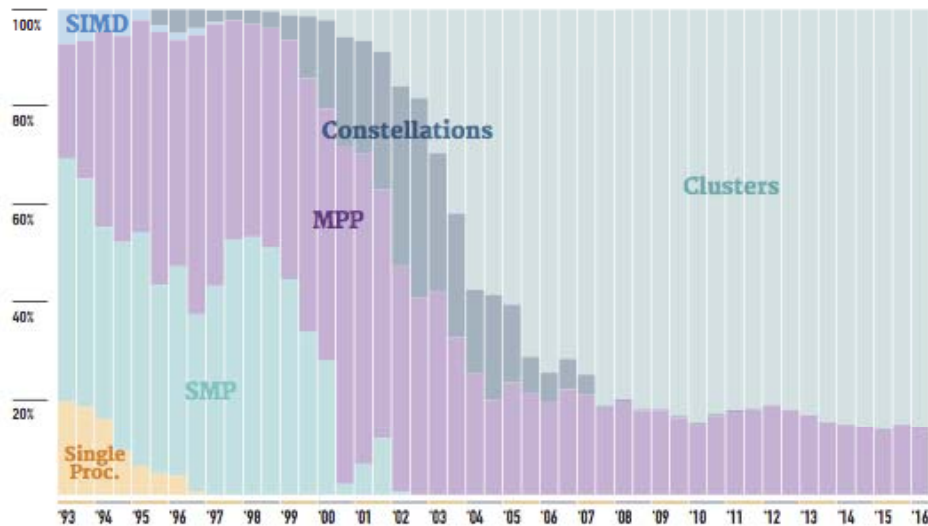
| | NAME | SPECS | SITE | COUNTRY | CORES | R _{MAX} PFLOP/S | POWER MW |
|---|-----------------------|---|-------------------|---------|------------|-----------------------------|-------------|
| 1 | Sunway TaihuLight | Shenwei SW26010 (260C 1.45 GHz) Custom interconnect | NSCC in Wuxi | China | 10,649,600 | 93.0 | 15.4 |
| 2 | Tianhe-2 (Milkyway-2) | Intel Ivy Bridge (12C 2.2 GHz) & Xeon Phi (57C 1.1 GHz), Custom interconnect | NSCC in Guangzhou | China | 3,120,000 | 33.9 | 17.8 |
| 3 | Titan | Cray XK7, Opteron 6274 (16C 2.2 GHz) + Nvidia Kepler GPU, Custom interconnect | DOE/SC/ORNL | USA | 560,640 | 17.6 | 8.2 |
| 4 | Sequoia | IBM BlueGene/Q, Power BQC (16C 1.60 GHz), Custom interconnect | DOE/NNSA/LLNL | USA | 1,572,864 | 17.2 | 7.9 |
| 5 | K computer | Fujitsu SPARC64 VIIIfx (8C 2.0 GHz), Custom interconnect | RIKEN AICS | Japan | 705,024 | 10.5 | 12.7 |

PERFORMANCE DEVELOPMENT

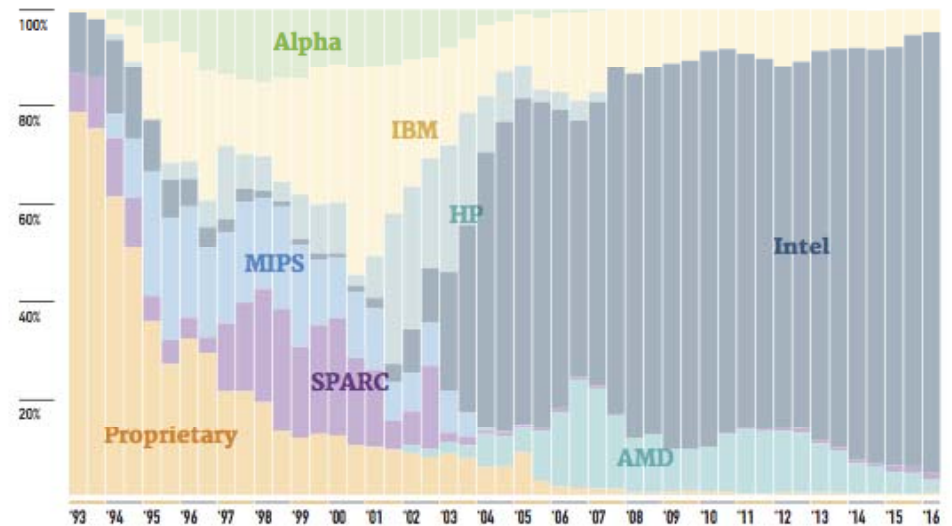


TOP-500

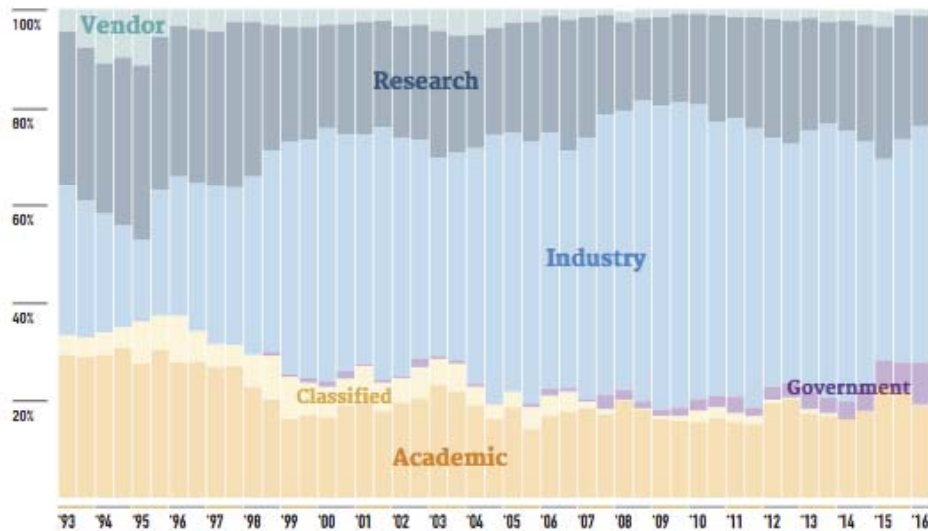
ARCHITECTURES



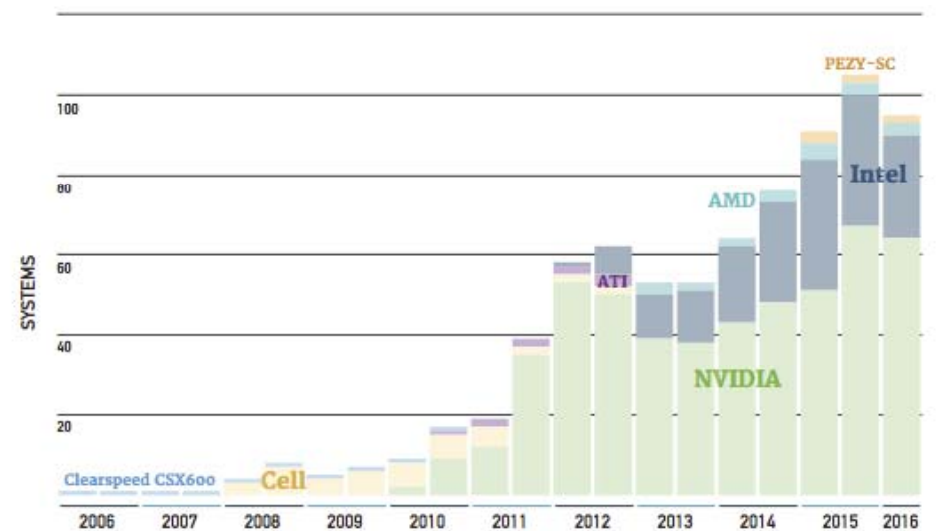
CHIP TECHNOLOGY



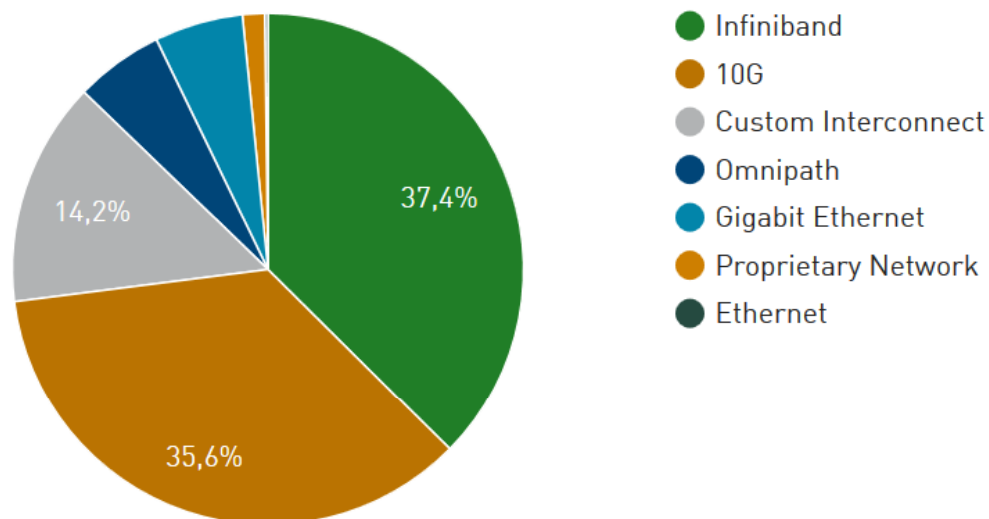
INSTALLATION TYPE



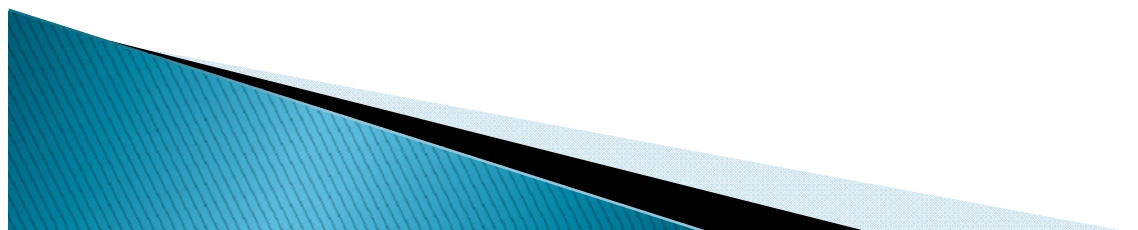
ACCELERATORS/CO-PROCESSORS



TOP-500



| Interconnect Family | Count | System Share (%) | Rmax (GFlops) | Rpeak (GFlops) | Cores |
|---------------------|-------|------------------|---------------|----------------|------------|
| Infiniband | 187 | 37,4 | 183,482,398 | 249,155,816 | 9,129,010 |
| 10G | 178 | 35,6 | 97,644,319 | 195,222,197 | 6,031,044 |
| Custom Interconnect | 71 | 14,2 | 320,842,652 | 444,875,261 | 24,288,812 |
| Omnipath | 28 | 5,6 | 43,672,750 | 64,890,729 | 1,617,076 |
| Gigabit Ethernet | 28 | 5,6 | 12,097,159 | 42,604,816 | 2,308,404 |
| Proprietary Network | 7 | 1,4 | 13,759,900 | 17,574,061 | 594,000 |
| Ethernet | 1 | 0,2 | 613,200 | 920,000 | 25,000 |



MPI

- ▶ MPI 1.1 Standard разрабатывался 92–94
- ▶ MPI 2.0 – 95–97
- ▶ MPI 2.1 – 2008 сентябрь 2008 г.
- ▶ MPI 3.0 – сентябрь 2012 г.
- ▶ MPI 3.1 – июнь 2015 г.

Более 450 процедур

- ▶ Стандарты

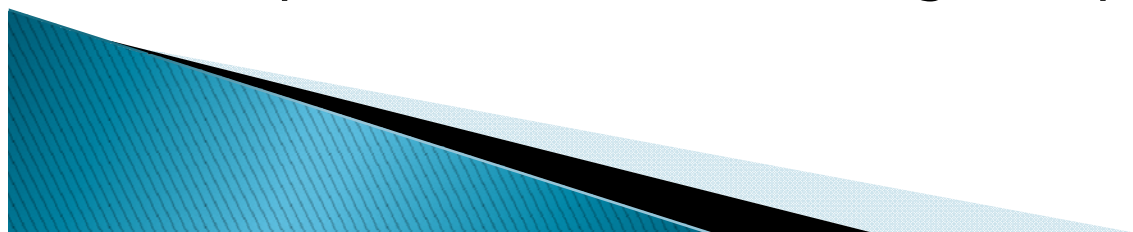
<http://www.mcs.anl.gov/mpi>

<http://www.mpi-forum.org/docs/docs.html>

<https://computing.llnl.gov/tutorials/mpi/>

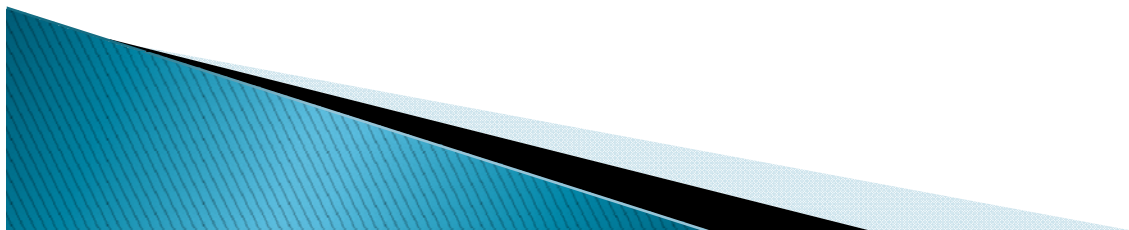
- ▶ Описание функций

<http://www-unix.mcs.anl.gov/mpi/www/>



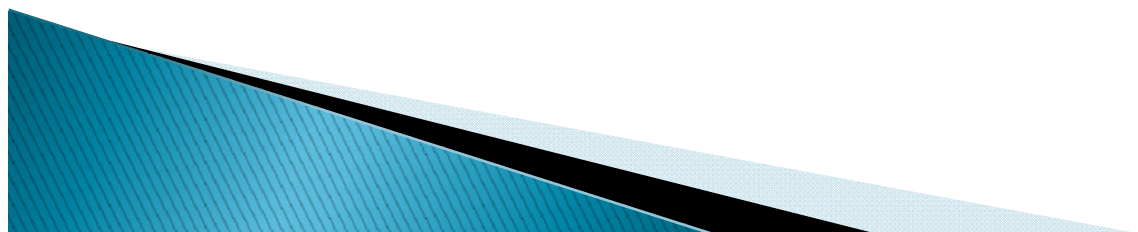
Цель MPI

- ▶ Основная цель:
 - Обеспечение переносимости исходных кодов
 - Эффективная реализация
- ▶ Кроме того:
 - Большая функциональность
 - Поддержка неоднородных параллельных архитектур



Реализации MPI

- ▶ MPICH
- ▶ LAM/MPI
- ▶ Mvarich
- ▶ OpenMPI
- ▶ Коммерческие реализации Intel, IBM и др.



Модель MPI

- ▶ Параллельная программа состоит из процессов, процессы могут быть многопоточными.
- ▶ MPI реализует передачу сообщений между процессами.
- ▶ Межпроцессное взаимодействие предполагает:
 - синхронизацию
 - перемещение данных из адресного пространства одного процесса в адресное пространство другого процесса.



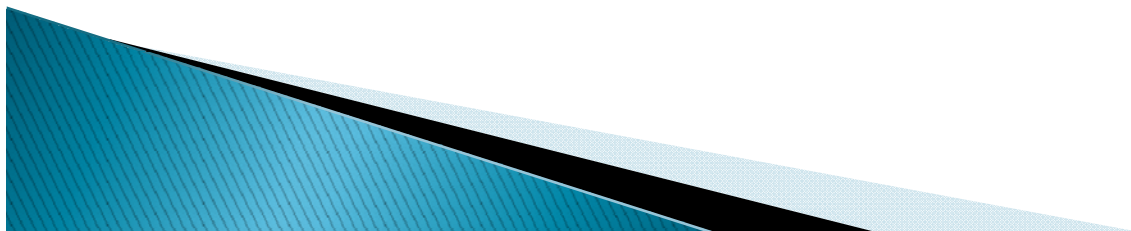
Основные понятия

- ▶ Процессы объединяются в группы.
- ▶ Каждое сообщение посылается в рамках некоторого контекста и должно быть получено в том же контексте.
- ▶ Группа и контекст вместе определяют коммуникатор.
- ▶ Процесс идентифицируется своим номером в группе, ассоциированной с коммуникатором.
- ▶ Коммуникатор, содержащий все начальные процессы, называется MPI_COMM_WORLD.



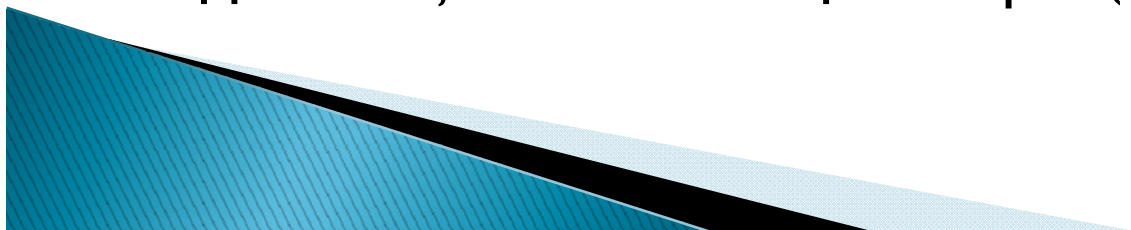
Понятие коммуникатора MPI

- ▶ Управляющий объект, представляющий группу процессов, которые могут взаимодействовать друг с другом
- ▶ Все обращения к MPI функциям содержат коммуникатор, как параметр
- ▶ Наиболее часто используемый коммуникатор `MPI_COMM_WORLD`
- ▶ Определяется при вызове `MPI_Init`
- ▶ Содержит ВСЕ процессы программы



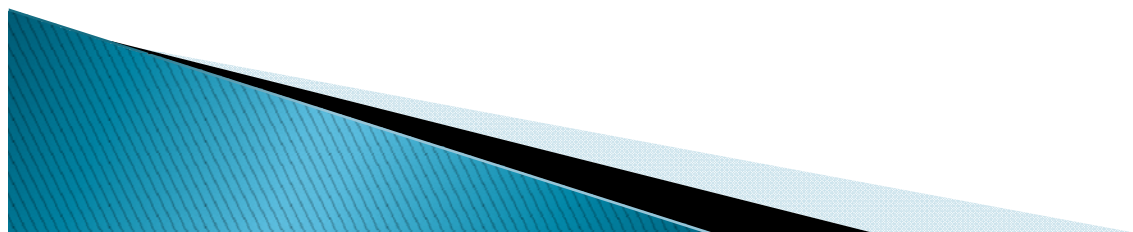
Типы данных MPI

- ▶ Данные в сообщении описываются тройкой: (address, count, datatype), где datatype определяется рекурсивно как:
 - Предопределенный базовый тип, соответствующий типу данных в базовом языке (например, MPI_INT, MPI_DOUBLE_PRECISION)
 - Непрерывный массив MPI типов
 - Векторный тип
 - Индексированный тип
 - Произвольные структуры
- ▶ MPI включает функции для построения пользовательских типов данных, например, типа данных, описывающих пары (int, float).



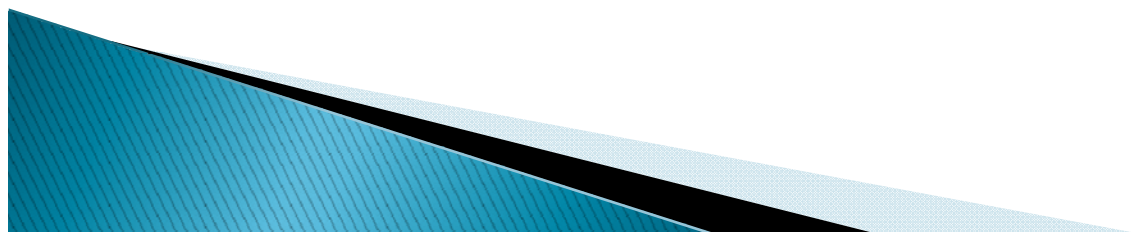
Базовые типы данных

| MPI datatype | C datatype |
|--------------------|--------------------|
| MPI_CHAR | signed char |
| MPI_SHORT | signed short int |
| MPI_INT | signed int |
| MPI_LONG | signed long int |
| MPI_UNSIGNED_CHAR | unsigned char |
| MPI_UNSIGNED_SHORT | unsigned short int |
| MPI_UNSIGNED | unsigned int |
| MPI_UNSIGNED_LONG | unsigned long int |
| MPI_FLOAT | float |
| MPI_DOUBLE | double |
| MPI_LONG_DOUBLE | long double |



Понятие тега

- ▶ Сообщение сопровождается определяемым пользователем признаком для идентификации принимаемого сообщения.
- ▶ Теги сообщений у отправителя и получателя должны быть согласованы.
- ▶ Можно указать в качестве значения тэга константу `MPI_ANY_TAG`.
- ▶ Некоторые не-MPI системы передачи сообщений называют тэг типом сообщения.
- ▶ MPI вводит понятие тэга, чтобы не путать это понятие с типом данных MPI.

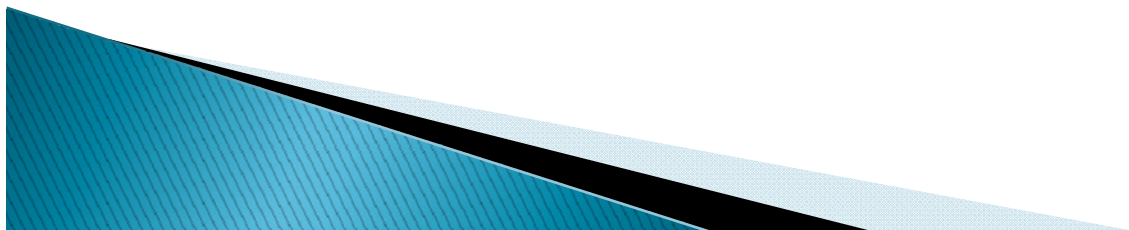


Формат MPI-функций

```
error = MPI_Xxxxx(parameter,...);  
MPI_Xxxxx(parameter,...);
```

- ▶ Возвращаемое значение – код ошибки. Определяется константой MPI_SUCCESS

```
int error;  
.....  
error = MPI_Init(&argc, &argv);  
if (error != MPI_SUCCESS)  
{  
    fprintf (stderr, " MPI_Init error \n");  
    return 1;  
}
```



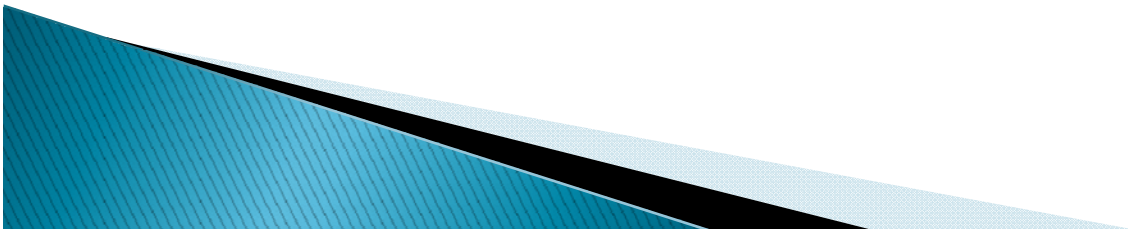
Выполнение MPI-программы

- ▶ При запуске указываем число требуемых процессоров `np` и название программы
`mpirun -np 3 prog`
- ▶ На выделенных для расчета узлах запускается `np` копий указанной программы
- ▶ Каждая копия программы получает два значения:
 - `np`
 - `rank` из диапазона `[0 ... np-1]`
- ▶ Любые две копии программы могут непосредственно обмениваться данными с помощью функций передачи сообщений



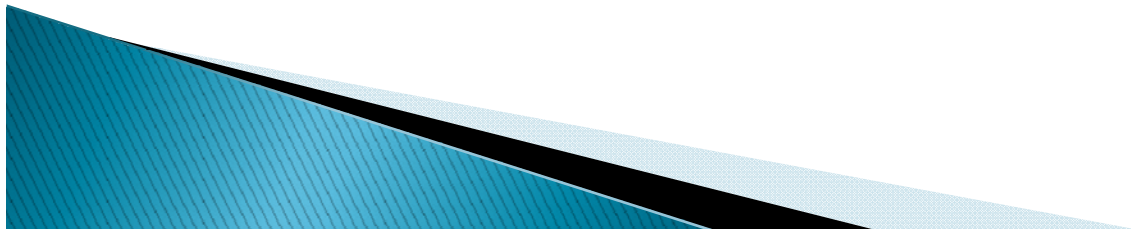
C: MPI helloworld.c

```
#include <mpi.h>
main(int argc, char **argv)
{
    int numtasks, rank;
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &
numtasks);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    printf("Hello World from process %d of %d\n",
rank, numtasks);
    MPI_Finalize();
}
```



Функции определения среды

- ▶ `int MPI_Init(int *argc, char ***argv)`
должна первым вызовом, вызывается только один раз
- ▶ `int MPI_Comm_size(MPI_Comm comm, int *size)`
число процессов в коммуникаторе
- ▶ `int MPI_Comm_rank(MPI_Comm comm, int *rank)`
номер процесса в коммуникаторе (нумерация с 0)
- ▶ `int MPI_Finalize()`
завершает работу процесса
- ▶ `int MPI_Abort (MPI_Comm comm, int*errorcode)`
завершает работу программы



Функции определения среды

- ▶ `int MPI_Initialized(int *flag)`

В аргументе `flag` возвращает 1, если вызвана после процедуры `MPI_Init`, и 0 в противном случае.

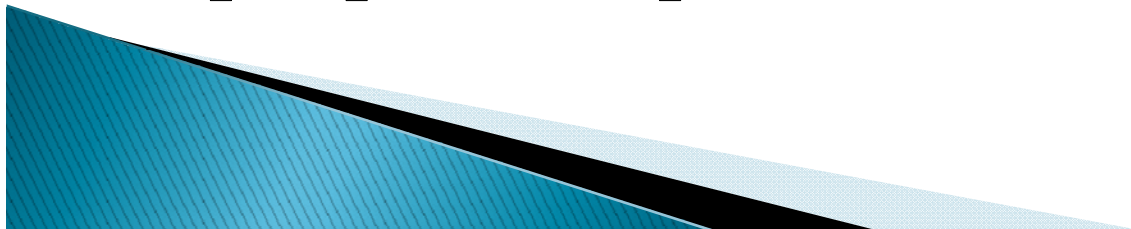
- ▶ `int MPI_Finalized(int *flag)`

В аргументе `flag` возвращает 1, если вызвана после процедуры `MPI_Finalize`, и 0 в противном случае.

Эти процедуры можно вызвать до `MPI_Init` и после `MPI_Finalize`.

- ▶ `int MPI_Get_processor_name(char *name, int *len)`

Возвращает в строке `name` имя узла, на котором запущен вызвавший процесс. В переменной `len` возвращается количество символов в имени, не превышающее константы `MPI_MAX_PROCESSOR_NAME`.



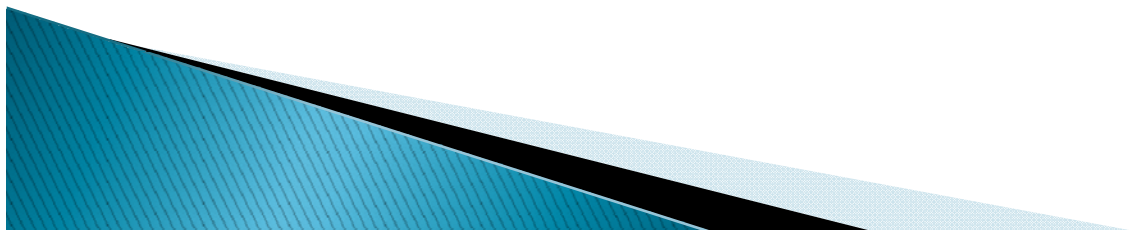
Функции работы с таймером

- ▶ **double MPI_Wtime(void)**

Возвращает для каждого вызвавшего процесса астрономическое время в секундах (вещественное число двойной точности), прошедшее с некоторого момента в прошлом. Момент времени, используемый в качестве точки отсчёта, не будет изменён за время существования процесса.

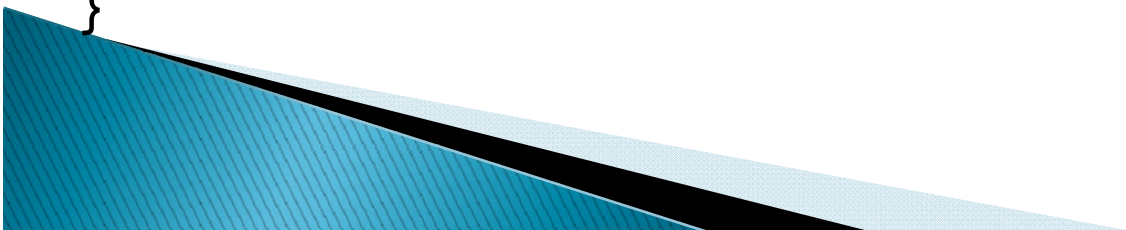
- ▶ **double MPI_Wtick(void)**

Возвращает разрешение таймера в секундах.



Использование таймера

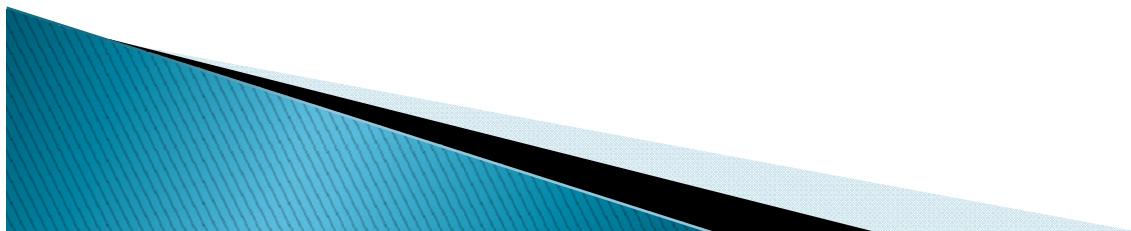
```
#include <stdio.h>
#include "mpi.h"
#define NTIMES 1000
int main(int argc, char **argv)
{
    double time_start, time_finish, tick;
    int rank, i;
    int len;
    char *name;
    name = (char*)malloc(MPI_MAX_PROCESSOR_NAME*sizeof(char));
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Get_processor_name(name, &len);
    tick = MPI_Wtick();
    time_start = MPI_Wtime();
    for (i = 0; i<NTIMES; i++)time_finish = MPI_Wtime();
    printf ("node %s, process %d: tick= %lf, time= %lf\n",
           name, rank, tick, (time_finish-time_start)/NTIMES);
    MPI_Finalize();
}
```



Информация о статусе сообщения

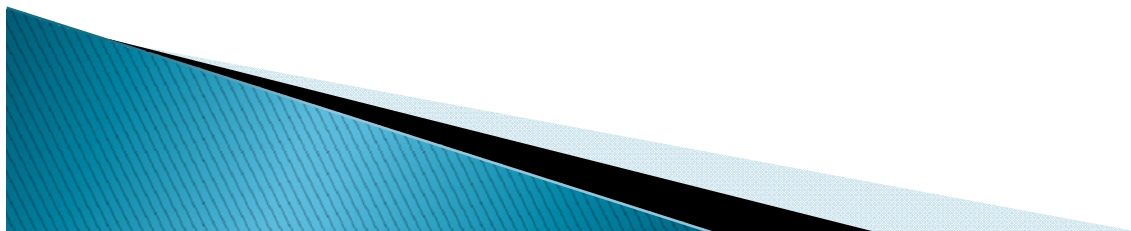
Содержит:

- ▶ Source: `status.MPI_SOURCE`
- ▶ Tag: `status.MPI_TAG`
- ▶ Код ошибки: `status.MPI_ERROR`
- ▶ Count: `MPI_Get_count`

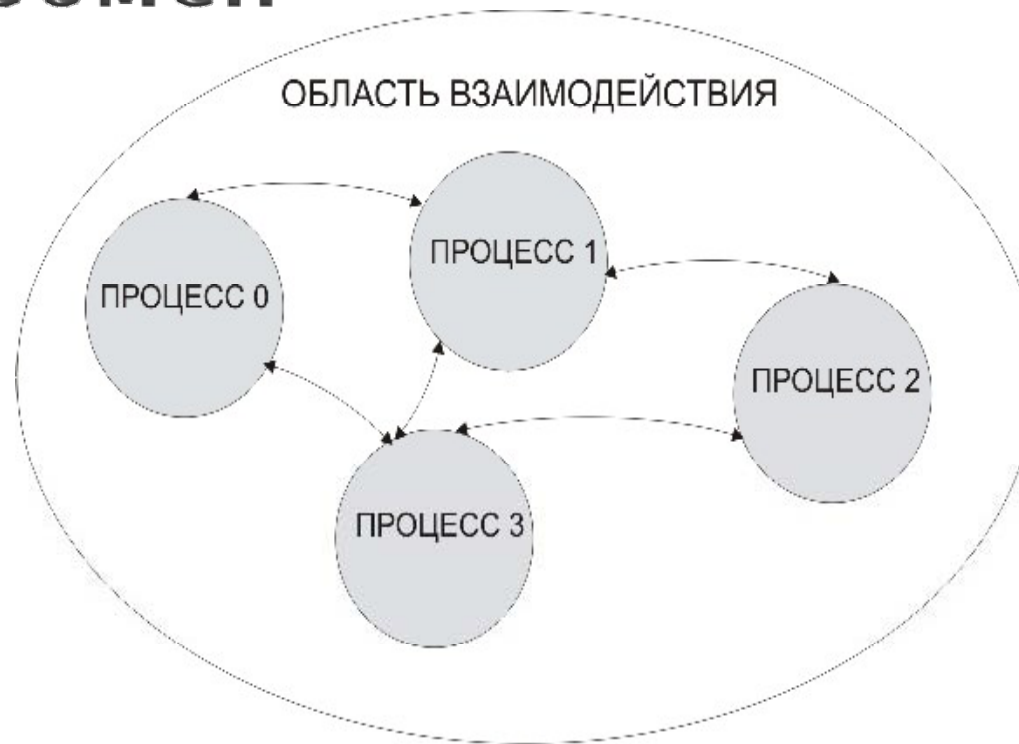


Двухточечный (point-to-point, p2p) обмен

- ▶ В двухточечном обмене участвуют только два процесса, процесс-отправитель и процесс-получатель (источник сообщения и адресат).
- ▶ Двухточечные обмены используются для организации локальных и неструктурированных коммуникаций.



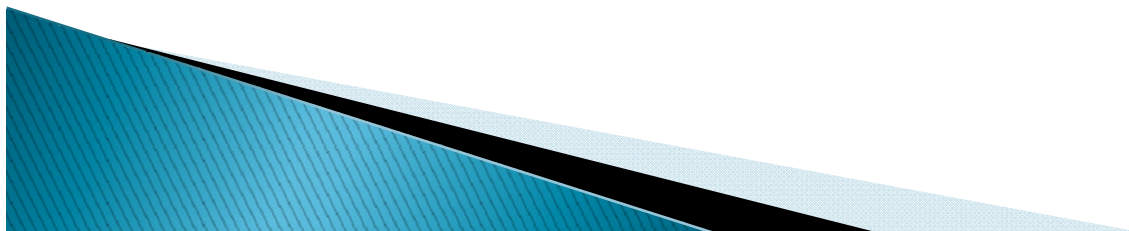
Двухточечный (point-to-point, p2p) обмен



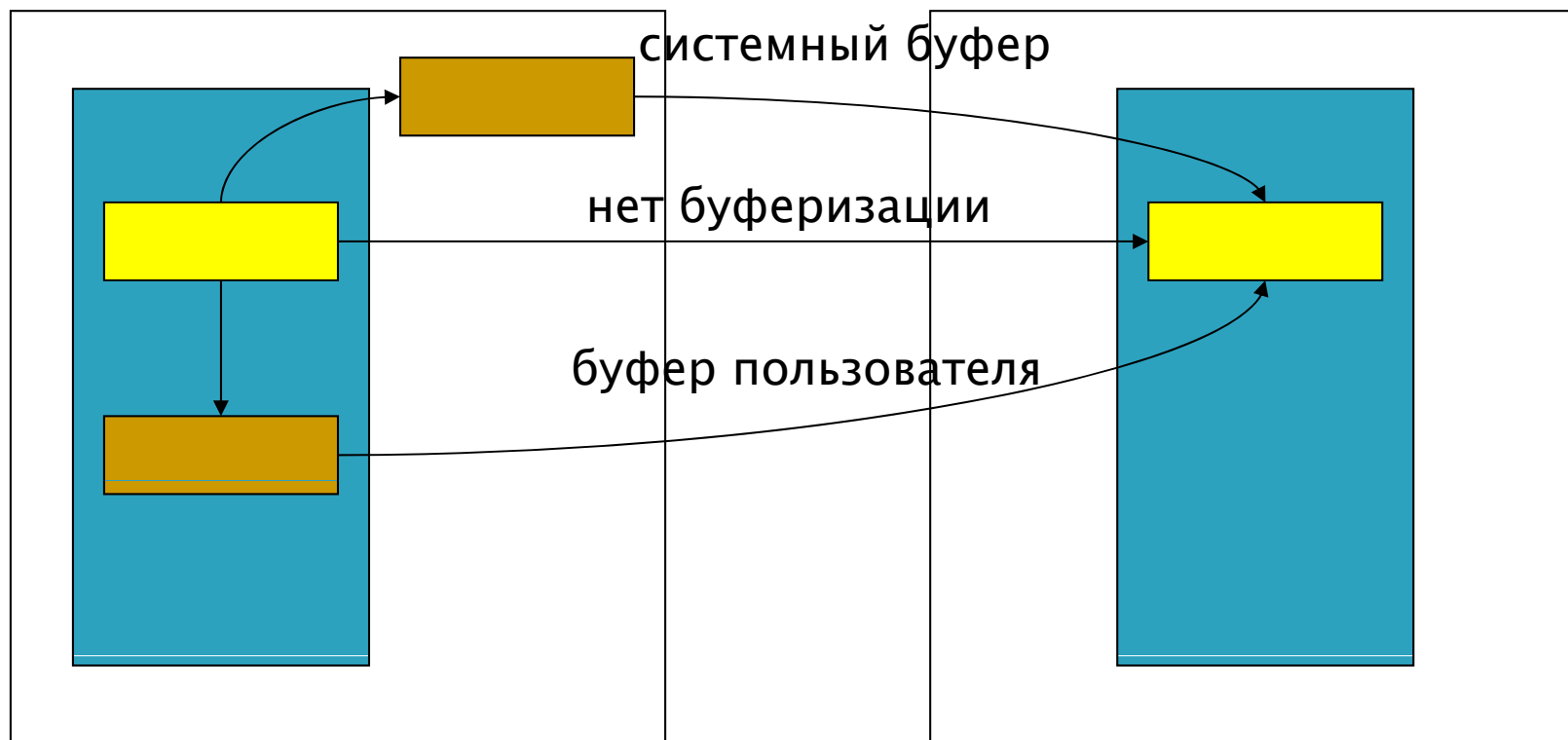
- ▶ Двухточечный обмен возможен только между процессами, принадлежащими одной области взаимодействия (одному коммутатору).

Условия успешного взаимодействия точка–точка

- ▶ Отправитель должен указать правильный rank получателя
- ▶ Получатель должен указать верный rank отправителя
- ▶ Одинаковый коммутатор
- ▶ Тэги должны соответствовать друг другу
- ▶ Буфер у процесса–получателя должен быть достаточного объема

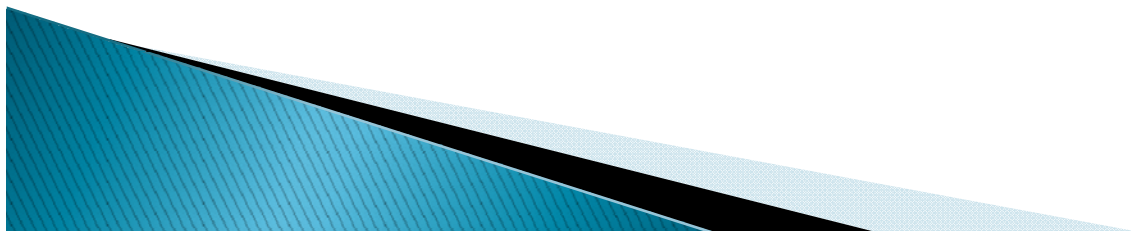


Выполнение двухточечных обменов



Разновидности двухточечного обмена

- ▶ *блокирующие* прием/передача, которые приостанавливают выполнение процесса на время приема или передачи сообщения;
- ▶ *неблокирующие* прием/передача, при которых выполнение процесса продолжается в фоновом режиме, а программа в нужный момент может запросить подтверждение завершения приема сообщения.



Двухточечный (point-to-point, p2p) обмен

- ▶ Правильно организованный двухточечный обмен сообщениями должен исключать возможность блокировки или некорректной работы параллельной MPI-программы.
- ▶ Примеры ошибок в организации двухточечных обменов:
 - ❑ выполняется передача сообщения, но не выполняется его прием;
 - ❑ процесс-источник и процесс-получатель одновременно пытаются выполнить блокирующие передачу или прием сообщения.



Двухточечный (point-to-point, p2p) обмен

В MPI приняты следующие соглашения об именах подпрограмм двухточечного обмена:

`MPI_[I][R, S, B]Send`

здесь префикс [I] (Immediate) обозначает неблокирующий режим.

Один из префиксов [R, S, B] обозначает режим обмена: по готовности, синхронный и буферизованный.

Отсутствие префикса обозначает подпрограмму стандартного обмена.

Имеется 8 разновидностей операции передачи сообщений.

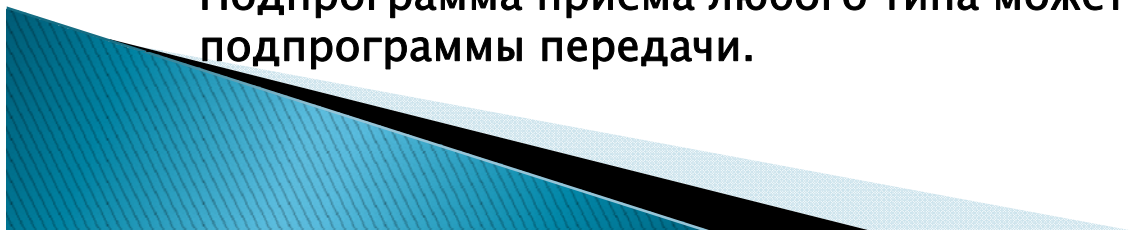
Для подпрограмм приема:

`MPI_[I]Recv`

то есть всего 2 разновидности приема.

Подпрограмма `MPI_Irsend`, например, выполняет передачу «по готовности» в неблокирующем режиме, `MPI_Bsend` буферизованную передачу с блокировкой, а `MPI_Recv` выполняет блокирующий прием сообщений.

Подпрограмма приема любого типа может принять сообщения от любой подпрограммы передачи.

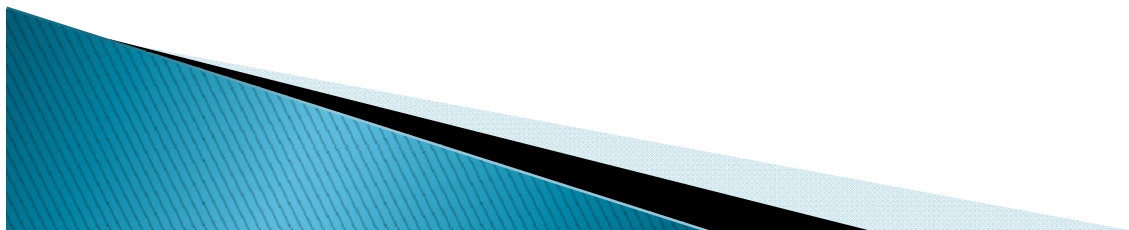


Стандартная блокирующая передача

```
int MPI_Send(void *buf, int count, MPI_Datatype datatype, int  
dest, int tag, MPI_Comm comm)
```

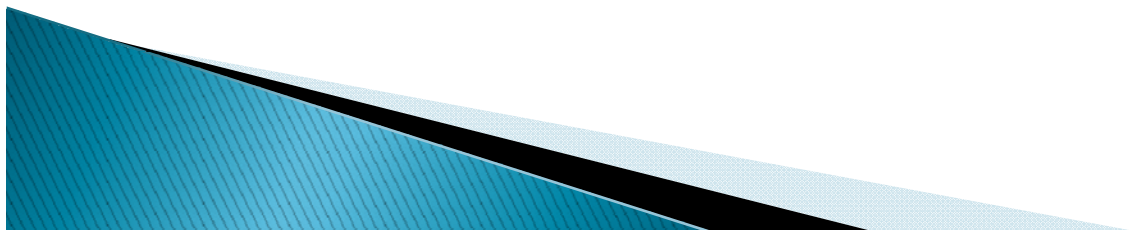
```
MPI_Send(buf, count, datatype, dest, tag, comm, ierr)
```

- ▶ buf – адрес первого элемента в буфере передачи;
- ▶ count – количество элементов в буфере передачи (допускается count = 0);
- ▶ datatype – тип MPI каждого пересылаемого элемента;
- ▶ dest – ранг процесса-получателя сообщения (целое число от 0 до n – 1, где n – число процессов в области взаимодействия);
- ▶ tag – тег сообщения;
- ▶ comm – коммуникатор;
- ▶ ierr – код завершения.



Стандартная блокирующая передача

- ▶ При стандартной блокирующей передаче после завершения вызова (после возврата из функции/процедуры передачи) можно использовать любые переменные, использовавшиеся в списке параметров. Такое использование не повлияет на корректность обмена.
- ▶ Дальнейшая «судьба» сообщения зависит от реализации MPI. Сообщение может быть сразу передано процессу-получателю или может быть скопировано в буфер передачи.
- ▶ Завершение вызова не гарантирует доставки сообщения по назначению.



Стандартный блокирующий прием

```
int MPI_Recv(void *buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status *status)
```

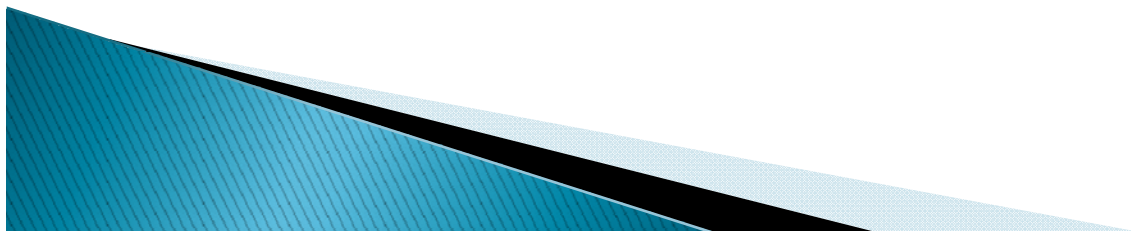
```
MPI_Recv(buf, count, datatype, dest, tag, comm, status, ierr)
```

- ▶ buf – адрес первого элемента в буфере приёма;
- ▶ count – количество элементов в буфере приёма;
- ▶ datatype – тип MPI каждого пересылаемого элемента;
- ▶ source – ранг процесса-отправителя сообщения (целое число от 0 до $n - 1$, где n – число процессов в области взаимодействия);
- ▶ tag – тег сообщения;
- ▶ comm – коммунитор;
- ▶ status – статус обмена;
- ▶ ierr – код завершения.



Стандартный блокирующий прием

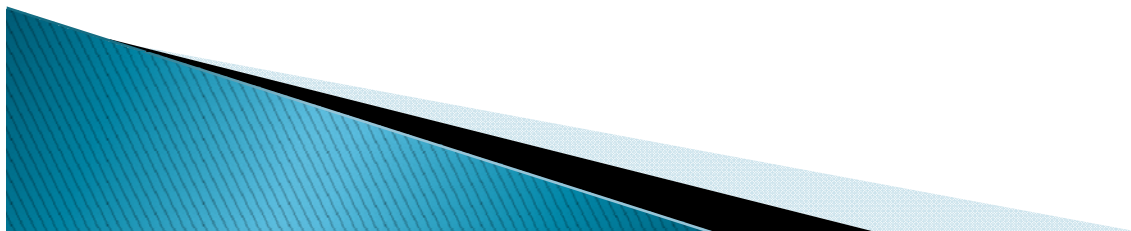
- ▶ Значение параметра `count` может оказаться больше, чем количество элементов в принятом сообщении. В этом случае после выполнения приёма в буфере изменится значение только тех элементов, которые соответствуют элементам фактически принятого сообщения.
- ▶ Для функции `MPI_Recv` гарантируется, что после завершения вызова сообщение принято и размещено в буфере приема.



Прием сообщения

Если один процесс последовательно посылает два сообщения, соответствующие одному и тому же вызову `MPI_Recv`, другому процессу, то первым будет принято сообщение, которое было отправлено раньше.

Если два сообщения были одновременно отправлены разными процессами, то порядок их получения принимающим процессом заранее не определён.



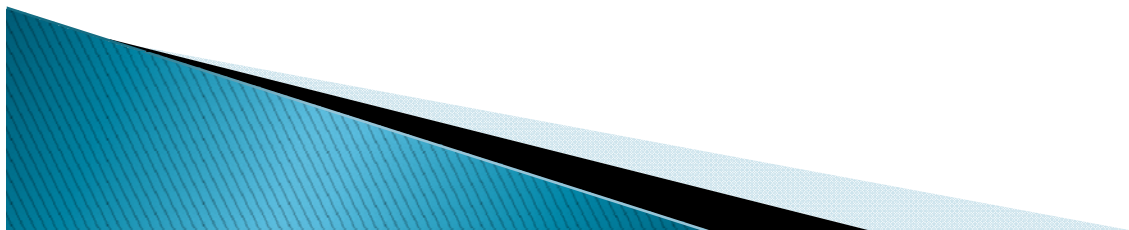
Коды завершения

- ▶ `MPI_ERR_COMM` – неправильно указан коммуникатор. Часто возникает при использовании «пустого» коммуникатора;
- ▶ `MPI_ERR_COUNT` – неправильное значение аргумента `count` (количество пересылаемых значений);
- ▶ `MPI_ERR_TYPE` – неправильное значение аргумента, задающего тип данных;
- ▶ `MPI_ERR_TAG` – неправильно указан тег сообщения;
- ▶ `MPI_ERR_RANK` – неправильно указан ранг источника или адресата сообщения;
- ▶ `MPI_ERR_ARG` – неправильный аргумент, ошибочное задание которого не попадает ни в один класс ошибок;
- ▶ `MPI_ERR_REQUEST` – неправильный запрос на выполнение операции.



Джокеры

- ▶ В качестве ранга источника сообщения и в качестве тега сообщения можно использовать «джокеры» :
 - MPI_ANY_SOURCE – любой источник;
 - MPI_ANY_TAG – любой тег.
- ▶ При использовании «джокеров» есть опасность приема сообщения, не предназначенного данному процессу

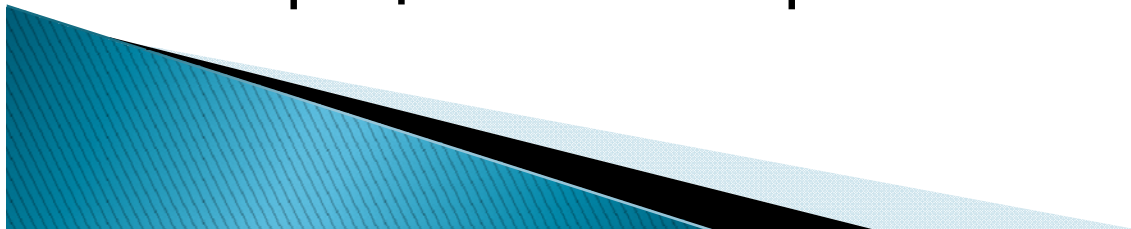


Двухточечные обмены

Подпрограмма `MPI_Recv` может принимать сообщения, отправленные в любом режиме.

Прием может выполняться от произвольного процесса, а в операции передачи должен быть указан вполне определенный адрес.

Приемник может использовать «джокеры» для источника и для тега. Процесс может отправить сообщение и самому себе, но следует учитывать, что использование в этом случае блокирующих операций может привести к «тупику».



Двухточечные обмены

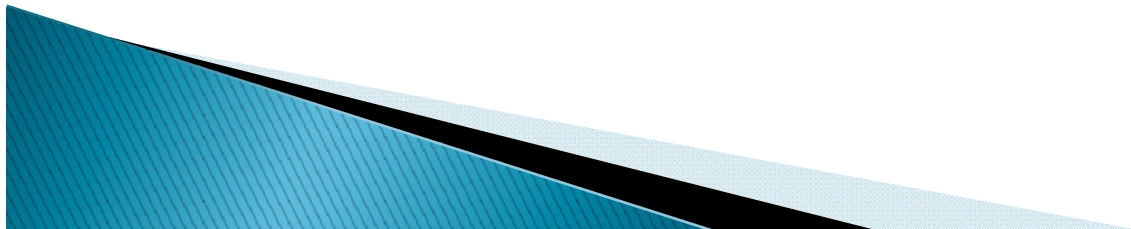
Размер полученного сообщения (count) можно определить с помощью вызова подпрограммы

```
int MPI_Get_count(MPI_Status *status, MPI_Datatype  
datatype, int *count)
```

```
MPI_Get_count(status, datatype, count, ierr)
```

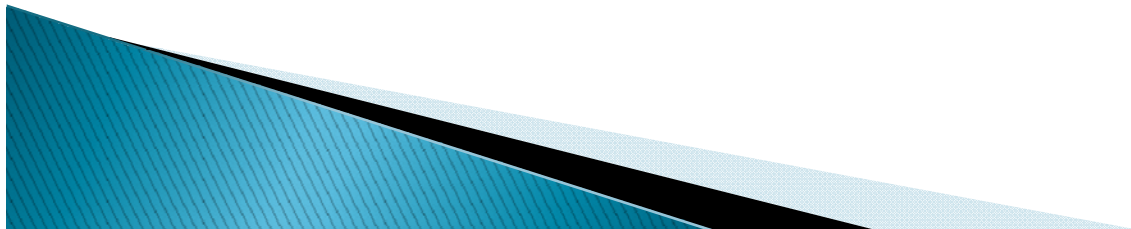
- ▶ count – количество элементов в буфере передачи;
- ▶ datatype – тип MPI каждого пересылаемого элемента;
- ▶ status – статус обмена;
- ▶ ierr – код завершения.

Аргумент datatype должен соответствовать типу данных, указанному в операции обмена



Двухточечный обмен с буферизацией

- ▶ Передача сообщения в буферизованном режиме может быть начата независимо от того, зарегистрирован ли соответствующий прием. Источник копирует сообщение в буфер, а затем передает его в неблокирующем режиме, так же как в стандартном режиме.
- ▶ Эта операция локальна, поскольку ее выполнение не зависит от наличия соответствующего приема.
- ▶ Если объем буфера недостаточен, возникает ошибка. Выделение буфера и его размер контролируются программистом.



Двухточечный обмен с буферизацией

- ▶ Размер буфера должен превосходить размер сообщения на величину `MPI_BSEND_OVERHEAD`. Это дополнительное пространство используется подпрограммой буферизованной передачи для своих целей.
- ▶ Если перед выполнением операции буферизованного обмена не выделен буфер, MPI ведет себя так, как если бы с процессом был связан буфер нулевого размера. Работа с таким буфером обычно завершается сбоем программы.
- ▶ Буферизованный обмен рекомендуется использовать в тех ситуациях, когда программисту требуется больший контроль над распределением памяти. Этот режим удобен и для отладки, поскольку причину переполнения буфера определить легче, чем причину тупика.



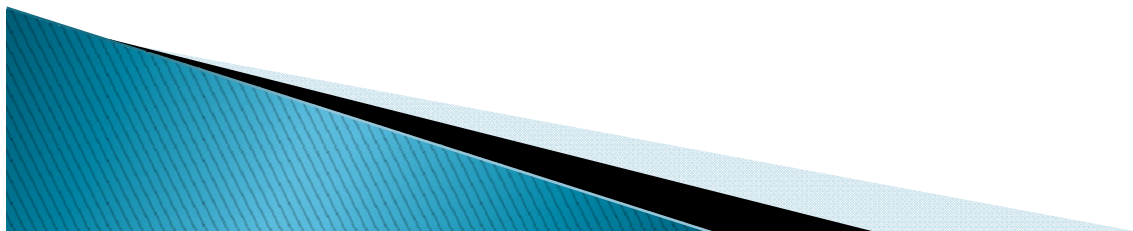
Двухточечный обмен с буферизацией

- ▶ При выполнении буферизованного обмена программист должен заранее создать буфер достаточного размера:

```
int MPI_Buffer_attach(void *buf, size)
```

```
MPI_Buffer_attach(buf, size, ierr)
```

- ▶ В результате вызова создается буфер `buf` размером `size` байтов. В программах на языке Fortran роль буфера может играть массив. За один раз к процессу может быть подключен только один буфер.

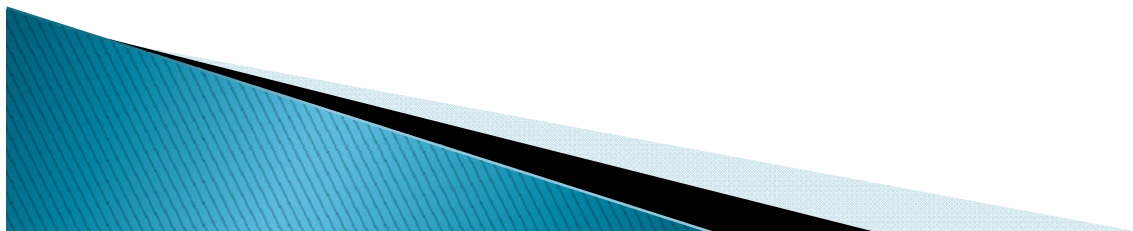


Двухточечный обмен с буферизацией

- ▶ Буферизованная передача завершается сразу, поскольку сообщение немедленно копируется в буфер для последующей передачи:

```
int MPI_Bsend(void *buf, int count,  
             MPI_Datatype datatype, int dest, int tag,  
             MPI_Comm comm)
```

```
MPI_Bsend(buf, count, datatype, dest, tag,  
          comm, ierr)
```



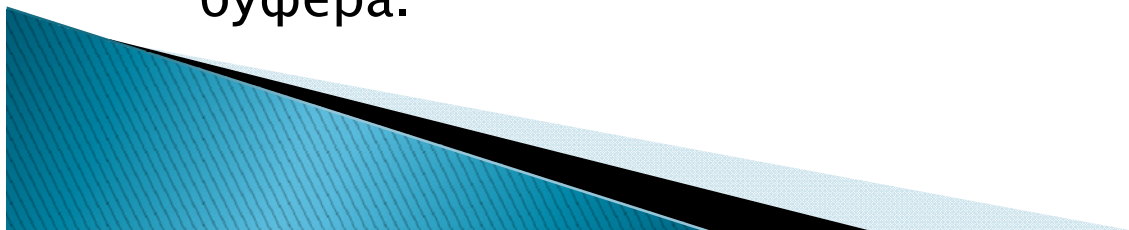
Двухточечный обмен с буферизацией

- ▶ После завершения работы с буфером его необходимо отключить:

```
int MPI_Buffer_detach(void *buf, int *size)
```

```
MPI_Buffer_detach(buf, size, ierr)
```

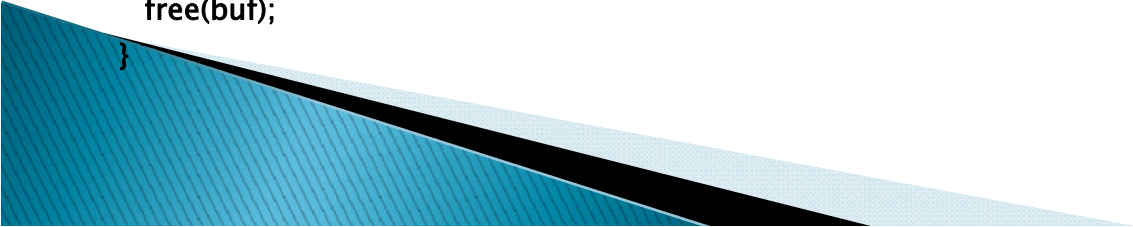
- ▶ Возвращается адрес (`buf`) и размер отключаемого буфера (`size`). Эта операция блокирует работу процесса до тех пор, пока все сообщения, находящиеся в буфере, не будут обработаны. Вызов данной подпрограммы можно использовать для форсированной передачи сообщений. После завершения вызова можно вновь использовать память, которую занимал буфер. В языке С данный вызов не освобождает автоматически память, отведенную для буфера.



Двухточечный обмен с буферизацией

```
#include <mpi.h>
#include <stdio.h>
#define M 10
int main( int argc, char **argv )
{
    int n;
    int rank, size;
    MPI_Status status;
    MPI_Init( &argc, &argv );
    MPI_Comm_size( MPI_COMM_WORLD, &size );
    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
    if ( rank == 0 ) {
        int blen = M * (sizeof(int) + MPI_BSEND_OVERHEAD);
        int *buf = (int*) malloc(blen);
        MPI_Buffer_attach (buf, blen);
        for(int i = 0; i < M; i ++ ) {
            n = i;
            MPI_Bsend (&n, 1, MPI_INT, 1, i, MPI_COMM_WORLD );
        }
        MPI_Buffer_detach(&buf, &blen);
        free(buf);
    }

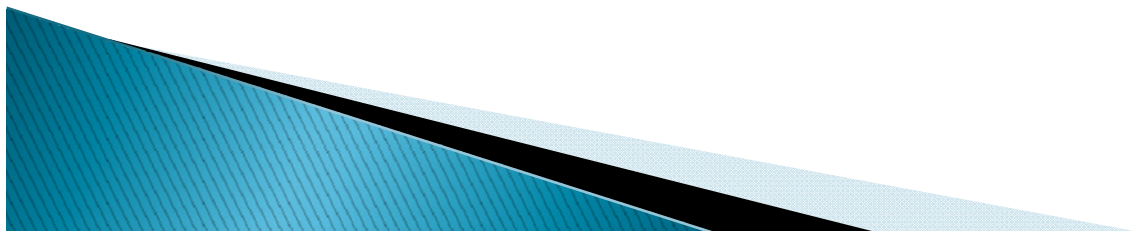
    else if ( rank == 1 ) {
        for(int i = 0; i < M; i ++ ) {
            MPI_Recv (&n, 1, MPI_INT, 0, i,
                    MPI_COMM_WORLD,&status );
        }
    }
    MPI_Finalize();
    return 0;
}
```



Синхронный режим

- ▶ Завершение передачи происходит только после того, как прием сообщения инициализирован другим процессом.
- ▶ Посылающая сторона запрашивает у принимающей стороны подтверждение выдачи операции receive – «КВИТАНЦИЮ».

```
int MPI_Ssend(void *buf, int count,  
MPI_Datatype datatype, int dest, int tag,  
MPI_Comm comm)
```



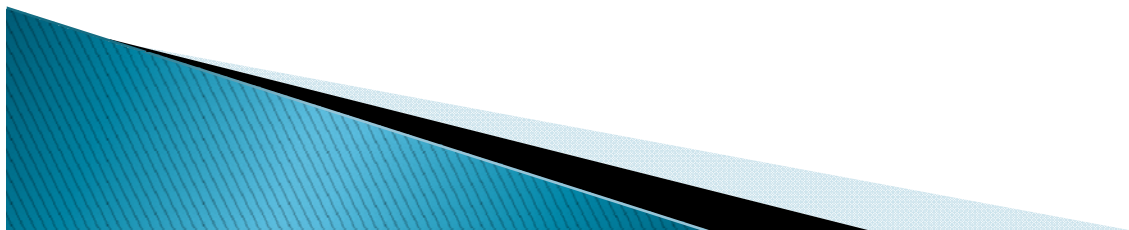
Режим «ПО ГОТОВНОСТИ»

Передача «по готовности» выполняется с помощью подпрограммы

```
int MPI_Rsend(void *buf, int count, MPI_Datatype  
datatype, int dest, int tag, MPI_Comm comm)
```

```
MPI_Rsend(buf, count, datatype, dest, tag, comm,  
ierr)
```

Передача «по готовности» должна начинаться, если уже зарегистрирован соответствующий прием. При несоблюдении этого условия результат выполнения операции не определен.

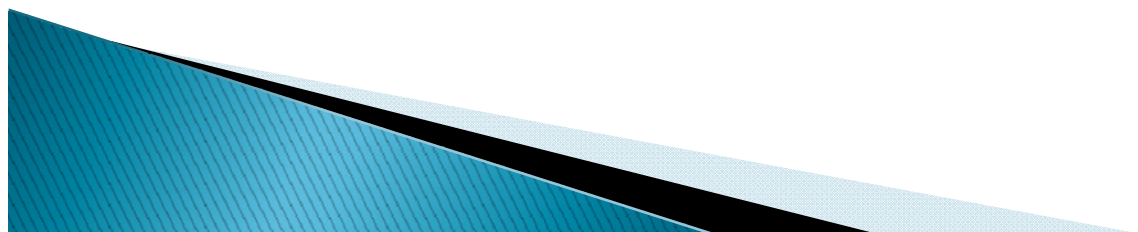


Режим «ПО ГОТОВНОСТИ»

Завершается она сразу же. Если прием не зарегистрирован, результат выполнения операции не определен.

Завершение передачи не зависит от того, вызвана ли другим процессом подпрограмма приема данного сообщения или нет, оно означает только, что буфер передачи можно использовать вновь.

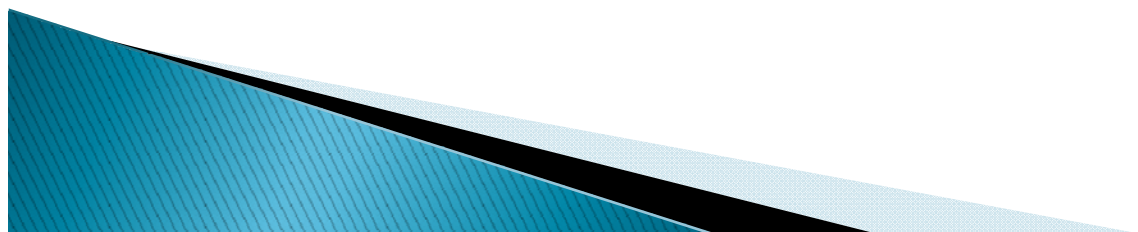
Сообщение просто выбрасывается в коммуникационную сеть в надежде, что адресат его получит. Эта надежда может и не сбыться.



Режим «ПО ГОТОВНОСТИ»

Обмен «по готовности» может увеличить производительность программы, поскольку здесь не используются этапы установки межпроцессных связей, а также буферизация.

Все это — операции, требующие времени. С другой стороны, обмен «по готовности» потенциально опасен, кроме того, он усложняет отладку, поэтому его рекомендуется использовать только в том случае, когда правильная работа программы гарантируется ее логической структурой, а выигрыша в быстродействии надо добиться любой ценой.



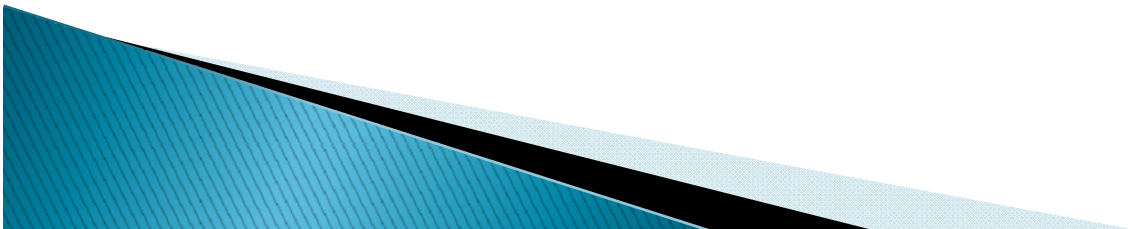
Совместные прием и передача

- ▶ Подпрограмма `MPI_Sendrecv` выполняет прием и передачу данных с блокировкой:

```
int MPI_Sendrecv(void *sendbuf, int sendcount,  
MPI_Datatype sendtype, int dest, int sendtag, void  
*recvbuf, int recvcount, MPI_Datatype recvtype,  
int source, int recvtag, MPI_Comm comm, MPI_Status  
*status)
```

- ▶ Подпрограмма `MPI_Sendrecv_replace` выполняет прием и передачу данных, используя общий буфер для передачи и приёма:

```
int MPI_Sendrecv_replace(void *buf, int count,  
MPI_Datatype datatype, int dest, int sendtag, int  
source, int recvtag, MPI_Comm comm, MPI_Status  
*status)
```



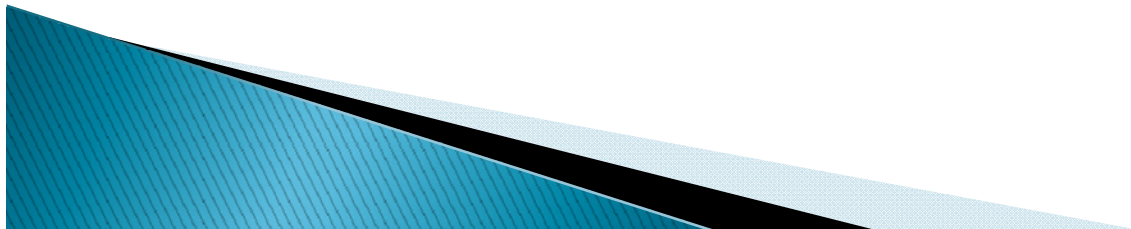
Неблокирующие обмены

- ▶ Вызов подпрограммы неблокирующей передачи инициирует, но не завершает ее. Завершиться выполнение подпрограммы может еще до того, как сообщение будет скопировано в буфер передачи.
- ▶ Применение неблокирующих операций улучшает производительность программы, поскольку в этом случае допускается перекрытие (то есть одновременное выполнение) вычислений и обменов. Передача данных из буфера или их считывание может происходить одновременно с выполнением процессом другой работы.



Неблокирующие обмены

- ▶ Для завершения неблокирующего обмена требуется вызов дополнительной процедуры, которая проверяет, скопированы ли данные в буфер передачи.
- ▶ **ВНИМАНИЕ!**
При неблокирующем обмене возвращение из подпрограммы обмена происходит сразу, но запись в буфер или считывание из него после этого производить нельзя – сообщение может быть еще не отправлено или не получено и работа с буфером может «испортить» его содержимое.



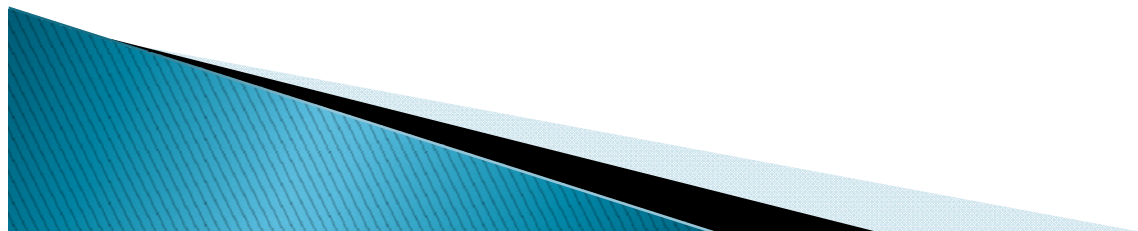
Неблокирующие обмены

Неблокирующий обмен выполняется в два этапа:

1. инициализация обмена;
2. проверка завершения обмена.

Разделение этих шагов делает необходимым *маркировку* каждой операции обмена, которая позволяет целенаправленно выполнять проверки завершения соответствующих операций.

Для маркировки в неблокирующих операциях используются *идентификаторы операций обмена*



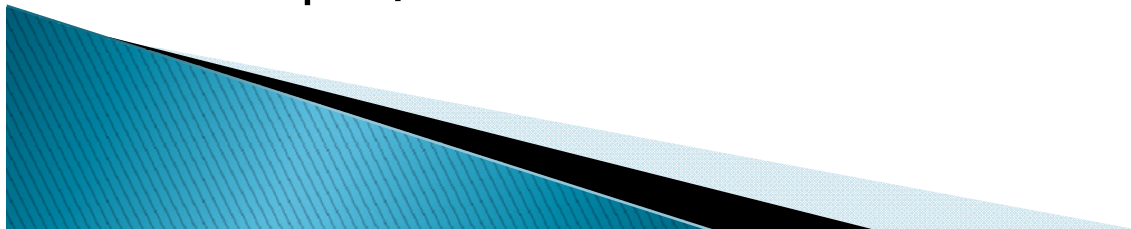
Неблокирующие обмены

- ▶ Инициализация неблокирующей стандартной передачи выполняется подпрограммами `MPI_I[S, B, R]send`.
Стандартная неблокирующая передача выполняется подпрограммой:

```
int MPI_Isend(void *buf, int count, MPI_Datatype  
datatype, int dest, int tag, MPI_Comm comm,  
MPI_Request *request)
```

```
MPI_Isend(buf, count, datatype, dest, tag, comm,  
request, ierr)
```

- ▶ Входные параметры этой подпрограммы аналогичны аргументам подпрограммы `MPI_Send`.
- ▶ Выходной параметр `request` – идентификатор операции.



Неблокирующие обмены

- ▶ Инициализация неблокирующего приема выполняется при вызове подпрограммы:

```
int MPI_Irecv(void *buf, int count,  
MPI_Datatype datatype, int source, int tag,  
MPI_Comm comm, MPI_Request *request)
```

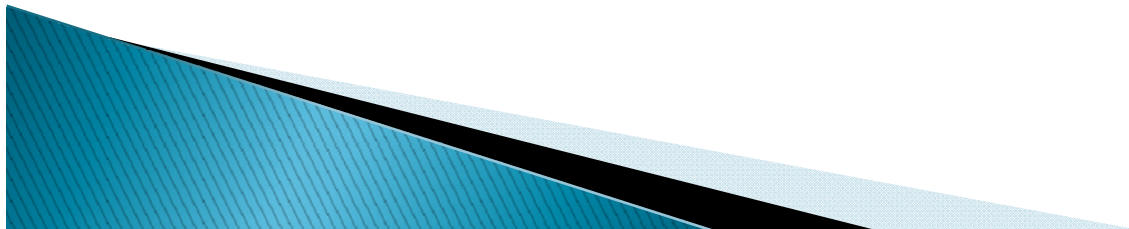
```
MPI_Irecv(buf, count, datatype, source, tag,  
comm, request, ierr)
```

- ▶ Назначение аргументов здесь такое же, как и в ранее рассмотренных подпрограммах, за исключением того, что указывается ранг не адресата, а источника сообщения (`source`).



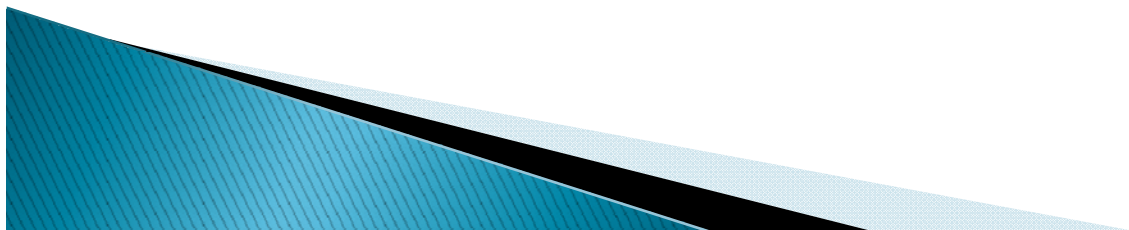
Неблокирующие обмены

- ▶ Вызовы подпрограмм неблокирующего обмена формируют *запрос* на выполнение операции обмена и связывают его с идентификатором операции `request`.
- ▶ Запрос идентифицирует свойства операции обмена:
 - ☐ режим;
 - ☐ характеристики буфера обмена;
 - ☐ контекст;
 - ☐ тег и ранг.
- ▶ Запрос содержит информацию о состоянии ожидающих обработки операций обмена и может быть использован для получения информации о состоянии обмена или для ожидания его завершения.



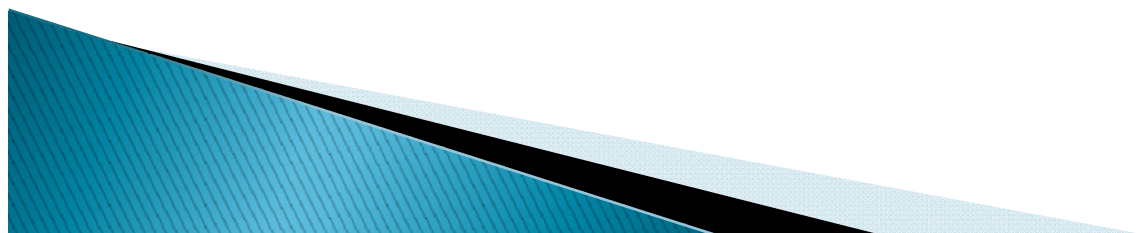
Проверка выполнения обмена

- ▶ Проверка фактического выполнения передачи или приема в неблокирующем режиме осуществляется с помощью вызова *подпрограмм ожидания*, блокирующих работу процесса до завершения операции или неблокирующих подпрограмм проверки, возвращающих логическое значение «истина», если операция выполнена



Проверка выполнения обмена

- ▶ В том случае, когда одновременно несколько процессов обмениваются сообщениями, можно использовать проверки, которые применяются одновременно к нескольким обменам.
- ▶ Есть три типа таких проверок:
 1. проверка завершения всех обменов;
 2. проверка завершения любого обмена из нескольких;
 3. проверка завершения заданного обмена из нескольких.
- ▶ Каждая из этих проверок имеет две разновидности:
 1. «ожидание»;
 2. «проверка».



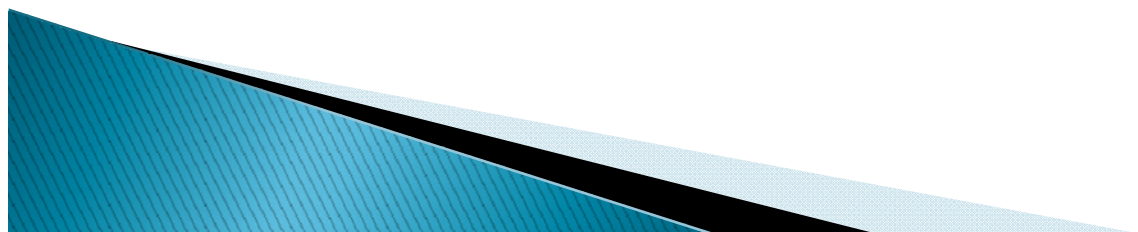
Блокирующие операции проверки

- ▶ Подпрограмма `MPI_Wait` блокирует работу процесса до завершения приема или передачи сообщения:

```
int MPI_Wait(MPI_Request *request, MPI_Status  
*status)
```

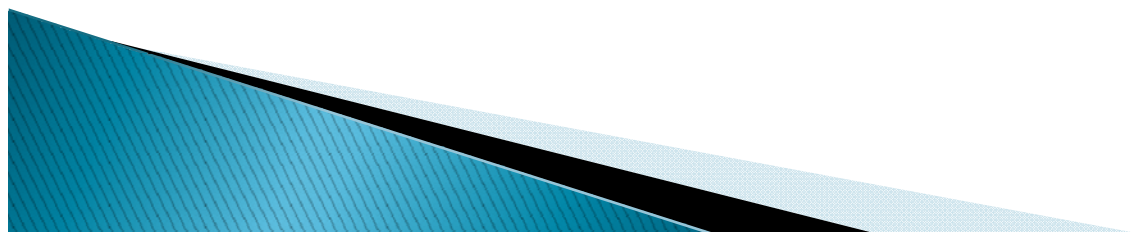
```
MPI_Wait(request, status, ierr)
```

- ▶ Входной параметр `request` — идентификатор операции обмена, выходной — статус (`status`).



Блокирующие операции проверки

- ▶ Успешное выполнение подпрограммы `MPI_Wait` после вызова `MPI_Ibsend` подразумевает, что буфер передачи можно использовать вновь, то есть пересылаемые данные отправлены или скопированы в буфер, выделенный при вызове подпрограммы `MPI_Buffer_attach`.
- ▶ В этот момент уже нельзя отменить передачу. Если не будет зарегистрирован соответствующий прием, буфер нельзя будет освободить. В этом случае можно применить подпрограмму `MPI_Cancel`, которая освобождает память, выделенную подсистеме коммуникаций.



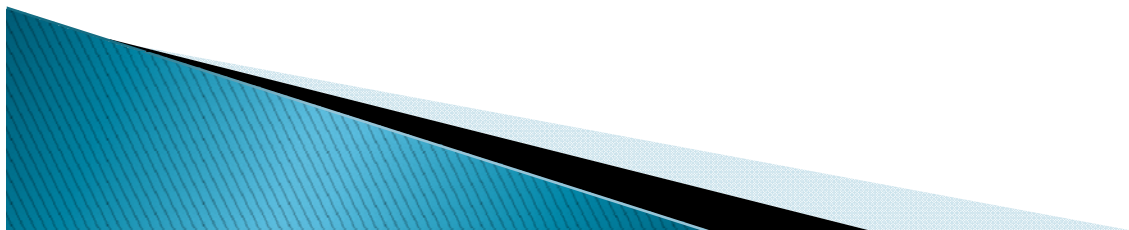
Проверка завершения всех обменов

- ▶ Проверка завершения всех обменов выполняется подпрограммой:

```
int MPI_Waitall(int count, MPI_Request  
requests[], MPI_Status statuses[])
```

```
MPI_Waitall(count, requests, statuses, ierr)
```

- ▶ При вызове этой подпрограммы выполнение процесса блокируется до тех пор, пока все операции обмена, связанные с активными запросами в массиве `requests`, не будут выполнены. Возвращается статус этих операций. Статус обменов содержится в массиве `statuses`. `count` – количество запросов на обмен (размер массивов `requests` и `statuses`).



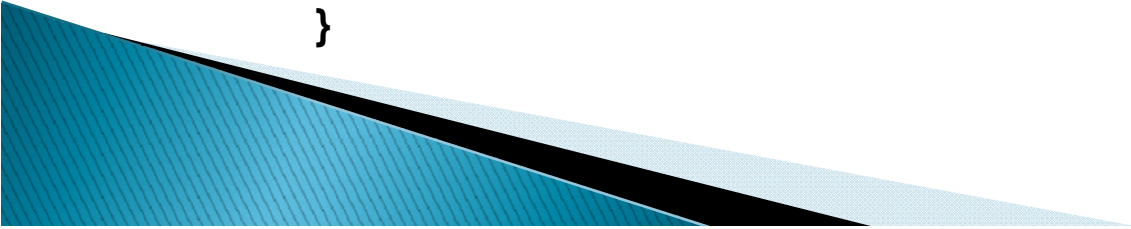
Проверка завершения всех обменов

- ▶ В результате выполнения подпрограммы `MPI_Waitall` запросы, сформированные неблокирующими операциями обмена, аннулируются, а соответствующим элементам массива присваивается значение `MPI_REQUEST_NULL`.
- ▶ В случае неуспешного выполнения одной или более операций обмена подпрограмма `MPI_Waitall` возвращает код ошибки `MPI_ERR_IN_STATUS` и присваивает полю ошибки статуса значение кода ошибки соответствующей операции.
- ▶ Если операция выполнена успешно, полю присваивается значение `MPI_SUCCESS`, а если не выполнена, но и не было ошибки – значение `MPI_ERR_PENDING`. Это соответствует наличию запросов на выполнение операции обмена, ожидающих обработки.



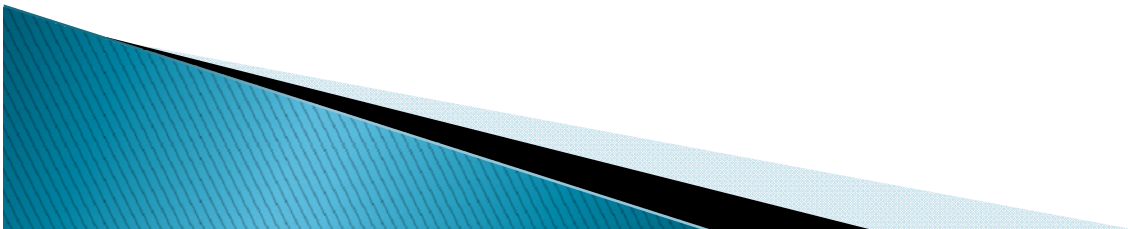
Алгоритм Якоби. Последовательная версия

```
/* Jacobi program */
#include <stdio.h>
#define L 1000
#define ITMAX 100
int i,j,it;
double A[L][L];
double B[L][L];
int main(int an, char **as)
{
    printf("JAC STARTED\n");
    for(i=0;i<=L-1;i++)
        for(j=0;j<=L-1;j++)
        {
            A[i][j]=0.;
            B[i][j]=1.+i+j;
        }
}
```



Алгоритм Якоби. Последовательная версия

```
/****** iteration loop *****/
for(it=1; it<ITMAX;it++)
{
    for(i=1;i<=L-2;i++)
        for(j=1;j<=L-2;j++)
            A[i][j] = B[i][j];
    for(i=1;i<=L-2;i++)
        for(j=1;j<=L-2;j++)
            B[i][j] = (A[i-1][j]+A[i+1][j]+A[i][j-1]+A[i][j+1])/4.;
}
return 0;
}
```



Алгоритм Якоби. MPI-версия

| | | | | | | | | |
|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| A ₀₀ | A ₀₁ | A ₀₂ | A ₀₃ | A ₀₄ | A ₀₅ | A ₀₆ | A ₀₇ | A ₀₈ |
| A ₁₀ | A ₁₁ | A ₁₂ | A ₁₃ | A ₁₄ | A ₁₅ | A ₁₆ | A ₁₇ | A ₁₈ |
| A ₂₀ | A ₂₁ | A ₂₂ | A ₂₃ | A ₂₄ | A ₂₅ | A ₂₆ | A ₂₇ | A ₂₈ |
| A ₃₀ | A ₃₁ | A ₃₂ | A ₃₃ | A ₃₄ | A ₃₅ | A ₃₆ | A ₃₇ | A ₃₈ |

| | | | | | | | | |
|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| A ₂₀ | A ₂₁ | A ₂₂ | A ₂₃ | A ₂₄ | A ₂₅ | A ₂₆ | A ₂₇ | A ₂₈ |
| A ₃₀ | A ₃₁ | A ₃₂ | A ₃₃ | A ₃₄ | A ₃₅ | A ₃₆ | A ₃₇ | A ₃₈ |
| A ₄₀ | A ₄₁ | A ₄₂ | A ₄₃ | A ₄₄ | A ₄₅ | A ₄₆ | A ₄₇ | A ₄₈ |
| A ₅₀ | A ₅₁ | A ₅₂ | A ₅₃ | A ₅₄ | A ₅₅ | A ₅₆ | A ₅₇ | A ₅₈ |
| A ₆₀ | A ₆₁ | A ₆₂ | A ₆₃ | A ₆₄ | A ₆₅ | A ₆₆ | A ₆₇ | A ₆₈ |

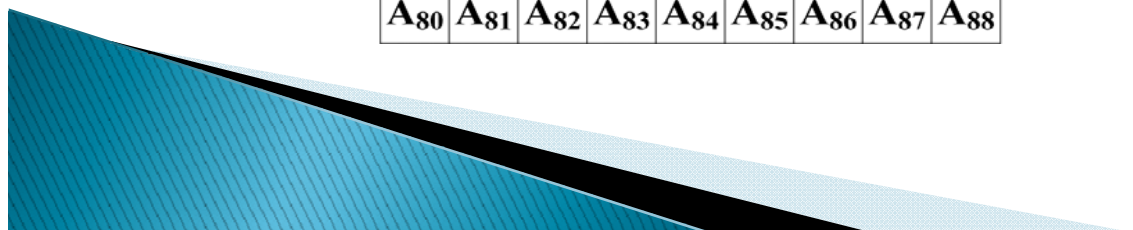
| | | | | | | | | |
|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| A ₅₀ | A ₅₁ | A ₅₂ | A ₅₃ | A ₅₄ | A ₅₅ | A ₅₆ | A ₅₇ | A ₅₈ |
| A ₆₀ | A ₆₁ | A ₆₂ | A ₆₃ | A ₆₄ | A ₆₅ | A ₆₆ | A ₆₇ | A ₆₈ |
| A ₇₀ | A ₇₁ | A ₇₂ | A ₇₃ | A ₇₄ | A ₇₅ | A ₇₆ | A ₇₇ | A ₇₈ |
| A ₈₀ | A ₈₁ | A ₈₂ | A ₈₃ | A ₈₄ | A ₈₅ | A ₈₆ | A ₈₇ | A ₈₈ |



Shadow edges

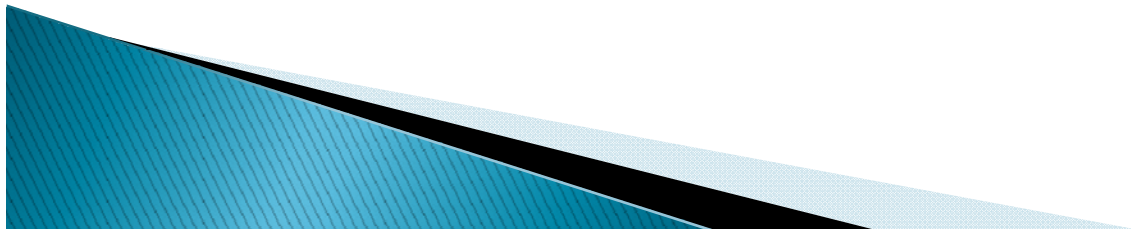


Imported elements



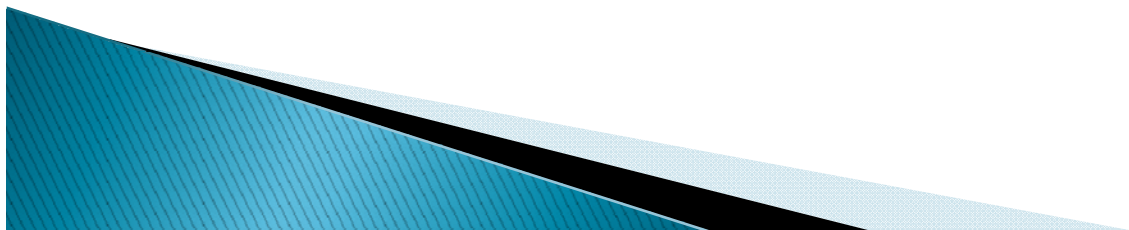
Алгоритм Якоби. MPI-версия

```
/* Jacobi-1d program */  
#include <math.h>  
#include <stdlib.h>  
#include <stdio.h>  
#include "mpi.h"  
#define m_printf if (myrank==0)printf  
#define L 1000  
#define ITMAX 100  
  
int i,j,it,k;  
int ll,shift;  
double (* A)[L];  
double (* B)[L];
```



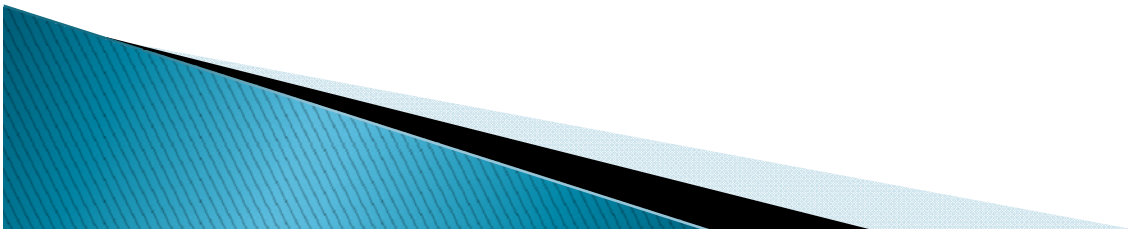
Алгоритм Якоби. MPI-версия

```
int main(int argc, char **argv)
{
    MPI_Request req[4];
    int myrank, ranksize;
    int startrow, lastrow, nrow;
    MPI_Status status[4];
    double t1, t2, time;
    MPI_Init (&argc, &argv); /* initialize MPI system */
    MPI_Comm_rank(MPI_COMM_WORLD, &myrank); /* my place in MPI system */
    MPI_Comm_size (MPI_COMM_WORLD, &ranksize); /* size of MPI system */
    MPI_Barrier(MPI_COMM_WORLD);
    /* rows of matrix I have to process */
    startrow = (myrank * L) / ranksize;
    lastrow = (((myrank + 1) * L) / ranksize) - 1;
    nrow = lastrow - startrow + 1;
    m_printf("JAC1 STARTED\n");
```



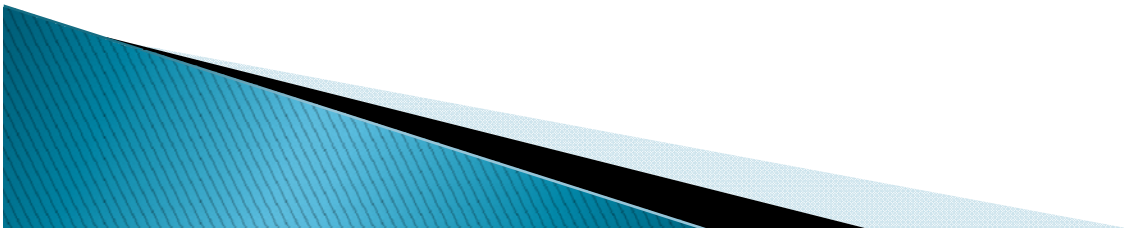
Алгоритм Якоби. MPI-версия

```
/* dynamically allocate data structures */  
A = malloc ((nrow+2) * L * sizeof(double));  
B = malloc ((nrow) * L * sizeof(double));  
for(i=1; i<=nrow; i++)  
    for(j=0; j<=L-1; j++)  
    {  
        A[i][j]=0.;  
        B[i-1][j]=1.+startrow+i-1+j;  
    }
```



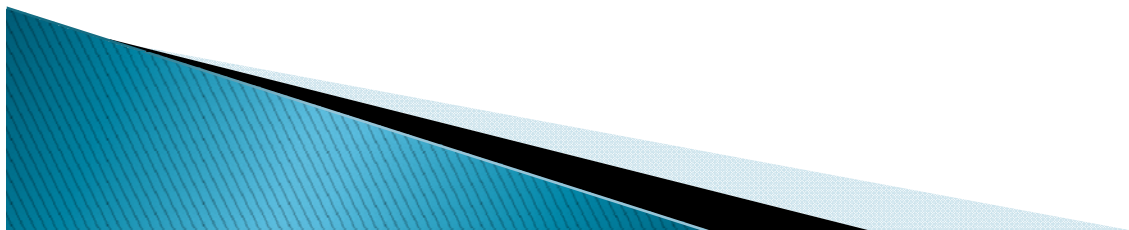
Алгоритм Якоби. MPI-версия

```
/****** iteration loop *****/
t1=MPI_Wtime();
for(it=1; it<=ITMAX; it++)
{
    for(i=1; i<=nrow; i++)
    {
        if (((i==1)&&(myrank==0))||((i==nrow)&&(myrank==ranksize-1)))
            continue;
        for(j=1; j<=L-2; j++)
        {
            A[i][j] = B[i-1][j];
        }
    }
}
```



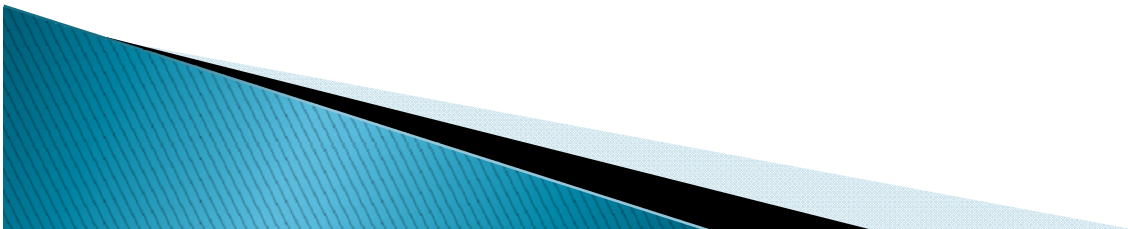
Алгоритм Якоби. MPI-версия

```
if(myrank!=0)
    MPI_Irecv(&A[0][0],L,MPI_DOUBLE, myrank-1, 1215,
              MPI_COMM_WORLD, &req[0]);
if(myrank!=ranksize-1)
    MPI_Isend(&A[nrow][0],L,MPI_DOUBLE, myrank+1, 1215,
              MPI_COMM_WORLD,&req[2]);
if(myrank!=ranksize-1)
    MPI_Irecv(&A[nrow+1][0],L,MPI_DOUBLE, myrank+1, 1216,
              MPI_COMM_WORLD, &req[3]);
if(myrank!=0)
    MPI_Isend(&A[1][0],L,MPI_DOUBLE, myrank-1, 1216,
              MPI_COMM_WORLD,&req[1]);
ll=4; shift=0;
if (myrank==0) {ll=2;shift=2;}
if (myrank==ranksize-1) {ll=2;}
MPI_Waitall(ll,&req[shift],&status[0]);
```



Алгоритм Якоби. MPI-версия

```
for(i=1; i<=nrow; i++)
{
    if (((i==1)&&(myrank==0))||((i==nrow)&&(myrank==ranksize-1))) continue;
    for(j=1; j<=L-2; j++)
        B[i-1][j] = (A[i-1][j]+A[i+1][j]+
                     A[i][j-1]+A[i][j+1])/4.;
}
/*DO it*/
printf("%d: Time of task=%lf\n",myrank,MPI_Wtime()-t1);
MPI_Finalize ();
return 0;
}
```



Проверка завершения любого числа обменов

- ▶ Проверка завершения любого числа обменов выполняется подпрограммой:


```
int MPI_Waitany(int count, MPI_Request requests[],  
int *index, MPI_Status *status)
```

- ▶ Выполнение процесса блокируется до тех пор, пока, по крайней мере, один обмен из массива запросов (`requests`) не будет завершен.

- ▶ Входные параметры:

- ☐ `requests` – запрос;
- ☐ `count` – количество элементов в массиве `requests`.

- ▶ Выходные параметры:

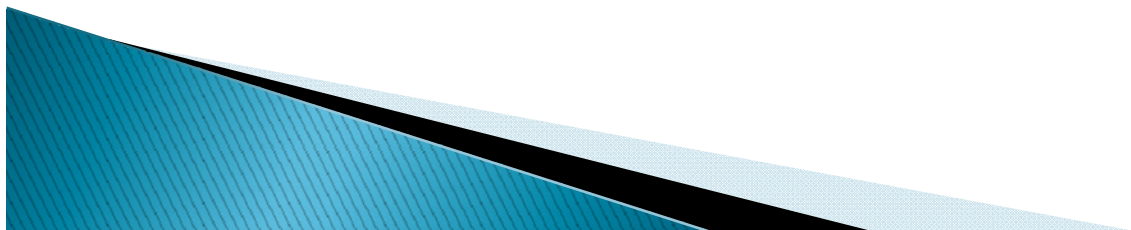
- ☐ `index` – индекс запроса (в языке C это целое число от 0 до `count - 1`) в массиве `requests`;
 - ☐ `status` – статус.
- 

Неблокирующие процедуры проверки

- ▶ Подпрограмма `MPI_Test` выполняет неблокирующую проверку завершения приема или передачи сообщения:

```
int MPI_Test(MPI_Request *request, int *flag,  
MPI_Status *status)
```

- ▶ Входной параметр: идентификатор операции обмена `request`.
- ▶ Выходные параметры:
 - ❑ `flag` — «истина», если операция, заданная идентификатором `request`, выполнена;
 - ❑ `status` — статус выполненной операции.

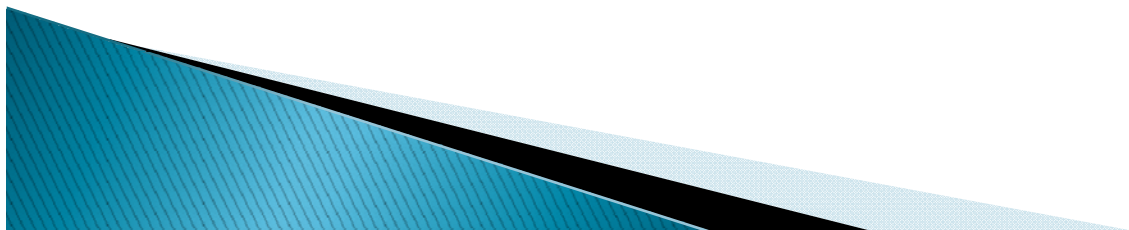


Неблокирующая проверка завершения всех обменов

- ▶ Подпрограмма `MPI_Testall` выполняет неблокирующую проверку завершения приема или передачи всех сообщений:

```
int MPI_Testall(int count, MPI_Request requests[],  
int *flag, MPI_Status statuses[])  
MPI_Testall(count, requests, flag, statuses, ierr)
```

- ▶ При вызове возвращается значение флага (`flag`) «истина», если все обмены, связанные с активными запросами в массиве `requests`, выполнены. Если завершены не все обмены, флагу присваивается значение «ложь», а массив `statuses` не определен.
- ▶ Параметр `count` – количество запросов.
- ▶ Каждому статусу, соответствующему активному запросу, присваивается значение статуса соответствующего обмена.



Неблокирующая проверка любого числа обменов

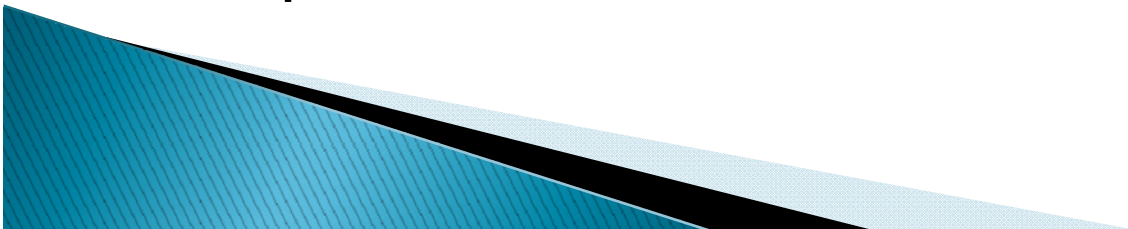
- ▶ Подпрограмма `MPI_Testany` выполняет неблокирующую проверку завершения приема или передачи сообщения:

```
int MPI_Testany(int count, MPI_Request  
requests[], int *index, int *flag, MPI_Status  
*status)
```
- ▶ Смысл и назначение параметров этой подпрограммы те же, что и для подпрограммы `MPI_Waitany`.
Дополнительный аргумент `flag`, принимает значение «истина», если одна из операций завершена.
- ▶ Блокирующая подпрограмма `MPI_Waitany` и неблокирующая `MPI_Testany` взаимозаменяемы, как и другие аналогичные пары.



Другие операции проверки

- ▶ Подпрограммы `MPI_Waitsome` и `MPI_Testsome` действуют аналогично подпрограммам `MPI_Waitany` и `MPI_Testany`, кроме случая, когда завершается более одного обмена.
- ▶ В подпрограммах `MPI_Waitany` и `MPI_Testany` обмен из числа завершенных выбирается произвольно, именно для него и возвращается статус, а для `MPI_Waitsome` и `MPI_Testsome` статус возвращается для всех завершенных обменов.
- ▶ Эти подпрограммы можно использовать для определения, сколько обменов завершено.

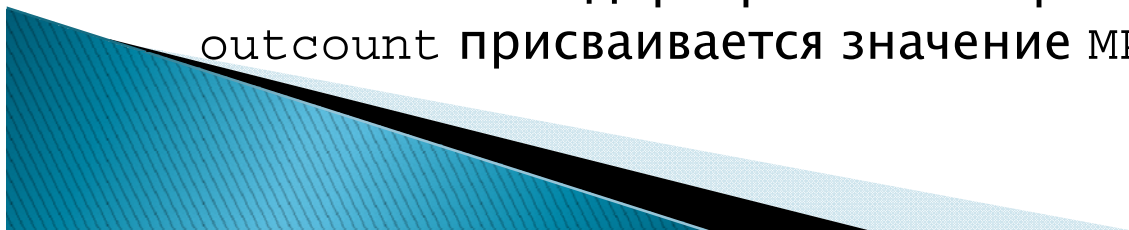


Другие операции проверки

- ▶ Блокирующая проверка выполнения обменов:

```
int MPI_Waitsome(int incount, MPI_Request  
requests[], int *outcount, int indices[],  
MPI_Status statuses[])
```

- ▶ Здесь `incount` – количество запросов. В `outcount` возвращается количество выполненных запросов из массива `requests`, а в первых `outcount` элементах массива `indices` возвращаются индексы этих операций. В первых `outcount` элементах массива `statuses` возвращается статус завершенных операций. Если выполненный запрос был сформирован неблокирующей операцией обмена, он аннулируется. Если в списке нет активных запросов, выполнение подпрограммы завершается сразу, а параметру `outcount` присваивается значение `MPI_UNDEFINED`.

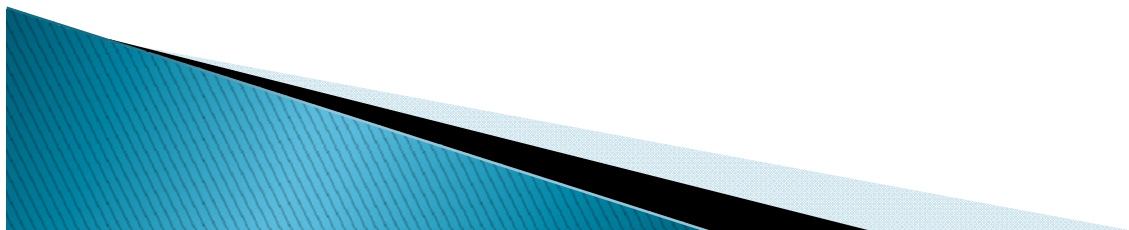


Другие операции проверки

- ▶ Неблокирующая проверка выполнения обменов:

```
int MPI_Testsome(int incount, MPI_Request  
requests[], int *outcount, int indices[],  
MPI_Status statuses[])
```

- ▶ Параметры такие же, как и у подпрограммы MPI_Waitsome. Эффективность подпрограммы MPI_Testsome выше, чем у MPI_Testany, поскольку первая возвращает информацию обо всех операциях, а для второй требуется новый вызов для каждой выполненной операции.



Проверка статуса операции приема сообщения

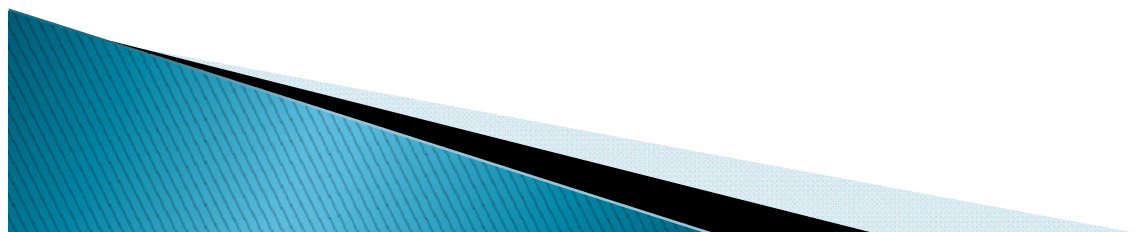
- ▶ Блокирующая проверка:

```
int MPI_Probe(int source, int tag, MPI_Comm comm,  
             MPI_Status* status)
```

- ▶ Неблокирующая проверка:


```
int MPI_Iprobe(int source, int tag, MPI_Comm comm,  
              int *flag, MPI_Status *status)
```

- ▶ Входные параметры этой подпрограммы те же, что и у подпрограммы MPI_Probe.
- ▶ Выходные параметры:
 - ☐ flag – флаг;
 - ☐ status – статус.
- ▶ Если сообщение уже поступило и может быть принято, возвращается значение флага «истина».



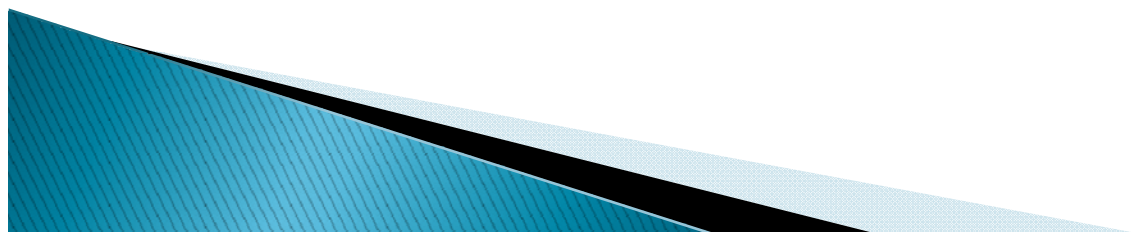
Проверка статуса операции приема сообщения (пример)

```
if (rank == 0) {
    MPI_Send(buf, size, MPI_INT, 1, 0, MPI_COMM_WORLD); // Send to process 1
    printf("0 sent %d numbers to 1\n", size);
} else if (rank == 1) {
    MPI_Status status;
    MPI_Probe(0, 0, MPI_COMM_WORLD, &status); // Probe for an incoming msg.
    // When probe returns, the status object has the size and other
    // attributes of the incoming message. Get the size of the message.
    MPI_Get_count(&status, MPI_INT, &size);
    // Allocate a buffer just big enough to hold the incoming numbers
    int* bufnumber = (int*)malloc(sizeof(int) * size);
    // Now receive the message with the allocated buffer
    MPI_Recv(bufnumber, size, MPI_INT, 0, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    printf("1 dynamically received %d numbers from 0.\n", size);
    free(bufnumber);
}
```



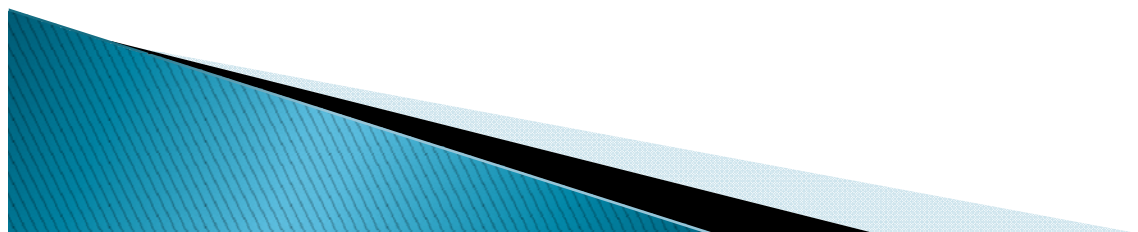
Коллективные операции

- ▶ Передача сообщений между группой процессов
- ▶ Вызываются ВСЕМИ процессами в коммутаторе
- ▶ Примеры:
 - Broadcast, scatter, gather (рассылка данных)
 - Global sum, global maximum, и.т.д. (редукционные операции)
 - Барьерная синхронизация



Характеристики коллективных передач

- ▶ Коллективные операции не являются помехой операциям типа точка–точка и наоборот
- ▶ Все процессы коммутатора должны вызывать коллективную операцию
- ▶ Синхронизация не гарантируется (за исключением барьера)
- ▶ Нет тэгов
- ▶ Принимающий буфер должен точно соответствовать размеру отсылаемого буфера

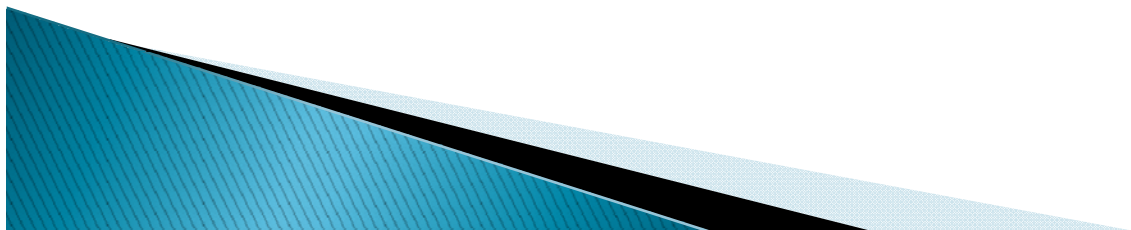


Широковещательная рассылка

- ▶ One-to-all передача: один и тот же буфер отсылается от процесса root всем остальным процессам в коммуникаторе

```
int MPI_Bcast (void *buffer, int, count,  
MPI_Datatype datatype,int root, MPI_Comm  
comm)
```

- ▶ Все процессы должны указать один тот же root и communicator

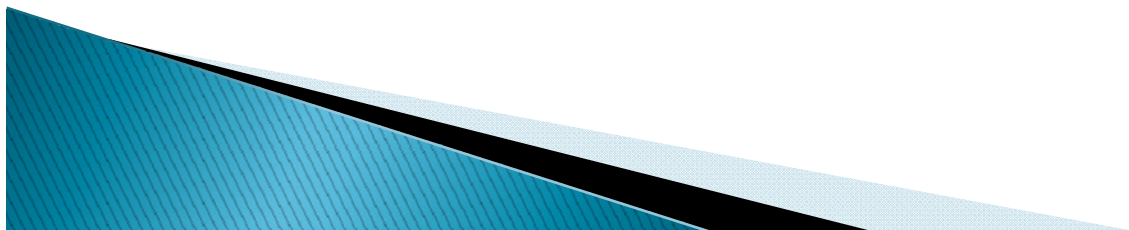


Scatter

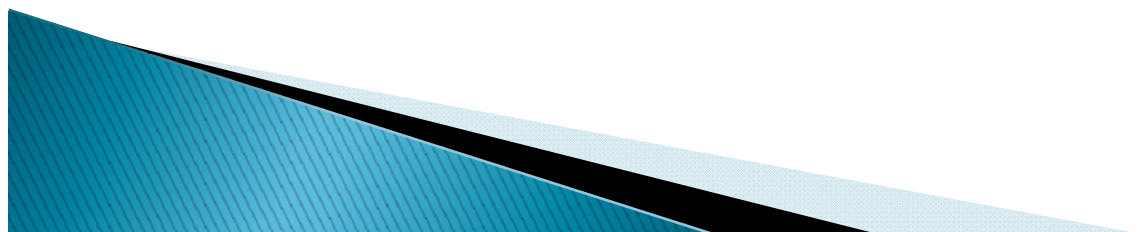
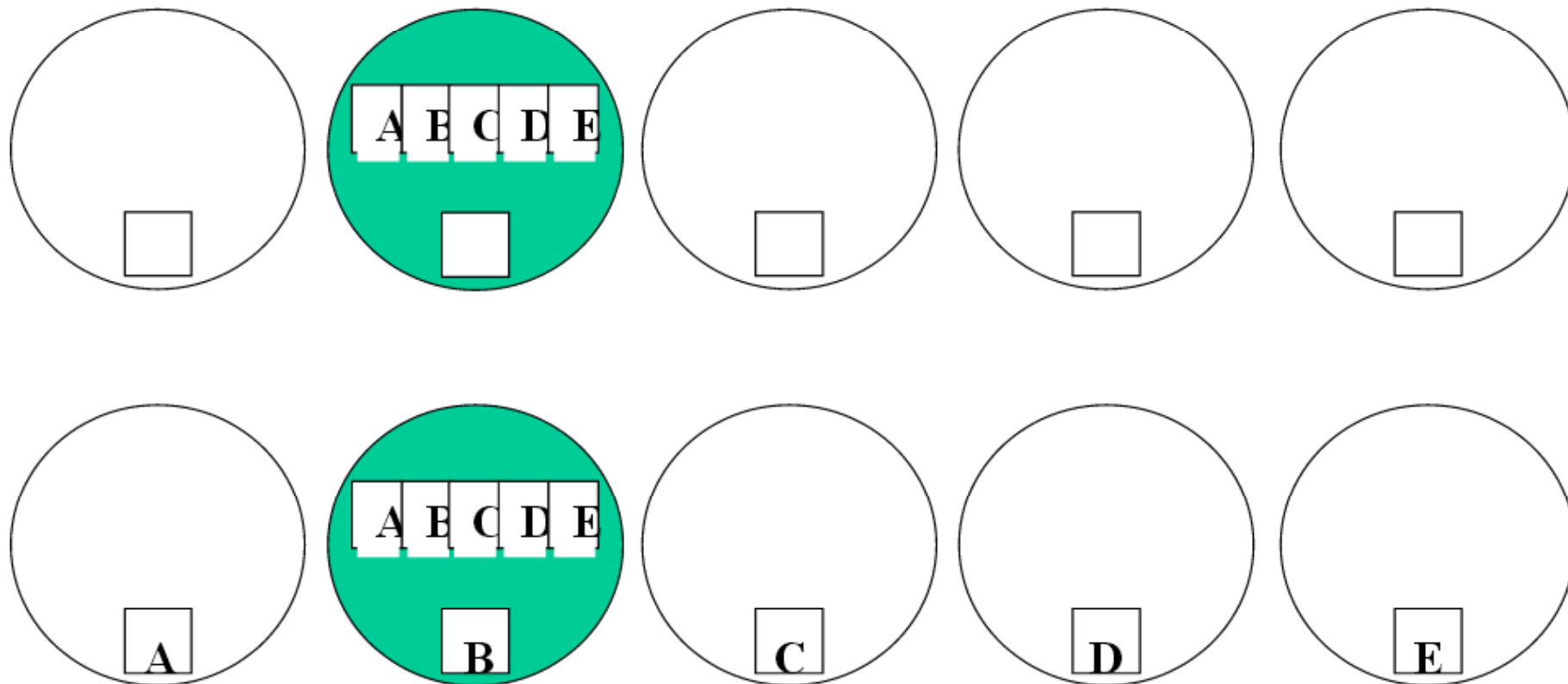
- ▶ One-to-all communication: различные данные из одного процесса рассылаются всем процессам коммутатора (в порядке их номеров)

```
int MPI_Scatter(void* sendbuf, int sendcount,  
               MPI_Datatype sendtype, void* recvbuf,  
               int recvcount, MPI_Datatype recvtype,  
               int root, MPI_Comm comm)
```

- ▶ sendcount – число элементов, посланных каждому процессу, не общее число отосланных элементов;
- ▶ send параметры имеют смысл только для процесса root



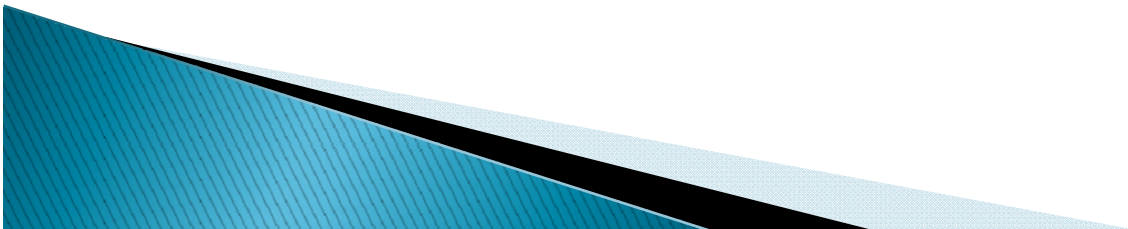
Scatter



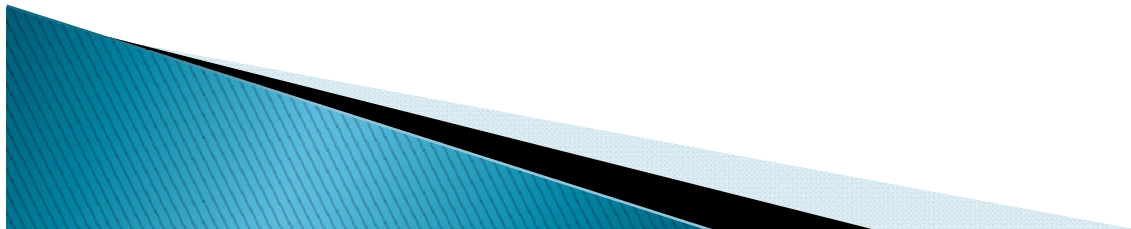
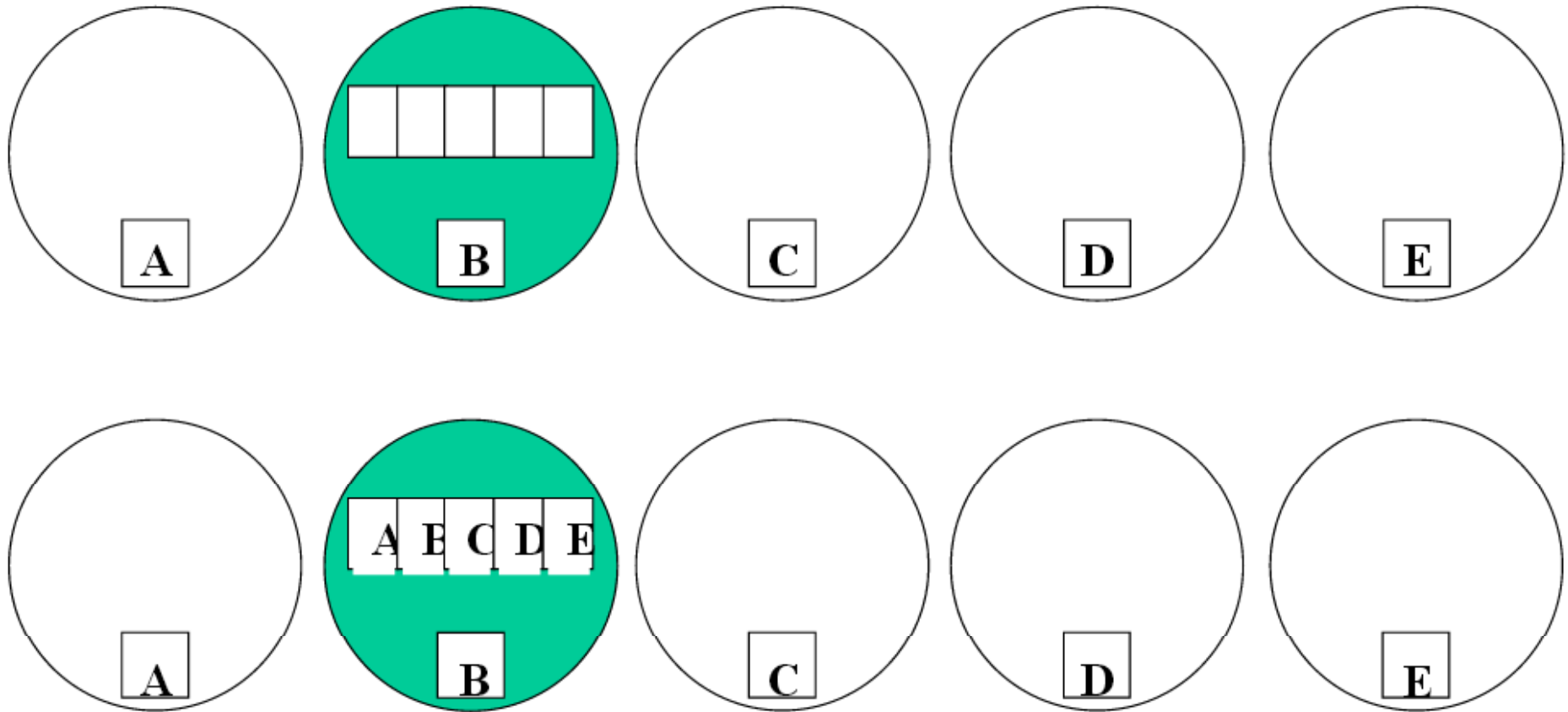
Gather

- ▶ All-to-one передачи: различные данные собираются процессом root
- ▶ Сбор данных выполняется в порядке номеров процессов.
- ▶ Длина блоков предполагается одинаковой, т.е. данные, посланные процессом i из своего буфера `sendbuf`, помещаются в i -ю порцию буфера `recvbuf` процесса root.
- ▶ Длина массива, в который собираются данные, должна быть достаточной для их размещения.

```
int MPI_Gather(void* sendbuf, int sendcount,  
              MPI_Datatype sendtype, void* recvbuf, int  
              recvcount,  
              MPI_Datatype recvtype, int root, MPI_Comm comm)
```

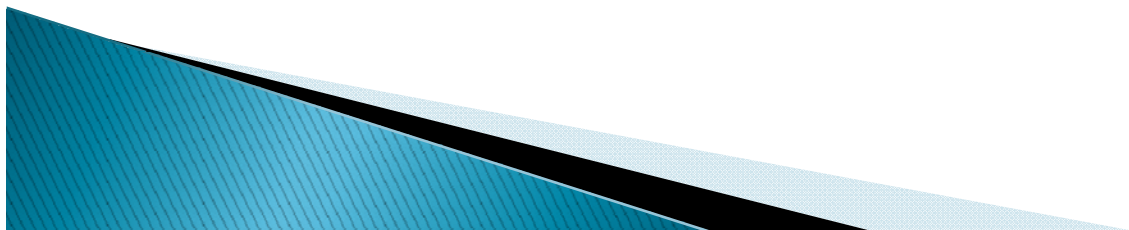


Gather



Глобальные операции редукции

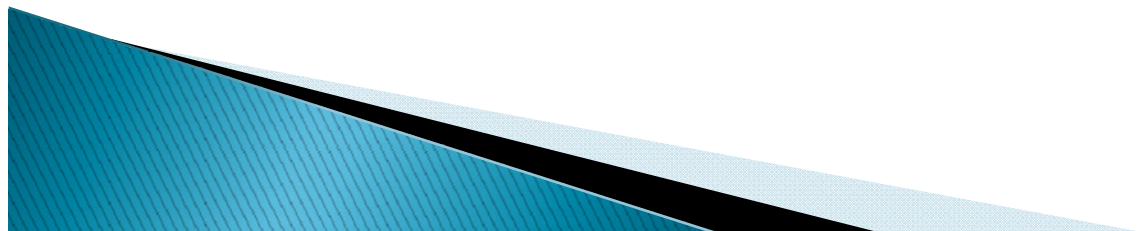
- ▶ Операции выполняются над данными, распределенными по процессам коммутатора
- ▶ Примеры:
 - Глобальная сумма или произведение
 - Глобальный максимум (минимум)
 - Глобальная операция, определенная пользователем



Глобальные операции редукции

```
int MPI_Reduce(void* sendbuf, void* recvbuf, int  
count, MPI_Datatype datatype, MPI_Op op, int root,  
MPI_Comm comm)
```

- ▶ count число операций "op", выполняемых над последовательными элементами буфера
- ▶ sendbuf (также размер recvbuf)
- ▶ op является ассоциативной операцией, которая выполняется над парой операндов типа datatype и возвращает результат того же типа:
MPI_MAX, MPI_MIN, MPI_SUM, MPI_PROD, MPI_LAND, MPI_BAND, MPI_LOR, MPI_BOR, MPI_LXOR, MPI_BXOR, MPI_MAXLOC, MPI_MINLOC
- ▶ MPI_ALLREDUCE – нет root процесса (все получают рез-т)

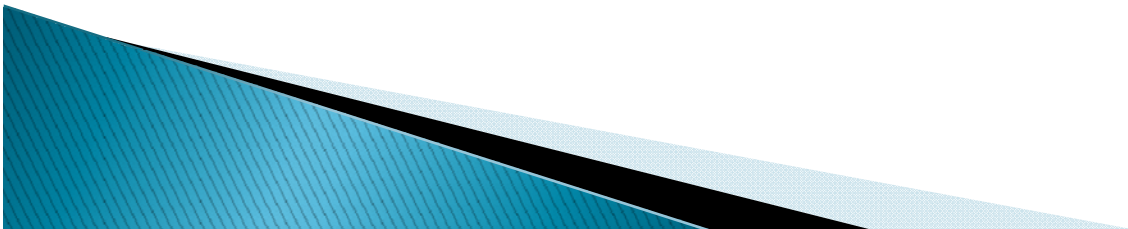


Рассылка MPI_Alltoall

```
int MPI_Alltoall (void* sendbuf,  
    int sendcount, /* in */  
    MPI_Datatype sendtype, /* in */  
    void* recvbuf, /* out */  
    int recvcount, /* in */  
    MPI_Datatype recvtype, /* in */  
    MPI_Comm comm);
```

Описание:

- ▶ Рассылка сообщений от каждого процесса каждому
- ▶ j-ый блок данных из процесса i принимается j-ым процессом и размещается в i-ом блоке буфера recvbuf



Рассылка MPI_Alltoall

