

# **Коллективные операции передачи данных**

# Коллективные операции передачи данных...

Под *коллективными операциями* в MPI понимаются операции данных, в которых принимают участие все процессы используемого коммутатора. Причем гарантировано, что эти операции будут выполняться гораздо эффективнее, поскольку MPI-функция реализована с использованием внутренних возможностей коммуникационной среды.

# Коллективные операции передачи данных...

## Обобщенная передача данных от одного процесса всем процессам...

- *Распределение данных* – ведущий процесс (*root*) передает процессам различающиеся данные

```
int MPI_Scatter(void *sbuf, int scount, MPI_Datatype stype,  
                void *rbuf, int rcount, MPI_Datatype rtype,  
                int root, MPI_Comm comm),
```

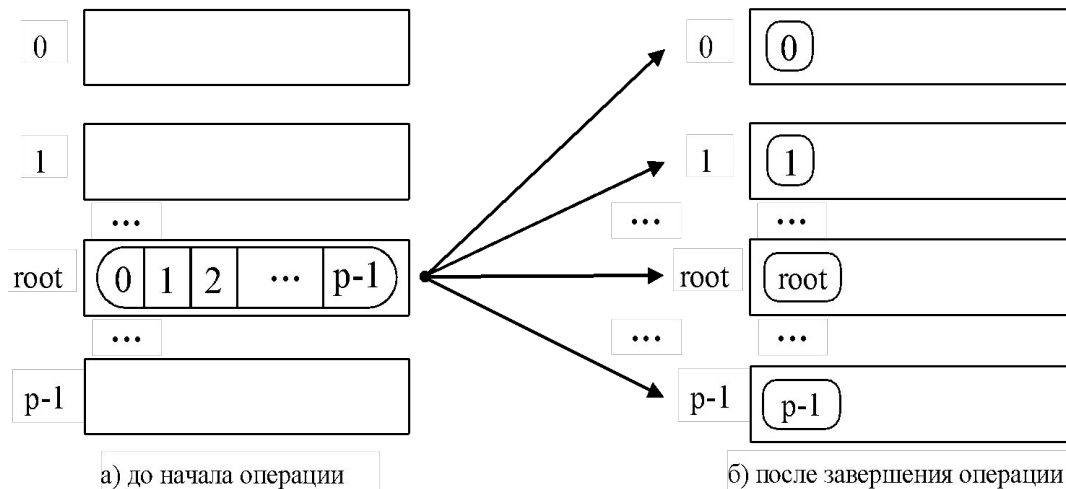
где

- **sbuf**, **scount**, **stype** – параметры передаваемого сообщения (**scount**,  
– определяет количество элементов, передаваемых на каждый процесс),
- **rbuf**, **rcount**, **rtype** – параметры сообщения, принимаемого в процессах,
- **root** – ранг процесса, выполняющего рассылку данных,
- **comm** – коммутатор, в рамках которого выполняется передача данных.

# Коллективные операции передачи данных...

## Обобщенная передача данных от одного процесса всем процессам

- Вызов *MPI\_Scatter* при выполнении рассылки данных должен быть обеспечен в каждом процессе коммутатора,
- *MPI\_Scatter* передает всем процессам сообщения одинакового размера. Если размеры сообщений для процессов могут быть разными, следует использовать функцию *MPI\_Scatterv*.

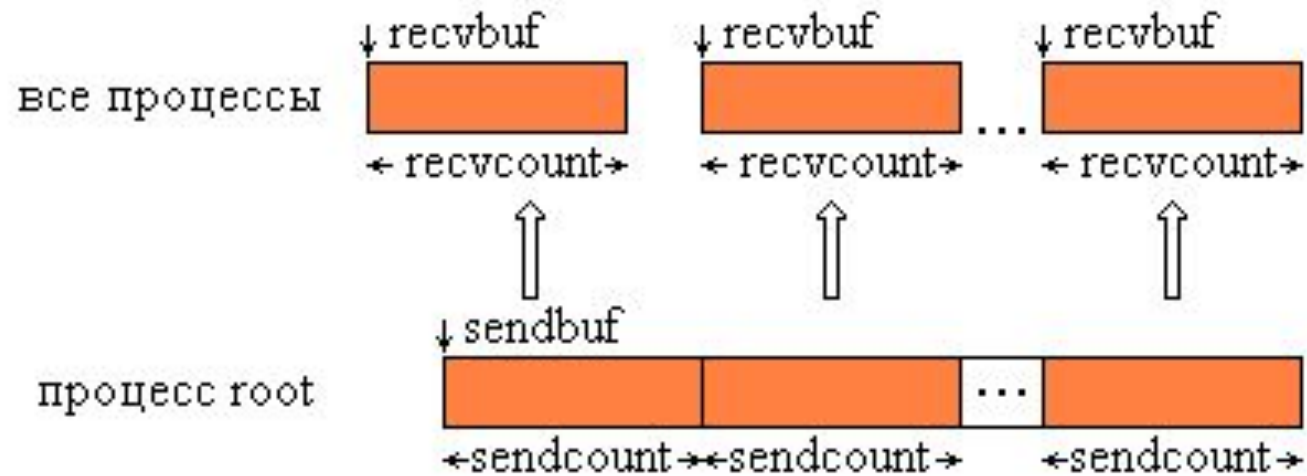


Основы параллельных вычислений:

*Моделирование и анализ  
параллельных вычислений*

© Гергель В. П.

# Графическая интерпретация операции Scatter



# Рассылка различного количества данных

```
int MPI_Scatterv(void *sbuf, int *scounts, int *displs, MPI_Datatype stype,  
                void *rbuf, int rcount, MPI_Datatype rtype, int root, MPI_Comm comm)
```

**sbuf** - адрес рассылаемого массива данных.

Начало рассылаемых порций задает массив **displs**, количество элементов в порции задает массив **scounts**.

**scounts** – целочисленный массив, содержащий количество элементов, передаваемых каждому процессу (индекс равен рангу адресата, длина равна числу процессов в коммутаторе).

**displs** – целочисленный массив, содержащий смещения относительно начала массива sbuf (индекс равен рангу адресата, длина равна числу процессов в коммутаторе).

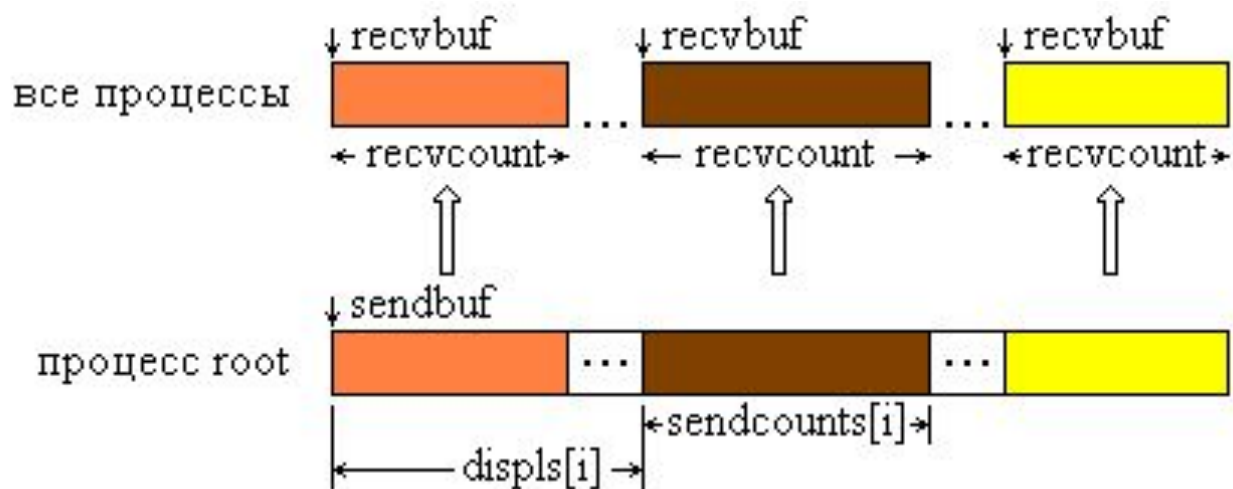
**rbuf** - адрес массива, принимающего порцию данных в i-ом процессе.

**rcount** - размер порции, принимаемой в ранге адресата .

**root** - ранг процесса, выполняющего рассылку данных.

**stype, rtype** – типы рассылаемых и принимаемых данных

# Графическая интерпретация операции Scatterv



# Определение частей массива в rank's

```
...
count=m / size;  ost=m % size;
/* Calculating parts of array for root rank */
if (rank==0) /* Process 0 - master */ {
/* Creation auxiliary arrays for data communication */
displs = (int *)malloc(size * sizeof(int));
rcounts = (int *)malloc(size * sizeof(int));
for(i=0;i < size;i++){
scol = i < ost ? count+1 : count;          rcounts[i] = scol;
nach = i*scol + (i >= ost ? ost : 0);      displs[i] = nach;
}
} /* End of work process 0 */
...
/* Calculating parts of vector vA for rank in others processes */
scol = rank < ost ? count+1 : count;
/* Offset (in strings) part for rank in vector vA */
nach = rank*scol + (rank >= ost ? ost : 0);
```



# Коллективные операции передачи данных...

## Обобщенная передача данных от всех процессов одному процессу...

- Передача данных от всех процессоров одному процессу (*сбор данных*) является обратной к операции распределения данных

```
int MPI_Gather(void *sbuf, int scount, MPI_Datatype stype,  
               void *rbuf, int rcount, MPI_Datatype rtype,  
               int root, MPI_Comm comm),
```

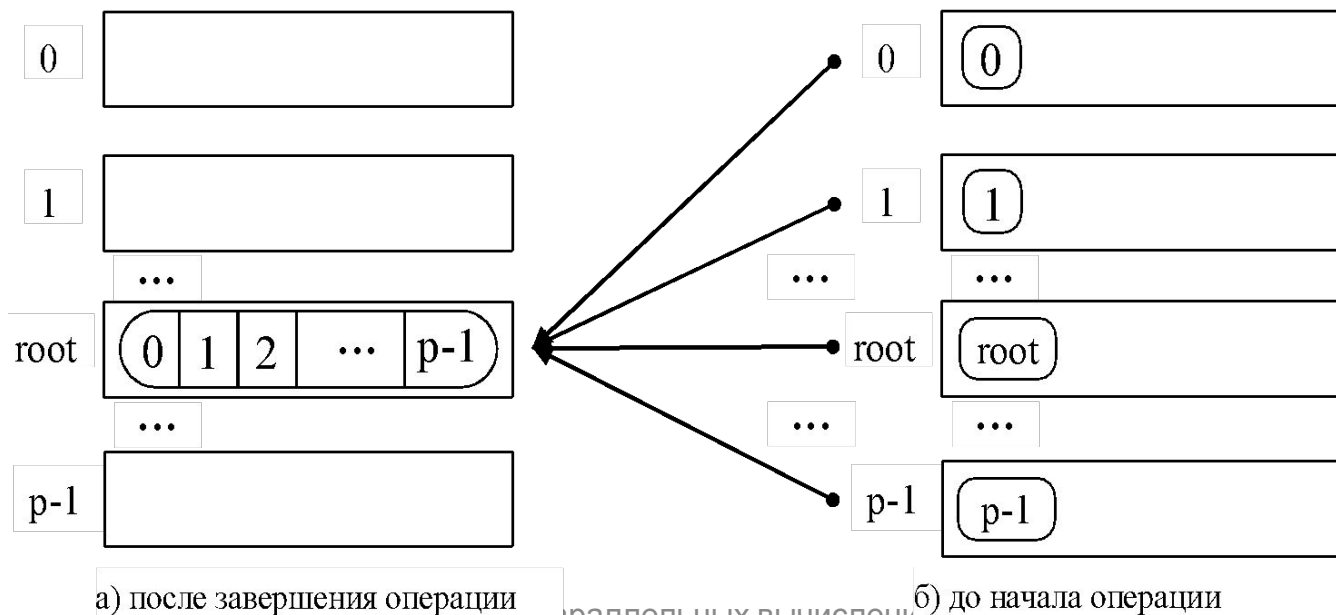
где

- **sbuf**, **scount**, **stype** - параметры передаваемого сообщения,
- **rbuf**, **rcount**, **rtype** - параметры принимаемого сообщения,
- **root** - ранг процесса, выполняющего сбор данных,
- **comm** - коммуникатор, в рамках которого выполняется передача данных.

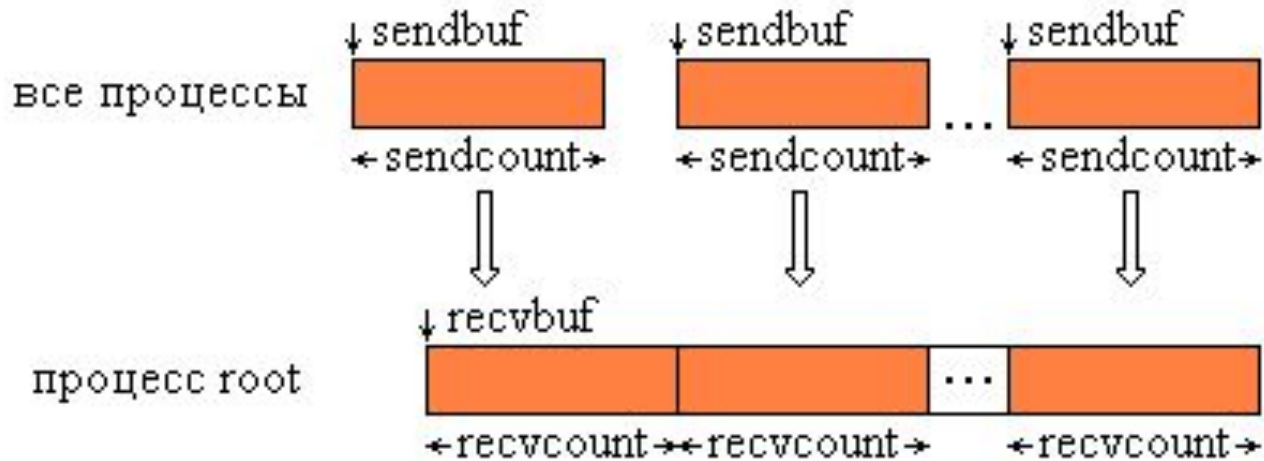
# Коллективные операции передачи данных...

## Обобщенная передача данных от всех процессов одному процессу...

- *MPI\_Gather* определяет коллективную операцию, и ее вызов при выполнении сбора данных должен быть обеспечен в каждом процессе коммутатора



# Графическая интерпретация операции Gather



# Сборка различного количества данных со всех процессов в один процесс

```
int MPI_Gatherv(void *sbuf, int scount, MPI_Datatype stype,  
               void *rbuf, int *rcounts, int *displs, MPI_Datatype rtype, int root,  
               MPI_Comm comm)
```

**sbuf** – адрес рассылаемых данных в процессе отправителе.

**scount** – число рассылаемых данных.

**rbuf** – адрес буфера в принимающем процессе

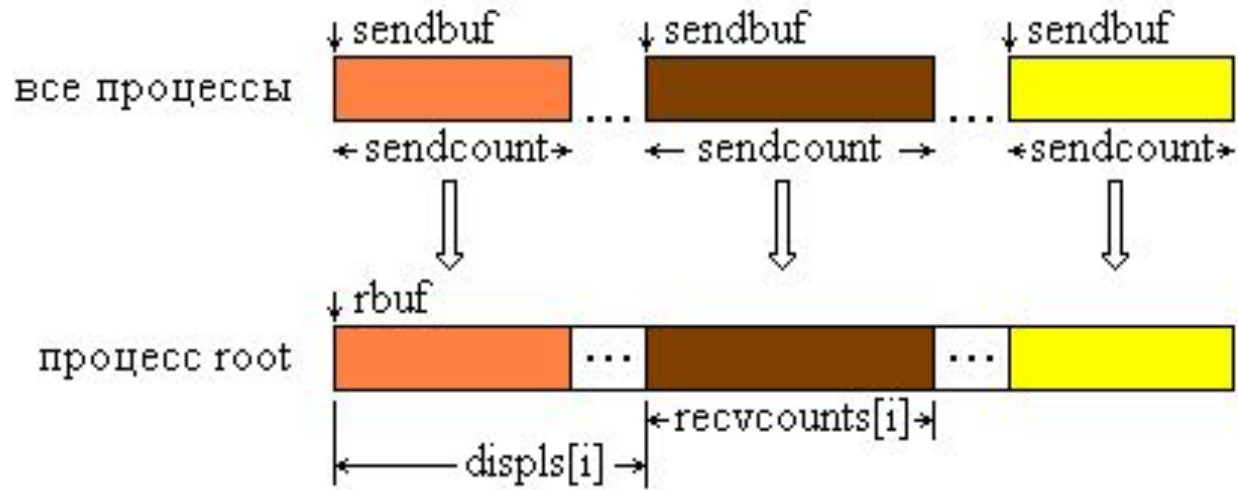
**scounts** – целочисленный массив, содержащий количество элементов, принимаемых от каждого процесса (индекс равен рангу процесса, длина равна числу процессов в коммутаторе).

**displs** – целочисленный массив, смещений относительно начала массива rbuf (индекс равен рангу адресата, длина равна числу процессов в коммутаторе).

**root** – ранг принимающего процесса.

**stype, rtype** – типы рассылаемых и принимаемых данных.

# Графическая интерпретация операции Gatherv

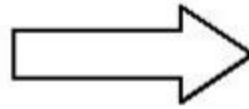


## данные

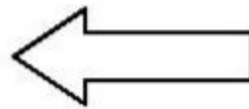
процессы

$A_0$	$A_1$	$A_2$	$A_3$	$A_4$	$A_5$

scatter



gather



$A_0$					
$A_1$					
$A_2$					
$A_3$					
$A_4$					
$A_5$					

# Коллективные операции передачи

## данных...

**Обобщенная передача данных от всех процессов одному процессу**

- *MPI\_Gather* собирает данные на одном процессе. Для получения всех собираемых данных на каждом процессе нужно использовать *функцию сбора и*

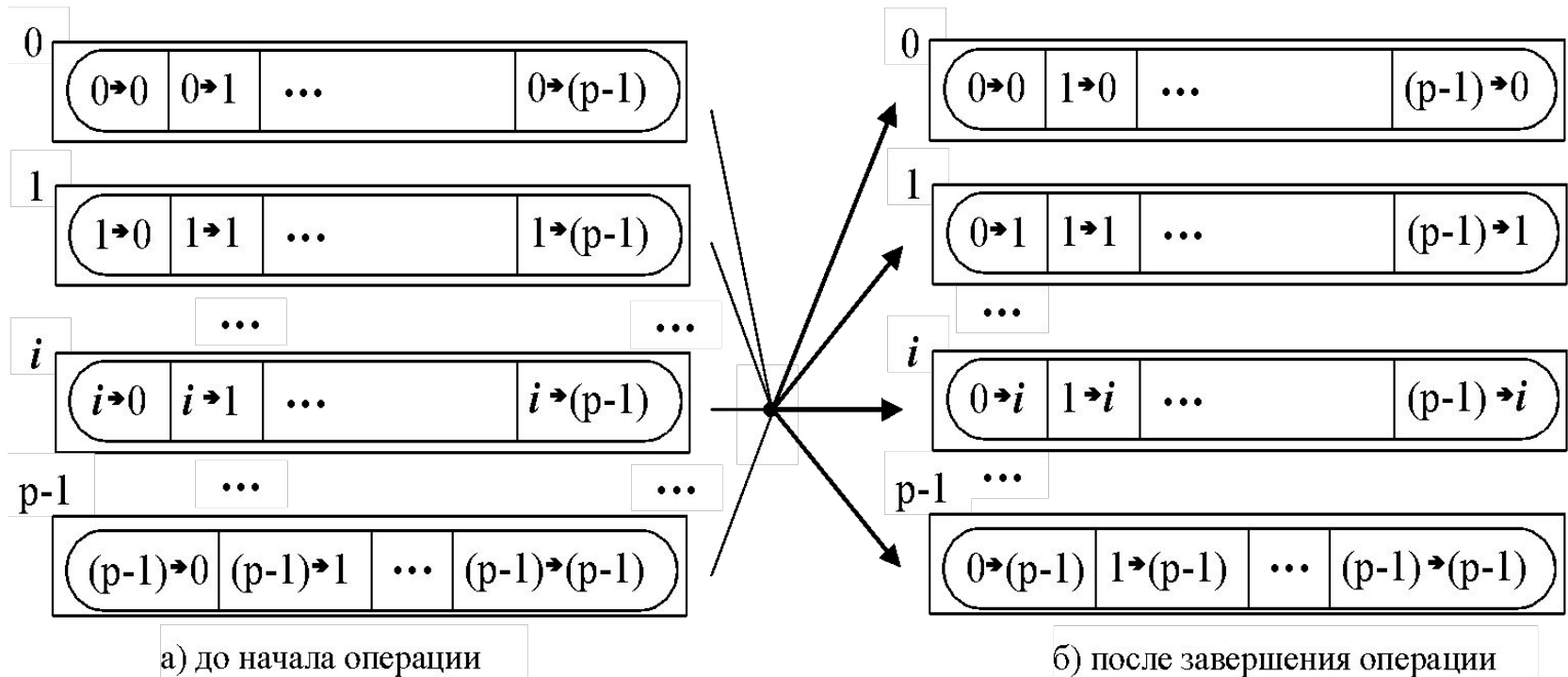
```
int MPI_Allgather(void *sbuf, int scount, MPI_Datatype
stype,
    void *rbuf, int rcount, MPI_Datatype rtype, MPI_Comm
comm) .
```

- В случае, когда размеры передаваемых процессами сообщений могут быть различны, для передачи данных необходимо использовать функции *MPI\_Gatherv* и *MPI\_Allgatherv*.

# Коллективные операции передачи данных...

## данных...

Общая передача данных от всех процессов всем процессам...



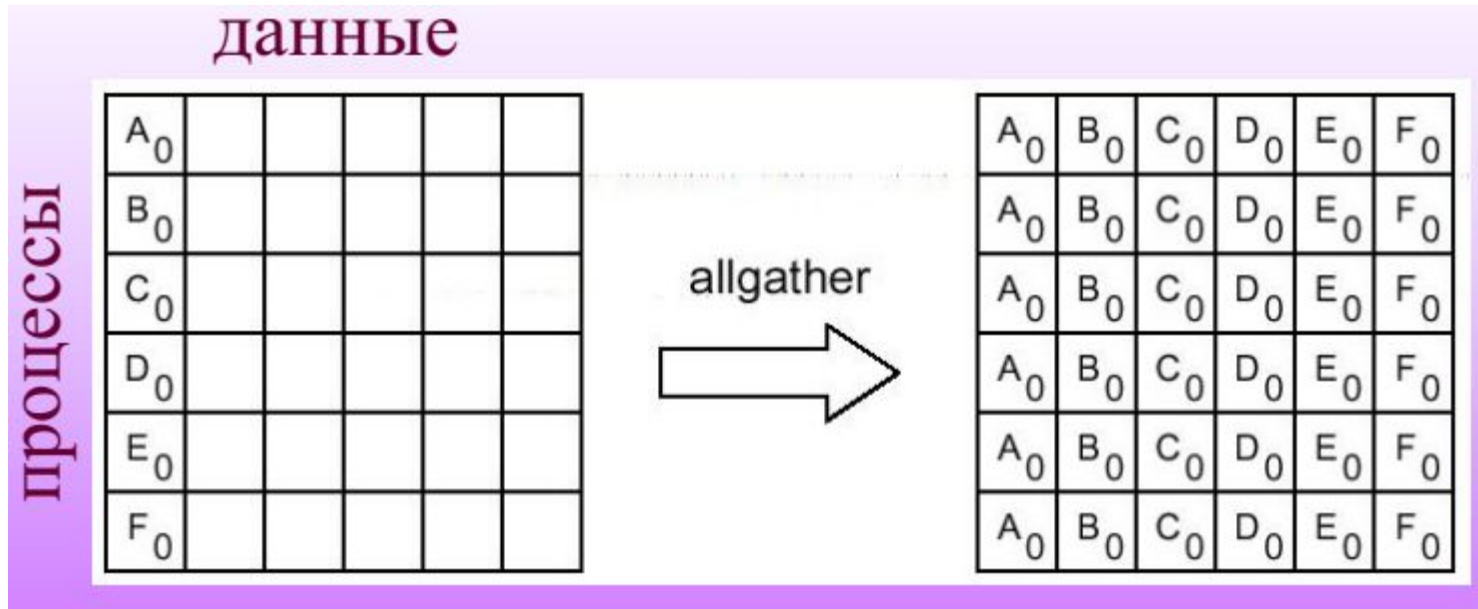
Основы параллельных вычислений:

*Моделирование и анализ  
параллельных вычислений*

© Гергель В. П.



# Коллективные операции передачи данных...



## Коллективные операции передачи данных...

```
int MPI_Alltoall(void *sbuf, int scount, MPI_Datatype stype,
void *rbuf, int rcount, MPI_Datatype rtype, MPI_Comm comm)
```

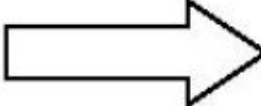
Рассылка каждым процессом коммуникатора **comm** различных порций данных всем другим процессам. **j**-й блок массива **sbuf** процесса **i** попадает в **i**-й блок массива **rbuf** процесса **j**.

данные

процессы

A <sub>0</sub>	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
B <sub>0</sub>	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	B <sub>4</sub>	B <sub>5</sub>
C <sub>0</sub>	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	C <sub>4</sub>	C <sub>5</sub>
D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>	D <sub>4</sub>	D <sub>5</sub>
E <sub>0</sub>	E <sub>1</sub>	E <sub>2</sub>	E <sub>3</sub>	E <sub>4</sub>	E <sub>5</sub>
F <sub>0</sub>	F <sub>1</sub>	F <sub>2</sub>	F <sub>3</sub>	F <sub>4</sub>	F <sub>5</sub>

alltoall



A <sub>0</sub>	B <sub>0</sub>	C <sub>0</sub>	D <sub>0</sub>	E <sub>0</sub>	F <sub>0</sub>
A <sub>1</sub>	B <sub>1</sub>	C <sub>1</sub>	D <sub>1</sub>	E <sub>1</sub>	F <sub>1</sub>
A <sub>2</sub>	B <sub>2</sub>	C <sub>2</sub>	D <sub>2</sub>	E <sub>2</sub>	F <sub>2</sub>
A <sub>3</sub>	B <sub>3</sub>	C <sub>3</sub>	D <sub>3</sub>	E <sub>3</sub>	F <sub>3</sub>
A <sub>4</sub>	B <sub>4</sub>	C <sub>4</sub>	D <sub>4</sub>	E <sub>4</sub>	F <sub>4</sub>
A <sub>5</sub>	B <sub>5</sub>	C <sub>5</sub>	D <sub>5</sub>	E <sub>5</sub>	F <sub>5</sub>

# Коллективные операции передачи

## данных...

### Дополнительные операции редукции данных...

- *MPI\_Reduce* обеспечивает получение результатов редукции данных только на одном процессе,
- Функция *MPI\_Allreduce* редукции и рассылки выполняет рассылку между процессами всех результатов операции редукции:

```
int MPI_Allreduce(void *sendbuf, void *recvbuf, int count,  
                 MPI_Datatype type, MPI_Op op, MPI_Comm comm)
```

- Возможность управления распределением этих данных между процессами предоставляется функций *MPI\_Reduce\_scatter*,
- Функция *MPI\_Scan* производит операцию сбора и обработки данных, при которой обеспечивается получение и всех частичных результатов редуцирования

```
int MPI_Scan(void *sendbuf, void *recvbuf, int count,  
            MPI_Datatype type, MPI_Op op, MPI_Comm comm)
```

Основы параллельных вычислений:

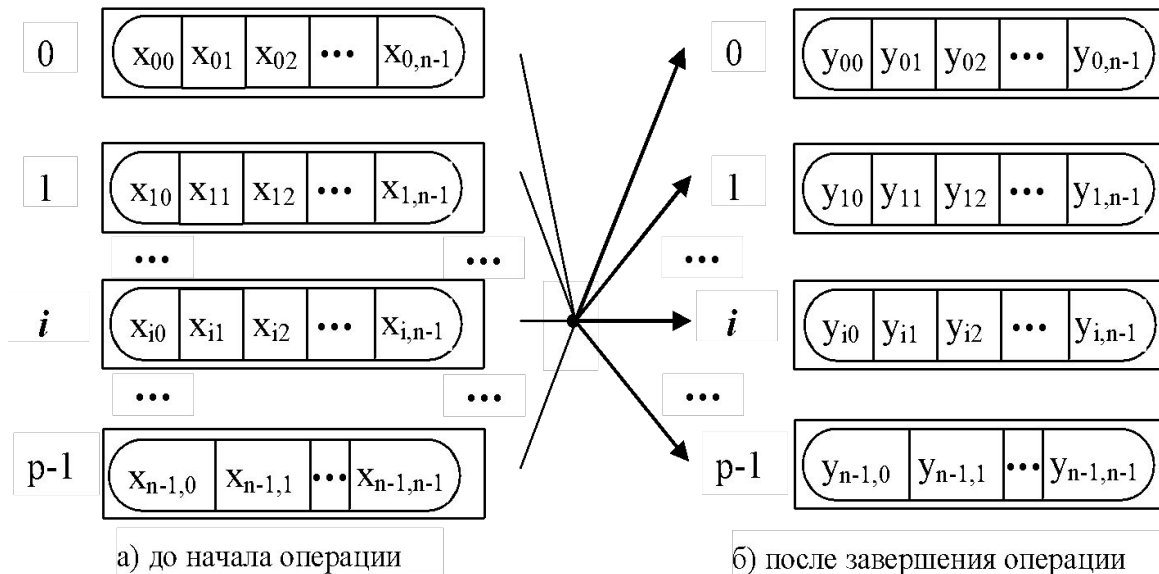
*Моделирование и анализ  
параллельных вычислений*

© Гергель В.П.

# Коллективные операции передачи данных

## Дополнительные операции редукции данных

- При выполнении функции *MPI\_Scan* элементы получаемых сообщений представляют собой результаты обработки соответствующих элементов передаваемых процессами сообщений, при этом для получения результатов на процессе с рангом  $i$ ,  $0 \leq i < n$ , используются данные от процессов, ранг которых меньше или равен  $i$



Основы параллельных вычислений:

*Моделирование и анализ  
параллельных вычислений*

© Гергель В. П.

# Графическая интерпретация операции Scan

