

Проект комиссии Президента  
по модернизации и технологическому развитию экономики России  
«Создание системы подготовки высококвалифицированных кадров  
в области суперкомпьютерных технологий и  
специализированного программного обеспечения»

УТВЕРЖДАЮ  
Председатель экспертного совета  
системы НОЦ СКТ, член-корр. РАН  
В.В. Воеводин

\_\_\_\_\_ 201\_\_ г.  
" \_\_\_\_\_ " \_\_\_\_\_

### **Конспект лекций дисциплины**

«Параллельное программирование для многопроцессорных систем  
с общей и распределенной памятью»

**«010400.62 – Прикладная математика и информатика»**

Разработчики: Лаева В.И. , Трунов А.А.  
Рецензент: проф. Старченко А.В.

**Москва**

## ДВУХТОЧЕЧНЫЙ НЕБЛОКИРУЮЩИЙ ОБМЕН СООБЩЕНИЯМИ

Библиотека MPI содержит набор функций для выполнения асинхронного (без блокировки) обмена сообщениями. В отличие от блокирующих, которые приостанавливают вызвавший их процесс до тех пор, пока операция передачи не будет завершена, неблокирующие подразумевают совмещение операций обмена с другими операциями. Использование неблокирующих функций более безопасно с точки зрения возникновения тупиковых ситуаций, а также может увеличить скорость работы программы за счет совмещения выполнения вычислительных и коммуникационных операций.

Вызов функции неблокирующей передачи инициирует, но не завершает ее. Завершиться выполнение неблокирующей передачи может еще до того, как сообщение будет скопировано в буфер передачи. Для завершения обмена требуется вызов дополнительной функции, которая проверяет, скопированы ли данные в буфер передачи. Таким образом, асинхронный обмен выполняется в два этапа:

- иницирование операции,
- проверка завершения операции.

Для связи между функциями обмена и функциями опроса их завершения неблокирующие операции используют специальный скрытый объект – *идентификатор асинхронной операции обмена* (*request*). Для прикладных программ доступ к идентификатору обмена возможен только через вызовы MPI функций. Если операция обмена завершена, функция проверки снимает «запрос обмена», устанавливая его значение ***MPI\_REQUEST\_NULL***.

### Инициализация неблокирующего обмена

Как было отмечено ранее, асинхронная передача сообщения может выполняться в четырех режимах:

- неблокирующий стандартный (***MPI\_Isend***),
  - неблокирующий синхронный (***MPI\_Issend***),
  - неблокирующий буферизованный (***MPI\_Ibsend***),
  - неблокирующий «по готовности» (***MPI\_Irsend***)
- (буква ***I*** в имени от *immediately* – немедленно).

Функция инициализации стандартной неблокирующей передачи сообщения:

***C:*** ***int MPI\_Isend(void \*buf, int count, MPI\_Datatype datatype, int dest, int tag, MPI\_Comm comm, MPI\_Request \*request)***

***FORTRAN:***

***INTEGER count, dest, tag, comm, error, request***

### ***MPI\_Isend(buf, count, datatype, dest, tag, comm, request, error)***

Функция инициирует посылку сообщения *buf* с идентификатором (тегом) *tag*, состоящего из *count* элементов типа *datatype*, процессу с номером *dest* в области связи с коммуникатором *comm*. Тип передаваемых элементов *datatype* должен указываться с помощью соответствующих типов MPI. Перечисленные параметры являются входными. Выходной параметр *request* – идентификатор запроса на выполнение операции имеет тип *MPI\_Request* в C и тип INTEGER в языке FORTRAN.

Возврат из функции *MPI\_Isend* происходит сразу без ожидания обработки всего сообщения, находящегося в буфере *buf*. Это означает, что нельзя что – то записывать в данный буфер без получения дополнительной информации, подтверждающей завершения пересылки. Определить, когда можно использовать повторно буфер передачи без опасения испортить передаваемое сообщение, можно с помощью параметра *request* и функций проверки выполнения обмена.

Функции *MPI\_Isend*, *MPI\_Issend*, *MPI\_Ibsend*, *MPI\_Irsend* имеют одинаковые параметры.

Функция инициализации неблокирующего приема сообщения:

**C:** *int MPI\_Irecv(void \*buf, int count, MPI\_Datatype datatype, int source, int tag, MPI\_Comm comm, MPI\_Request \*request)*

**FORTRAN:**

*INTEGER count, source, tag, comm, error, request*

*MPI\_Irecv(buf, count, datatype, source, tag, comm, request, error)*

Функция инициирует прием сообщения *buf* (выходной параметр) с идентификатором *tag* от процесса *source*. Число элементов в принимаемом сообщении не должно превосходить значения *count*. Тип получаемых данных *datatype* должен указываться с помощью соответствующих типов MPI. Выходной параметр *request* – идентификатор запроса на выполнение операции.

Возврат из функции *MPI\_Irecv* происходит сразу без ожидания получения и записи всего сообщения в буфере *buf*. Окончание процесса приема можно определить с помощью параметра *request* и функций проверки выполнения обмена.

Сообщение, отправленное функцией *MPI\_Send*, *MPI\_Isend* и любой из трех их модификаций, может быть принято любой функцией *MPI\_Recv*, *MPI\_Irecv*.

### **Функции проверки выполнения обмена**

Проверка фактического выполнения передачи/приема сообщения в неблокирующем режиме осуществляется с помощью вызова функций проверки. Блокирующие функции проверки приостанавливают работу

процесса до завершения операции обмена (*MPI\_Wait*, *MPI\_Waitall*, *MPI\_Waitany*, *MPI\_Waitsome*). Неблокирующие функции проверки возвращают логическое значение «истина», если операция выполнена (*MPI\_Test*, *MPI\_Testall*, *MPI\_Testany*), а функция *MPI\_Testsome* возвращает количество завершившихся обменов.

Функция проверки с блокировкой:

**C:** *int MPI\_Wait (MPI\_Request \*request, MPI\_Status \*status)*

**FORTRAN:**

*INTEGER status(MPI\_STATUS\_SIZE), request, error*

*MPI\_Wait (request, status, error)*

Функция блокирует работу процесса до тех пор, пока не будет завершен обмен с запросом *request*. Параметр *request* – входной и выходной, после завершения этой функции, он становится равным *MPI\_REQUEST\_NULL*. Выходной параметр *status* возвращает информацию о выполненной операции.

Функция проверки без блокировки:

**C:** *int MPI\_Test(MPI\_Request \*request, int \*flag, MPI\_Status \*status)*

**FORTRAN:**

*LOGICAL flag*

*INTEGER request, error, status(MPI\_STATUS\_SIZE)*

*MPI\_Test(request, flag, status, error)*

Входной параметр :

*request* – идентификатор запроса на выполнения операции асинхронного обмена.

Выходные параметры:

*flag* – признак завершенности операции обмена (если *flag=true*, то операция обмена завершилась успешно),

*status*– статус (информация о полученном сообщении),

*error* – код ошибки.

Когда сразу несколько процессов обмениваются сообщениями, можно использовать проверки, которые применяются одновременно к нескольким обменам:

*MPI\_Waitall*, *MPI\_Testall* – завершились все операции обмена,

*MPI\_Waitany*, *MPI\_Testany* – завершилась по крайней мере одна операция обмена,

*MPI\_Waitsome*, *MPI\_Testsome* – завершилась одна из списка проверяемых операций обмена.

Функция блокирующей проверки завершения всех обменов:

**C:** *int MPI\_Waitall (int count, MPI\_Request \*requests[], MPI\_Status \*statuses[])*

**FORTRAN:**

***INTEGER count, requests(count), error, statuses(count)***

***MPI\_Waitall(count, requests, statuses, error)***

**Входные параметры:**

***count*** – количество запросов на обмен (размерность массивов ***requests*** и ***statuses*** );

***requests*** – массив идентификаторов операции асинхронного обмена.

**Выходные параметры:**

***statuses*** – массив статусов сообщений,

***error*** – код ошибки.

В результате выполнения функции ***MPI\_Waitall*** аннулируются запросы, сформированные неблокирующими операциями обмена, и соответствующим элементам массива ***requests*** присваивается значение ***MPI\_REQUEST\_NULL***. Если операция выполнена успешно, полю (элементу массива) ошибки статуса присваивается ***MPI\_SUCCESS***. В случае неуспешного выполнения одной или более операций возвращается код ошибки ***MPI\_ERR\_IN\_STATUS*** и полю присваивается значение кода ошибки соответствующей операции. Если операция не выполнена, но и не было ошибки, полю ошибки присваивается ***MPI\_ERR\_PENDING*** (это соответствует наличию запросов на выполнение операции обмена, ожидающих обработки).

Функция неблокирующей проверки завершения всех обменов:

***C: int MPI\_Testall (int count, MPI\_Request \*requests[], int \*flag, MPI\_Status \*statuses[])***

**FORTRAN:**

***INTEGER count, requests(count), error, statuses(count)***

***LOGICAL flag***

***MPI\_Testall(count, requests, flag, statuses, error)***

Смысл и значение параметров ***count***, ***requests*** и ***statuses*** тот же, что и для функции ***MPI\_Waitall***. Выходной параметр ***flag=true***, если все обмены, связанные с активными запросами в массиве ***requests*** выполнены. Если не все обмены завершены, то ***flag=false***, а массив ***statuses*** не определен.

Функция блокирующей проверки завершения любого числа обменов:

***C: int MPI\_Waitany(int count, MPI\_Request \*requests[], int \*index, MPI\_Status \*statuses[])***

**FORTRAN:**

***INTEGER count, requests(count), error, statuses(count), index***

***MPI\_Waitany(count, requests, index, statuses, error)***

Смысл и значение параметров ***count***, ***requests*** и ***statuses*** тот же, что и для функции ***MPI\_Waitall***. Выходной параметр ***index*** – номер элемента в массиве ***requests***, содержащего идентификатор завершенной операции.

Если в списке вообще нет активных запросов или он пуст, вызовы завершаются сразу со значением индекса ***MPI\_UNDEFINED***.

Функция блокирует работу процесса до тех пор, пока, по крайней мере, один обмен из массива запросов ***requests***, не будет завершен.

Функция *неблокирующей проверки завершения любого числа обменов*:

***C: int MPI\_Testany(int count, MPI\_Request \*requests[], int \*index, int \*flag, MPI\_Status \*statuses[])***

***FORTRAN:***

***INTEGER count, requests(count), error, statuses(count), index***

***LOGICAL flag***

***MPI\_Testany(count, requests, index, flag, statuses, error)***

Смысл и значение параметров ***count***, ***index***, ***requests*** и ***statuses*** тот же, что и для функции ***MPI\_Waitany***, но выходной параметр ***flag=true***, если одна из операций завершена.

Функция *блокирующей проверки завершения обменов из указанного списка*:

***C: int MPI\_Waitsome(int count, MPI\_Request \*requests[], int \*outcount, int \*indexes[], MPI\_Status \*statuses[])***

***FORTRAN:***

***INTEGER count, requests(count), statuses(count), indexes(count), error, outcount***

***MPI\_Waitsome(count, requests, outcount, indexes, statuses, error)***

Смысл и значение параметров ***count***, ***requests*** и ***statuses*** тот же, что и для функции ***MPI\_Waitany***. Выходные параметры: ***outcount*** – количество выполненных запросов из массива ***requests***, в первых ***outcount*** элементах массива ***indexes*** – возвращаются индексы этих операций, в первых ***outcount*** элементах массива ***statuses*** – возвращается статус завершенных операций. Если в списке нет активных запросов, выполнение функции завершается сразу, а параметру ***outcount*** присваивается значение ***MPI\_UNDEFINED***.

Функция *неблокирующей проверки завершения обменов из указанного списка*:

***C: int MPI\_Testsome(int count, MPI\_Request \*requests[], int \*outcount, int \*indexes[], MPI\_Status \*statuses[])***

***FORTRAN:***

***INTEGER count, requests(count), statuses(count), indexes(count), outcount, error***

***MPI\_Testsome(count, requests, outcount, indexes, statuses, error)***

Функция работает так же, как *MPI\_Waitsome*, за исключением того, что возврат из функции происходит немедленно. Если ни одна из указанных операций не завершилась, то *outcount=0*.

В функциях *MPI\_Waitany* и *MPI\_Testany* обмен из числа завершенных выбирается произвольно и для него возвращается статус. В *MPI\_Waitsome* и *MPI\_Testsome* статус возвращается для всех завершенных обменов, поэтому эти функции можно использовать для того, чтобы определить количество завершенных обменов.

**Пример** программы обмена сообщениями по «кольцу» с использованием неблокирующих функций обмена.

#### Program Example

Implicit none

Include 'mpif.h'

Integer Rank,Size,Comm,Ierr,n,tag,k

Parameter (n=2)

Integer next,secc,requests(n),statuses(n)

Real sbuf,rbuf

tag=0

Call MPI\_Init(Ierr)

comm=MPI\_COMM\_WORLD

Call MPI\_Comm\_size(Comm,size,Ierr)

Call MPI\_Comm\_rank(Comm,Rank,Ierr)

next=rank+1

secc=rank-1

if(rank==0) secc=size-1

if(rank==size-1) next=0

Call Random\_seed()

do k=1,rank+1

Call Random\_number(sbuf)

enddo

Call MPI\_Irecv(rbuf,1,MPI\_REAL,secc,tag,comm,requests(1),Ierr)

Call MPI\_Isend(sbuf,1,MPI\_REAL,next,tag,comm,requests(2),Ierr)

write(\*,\*) 'process=',rank,' send sbuf=',sbuf

write(\*,\*) 'process=',rank,' before Wait recv rbuf=',rbuf

Call MPI\_Waitall(n,requests,statuses,Ierr)

write(\*,\*) 'process=',rank,' after Wait recv rbuf=',rbuf

Call MPI\_Finalize(Ierr)

end

### Объединение запросов на взаимодействие

MPI позволяет для неблокирующих операций формировать пакеты запросов на коммуникационные операции. Несколько запросов на обмен можно объединить вместе для того, чтобы далее их можно было запустить одной командой. Это позволяет минимизировать накладные расходы на организацию связи между процессором и сетевым контроллером.

Для формирования отложенного запроса на передачу сообщения можно использовать функции: *MPI\_Send\_init*, *MPI\_Ssend\_init*,

*MPI\_Bsend\_init, MPI\_Rsend\_init.* Сама операция пересылки при вызове этих функций не начинается.

Функция *формирование запроса на выполнение стандартной передачи сообщения:*

**C:** *int MPI\_Send\_init(void \*buf, int count, MPI\_Datatype datatype, int dest, int tag, MPI\_Comm comm, MPI\_Request request)*

**FORTRAN:**

*INTEGER count, datatype, dest, tag, comm, error,request*

*MPI\_Send\_init(buf, count, datatype, dest, tag, comm, request,error)*

Все параметры этой функции такие же, как у функции *MPI\_Isend*.

Функция *формирование запроса на выполнение приема сообщения:*

**C:** *int MPI\_Recv\_init(void \*buf, int count, MPI\_Datatype datatype, int source, int tag, MPI\_Comm comm, MPI\_Request request)*

**FORTRAN:**

*INTEGER count, datatype, source, tag, comm, error,request*

*MPI\_Recv\_init(buf, count, datatype, source, tag, comm, request,error)*

Все параметры этой функции такие же, как у функции *MPI\_Irecv*. Функция формирует отложенный запрос на прием сообщения, но операцию приема не выполняет.

Отложенный обмен сообщениями инициализируется вызовом специальных функций.

Функция *инициализации отложенного обмена:*

**C:** *int MPI\_Start(MPI\_Request \*request)*

**FORTRAN:**

*INTEGER error,request*

*MPI\_Start (request,error)*

Входным параметром *request* является идентификатор запроса на выполнение операции обмена. Вызов *MPI\_Start* с запросом на обмен, созданным *MPI\_Send\_init*, инициализирует обмен с теми же свойствами, что и вызов *MPI\_Isend, MPI\_Bsend\_init* –вызову *MPI\_Ibsend* и т.д.

Функция *инициализации всех отложенных обменов:*

**C:** *int MPI\_Startall(int count,MPI\_Request \*requests)*

**FORTRAN:**

*INTEGER count,error,requests(count)*

*MPI\_Startall (count,requests,error)*

Функция инициализирует *count* отложенных запросов на обмен, соответствующих значениям элементов массива *requests*.



Сообщение, отправленное с помощью объединения запросов или обычным способом, может быть принято обычным способом или с помощью объединения запросов.

Значение параметра *request* после завершения коммуникационной операции, запущенной при помощи отложенного запроса, *сохраняется и может быть использовано* в дальнейшем.

### **Функции отмены «ждущих» асинхронных обменов**

Аннулировать неблокирующие ожидающие обработки («ждущие») коммуникационные операции можно с помощью функций *MPI\_Request\_free* и *MPI\_Cancell*.

Функция *снятия запроса без ожидания завершения неблокирующей операции*:

**C:** *int MPI\_Request\_free(MPI\_Request \*request)*

**FORTTRAN:**

*INTEGER error,request*

*MPI\_Request\_free(request,error)*

При вызове эта функция помечает запрос на обмен *request* для удаления и присваивает ему значение *MPI\_REQUEST\_NULL*. Операция обмена, связанная с этим запросом завершается, а сам запрос удаляется только после завершения обмена.

Функция *аннулирования неблокирующего обмена*:

**C:** *int MPI\_Cancel(MPI\_Request \*request)*

**FORTTRAN:**

*INTEGER error,request*

*MPI\_Cancel(request,error)*

Функция отменяет неблокирующий обмен с запросом *request*. Вызов этой функции завершается сразу, возможно еще до реального аннулирования обмена, поэтому этот обмен следует завершить с помощью функций проверки выполнения асинхронного обмена. После вызова *MPI\_Cancel* и следующего за ним *MPI\_Wait* или *MPI\_Test*, запрос на выполнение операции обмена становится неактивным и может быть активирован для нового обмена. *MPI\_Cancel* может быть использована для отмены отложенных и обычных обменов.

Функция *проверки аннулирования неблокирующего обмена*:

**C:** *int MPI\_Test\_cancelled(MPI\_Status \*status, int \*flag)*

**FORTTRAN:**

*INTEGER error, status(MPI\_STATUS\_SIZE)*

*LOGICAL flag*

*MPI\_Test\_cancelled(status, flag, error)*

Функция возвращает значение *flag=true*, если обмен с указанным статусом *status* успешно аннулирован.

Пример программы с функциями *MPI\_Send\_init* и *MPI\_Recv\_init*.

```
Program Example
Implicit none
Include 'mpif.h'
Integer Rank,Size,Comm,Ierr,n,tag,m,k,i
Parameter (n=2,m=3)
Integer next,secc,rqs(2),sts(2)
Real sbuf,rbuf
tag=0
Call MPI_Init(Ierr)
comm=MPI_COMM_WORLD
Call MPI_Comm_size(Comm,size,Ierr)
Call MPI_Comm_rank(Comm,Rank,Ierr)
next=rank+1
secc=rank-1
if(rank==0) secc=size-1
if(rank==size-1) next=0
Call MPI_Recv_init(rbuf,1,MPI_REAL,secc,tag,comm,rqs(1),Ierr)
Call MPI_Send_init(sbuf,1,MPI_REAL,next,tag,comm,rqs(2),Ierr)
Call Random_seed()
do k=1,m
do i=1,rank+k
Call Random_number(sbuf)
enddo
Call MPI_Startall(n,rqs,Ierr)
write(*,*) 'process=',rank,' k= ',k,' send sbuf=',sbuf
Call MPI_Waitall(n,rqs,sts,Ierr)
write(*,*) 'process=',rank,' k=',k,
&      ' after Wait recv rbuf=',rbuf
enddo
call MPI_Request_free(rqs(1),Ierr)
call MPI_Request_free(rqs(2),Ierr)
Call MPI_Finalize(Ierr)
end
```