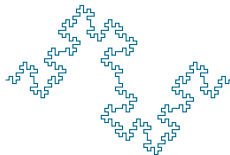


Programowanie współbieżne

Wykład 12

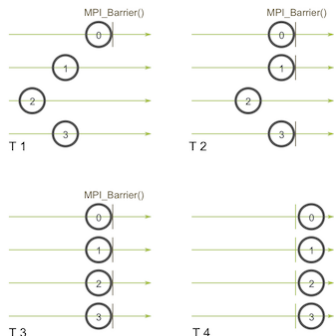
MPI c.d.

Rafał Skinderowicz



MPI_BARRIER

Do synchronizacji procesów w MPI służy funkcja:
 MPI_Barrier(MPI_Comm komunikator)



Źródło: <http://www.mpitutorial.com/>

Jawna synchronizacja stosowana jest zazwyczaj do podziału wykonania programu na sekcje.

MPI_BCAST

W MPI rozgłaszanie realizowane jest za pomocą funkcji:
`MPI_Bcast(void *data, int count, MPI_Datatype datatype, int root, MPI_Comm comm)`

Zarówno proces nadawca, jak i procesy odbiorcy wywołują tę samą metodę, różnica jest w **znaczeniu** parametrów

- ▶ `data`
 - ▶ w procesie–nadawcy wsk. bufora z danymi do wysłania
 - ▶ w procesie–odbiorcy wsk. bufora na rozgłaszane dane
- ▶ `count` – rozmiar danych do wysłania
- ▶ `datatype` – typ wysyłanych danych, np. `MPI_DOUBLE`
- ▶ `root` – identyfikator (ranga) procesu rozgłaszającego
- ▶ `comm` – komunikator, może być `MPI_COMM_WORLD`

Wykonanie rozgłaszania powoduje **niejawną** synchronizację procesów

MPI_BCAST – UWAGI

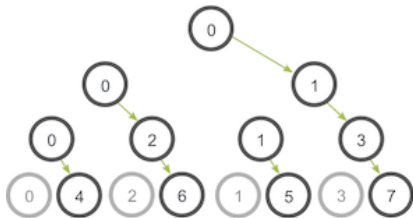
Wykonanie funkcji MPI_Bcast jest „w przybliżeniu” równoważne wykonaniu poniższej „naiwnej” implementacji:

```
void my_bcast(void* data, int count, MPI_Datatype datatype,
              int root, MPI_Comm comm) {
    int proc_rank, proc_num;
    MPI_Comm_rank(comm, &proc_rank);
    MPI_Comm_size(comm, &proc_num);

    if (proc_rank == root) { // proces rozgłaszający
        for (int i = 0; i < proc_num; i++) {
            if (i != proc_rank) {
                MPI_Send(data, count, datatype, i, 0, communicator);
            }
        }
    } else { // proces odbiorca, odbierz dane od nadawcy (root)
        MPI_Recv(data, count, datatype, root, 0, communicator,
                 MPI_STATUS_IGNORE);
    }
}
```

MPI_BCAST – UWAGI

- ▶ Przedstawiona implementacja jest bardzo nieefektywna – dla n procesów proces–nadawca musi wysłać $n - 1$ komunikatów
- ▶ Na szczęście można zrobić to znacznie szybciej, zgodnie z porządkiem „drzewiastym”



Źródło: <http://www.mpitutorial.com/>

- ▶ Teraz pojedynczy proces musi wysłać, co najwyżej $O(\log n)$ komunikatów

MPI_BCAST – PRZYKŁAD

Rozgłaszanie dla macierzy dynamicznej najłatwiej wykonać umieszczając ją w spójnym obszarze pamięci.

```
const int N = 10;
int *bufor = new int[N*N];
if (proc_id == 0) {
    int **m = new int*[N]; // inicjuj macierz
    for (int i = 0; i < N; i++) {
        m[i] = bufor + i * N;
        for (int j = 0; j < N; j++)
            m[i][j] = (i+1) * (j+1); // przykł. zawartość
    }
}
MPI_Bcast(bufor, N*N, MPI_INT, 0, MPI_COMM_WORLD);
if (proc_id != 0) {
    int **m = new int*[N];
    // poskładaj macierz
    for (int i = 0; i < N; i++)
        m[i] = bufor + i * N;
    // teraz można wykonywać operacje na elementach m[i][j]
}
```

MPI_SCATTER

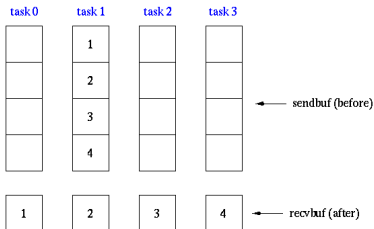
- Czasem zachodzi potrzeba wysłania do wszystkich procesów różnych porcji danych tego samego rodzaju, np. wierszy macierzy – służy do tego funkcja `MPI_Scatter`

MPI_Scatter

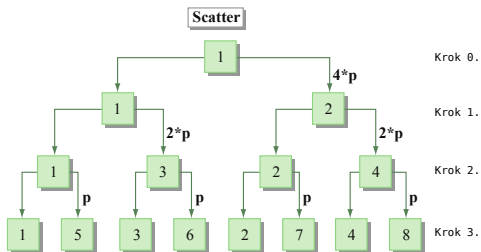
Sends data from one task to all other tasks in a group

```
sendcnt = 1;
recvcnt = 1;
src = 1;
MPI_Scatter(sendbuf, sendcnt, MPI_INT,
            recvbuf, recvcnt, MPI_INT,
            src, MPI_COMM_WORLD);
```

task 1 contains the message to be scattered



MPI_SCATTER – EFEKTYWNOŚĆ



Źródło: MIT OpenCourseWare

Podobnie jak MPI_Bcast operacja MPI_Scatter wykonywana jest efektywnie, tak by zaangażować w rozsyłanie danych jak najwięcej procesów.

MPI_GATHER

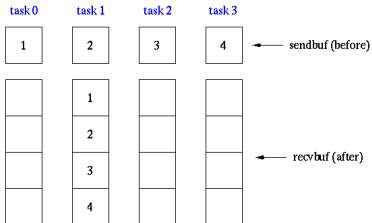
- ▶ Odwrotnością funkcji MPI_Scatter jest funkcja MPI_Gather, która polega na *zebraniu* przez jeden proces danych od pozostałych procesów

MPI_Gather

Gathers together values from a group of processes

```
sendcnt = 1;
recvnt = 1;
src = 1;
MPI_Gather(sendbuf, sendcnt, MPI_INT,
           recvbuf, recvnt, MPI_INT,
           src, MPI_COMM_WORLD);
```

messages will be gathered in task 1



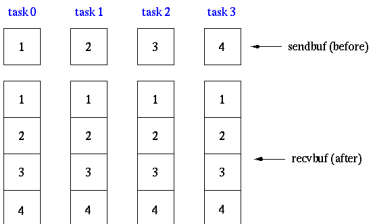
MPI_ALLGATHER

- ▶ Jeżeli *każdy* proces ma uzyskać dostęp do danych wszystkich procesów, to odpowiednią metodą jest MPI_Allgather

MPI_Allgather

Gathers together values from a group of processes and distributes to all

```
sendcnt = 1;
recvcnt = 1;
MPI_Allgather(sendbuf, sendcnt, MPI_INT,
              recvbuf, recvcnt, MPI_INT,
              MPI_COMM_WORLD);
```



OPERACJE REDUKCJI

- ▶ MPI udostępnia możliwość efektywnego wykonywania operacji **redukcji**, tj. łącznej operacji dwuargumentowej na zbiorze danych rozproszonych między procesami w komunikatorze:
`MPI_Reduce(void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)`
- ▶ Zdefiniowane są następujące rodzaje operacji (MPI_Op):
 - ▶ MPI_MAX, MPI_MIN, MPI_MAXLOC, MPI_MINLOC
 - ▶ MPI_SUM, MPI_PROD
 - ▶ MPI_LAND, MPI_BAND, MPI_LOR, MPI_BOR, MPI_LXOR, MPI_BXOR
- ▶ MPI_MINLOC i MPI_MAXLOC zwracają odpowiednio minimaklną i maksymalną wartość odraz indeks (rangę) procesu, który je posiada

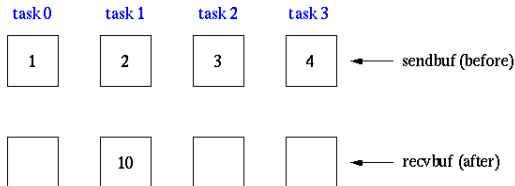
MPI_REDUCE

Prosty przykład sumowania liczb:

MPI_Reduce

Perform and associate reduction operation across all tasks in the group and place the result in one task

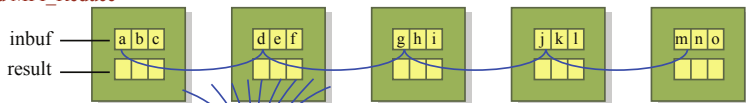
```
count = 1;  
dest = 1;           result will be placed in task 1  
MPI_Reduce(sendbuf, rcvbuf, count, MPI_INT, MPI_SUM,  
           dest, MPI_COMM_WORLD);
```



MPI_REDUCE

W ogólnym przypadku zadana operacja wykonywana jest na wszystkich elementach podanych tablic.

Przed MPI_Reduce



Po MPI_Reduce



a @ d @ g @ j @ m

gdzie @ oznacza wykonywaną operację

Źródło: MIT OpenCourseWare

Analogicznie do operacji rozgłaszania, wykonanie redukcji wymaga $O(\log n)$ etapów komunikacji

PRZYKŁAD – ZLICZANIE LICZB PIERWSZYCH

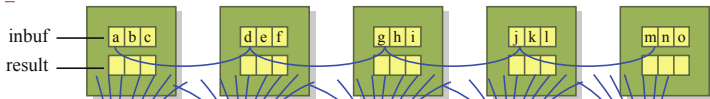
```
int n, proc_id, proc_num;
MPI_Comm_size(MPI_COMM_WORLD, &proc_num);
MPI_Comm_rank(MPI_COMM_WORLD, &proc_id);
if (proc_id == 0)
    printf("Podaj n: "); scanf("%d", &n);
MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);
double start = MPI_Wtime();
int sum = proc_id == 0 ? 1 : 0;
for (int i = 2*(proc_id + 1)+1; i <= n; i += 2*proc_num) {
    if (is_prime(i))
        sum++;
}
double koniec = MPI_Wtime();
int wynik = 0;
MPI_Reduce(&sum, &wynik, 1, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
if (proc_id == 0) {
    printf("Jest %d liczb pierwszych <= %d", wynik, n);
    printf("Czas obliczen: %.2lf sek.", koniec - start);
}
```

MPI_ALLREDUCE

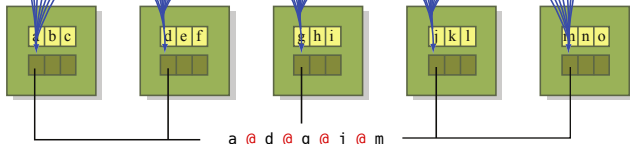
Możliwe jest również wykonywanie operacji redukcji dla wszystkich procesów jednocześnie za pomocą:

```
MPI_Allreduce(void *sendbuf, void *recvbuff, int
count, MPI_Datatype datatype, MPI_Op op, MPI_Comm
comm)
```

Przed MPI_Allreduce



Po MPI_Allreduce



gdzie @ oznacza wykonywaną operację

MPI_REDUCE_SCATTER

Jak nazwa wskazuje polega na wykonaniu:

- ▶ operacji redukcji na tablicach (wektorach) wszystkich procesów w komunikatorze (reduce)
- ▶ rezultat dzielony jest na rozłączne fragmenty rozsyłane do poszczególnych procesów (scatter)

MPI_Reduce_scatter

Perform reduction operation on vector elements across all tasks in the group, then distribute segments of result vector to tasks

```
recvcount = 1;
MPI_Reduce_scatter(sendbuf, recvbuf, recvcount, MPI_INT, MPI_SUM,
MPI_COMM_WORLD);
```

task 0	task 1	task 2	task 3	
0	0	0	0	← sendbuf (before)
1	1	1	1	
2	2	2	2	
3	3	3	3	
0	4	8	12	← recvbuf (after)

MPI_ALLTOALL

- ▶ Ciekawą operacją jest `MPI_Alltoall`, która polega na wysłaniu przez każdy proces unikalnej porcji danych do wszystkich pozostałych procesów
- ▶ $n(n - 1)$ operacji wysyłania komunikatów

MPI_Alltoall

Sends data from all to all processes. Each process performs a scatter operation.

```
sendcnt = 1;
recvcnt = 1;
```

```
MPI_Alltoall(sendbuf, sendcnt, MPI_INT,
             recvbuf, recvcnt, MPI_INT,
             MPI_COMM_WORLD);
```

task 0	task 1	task 2	task 3	task 0	task 1	task 2	task 3
1	5	9	13	1	2	3	4
2	6	10	14	5	6	7	8
3	7	11	15	9	10	11	12
4	8	12	16	13	14	15	16

sendbuf recvbuf after

OPERACJE PREFIKSOWE

- ▶ Operacje prefiksowe realizowane są przez funkcję

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

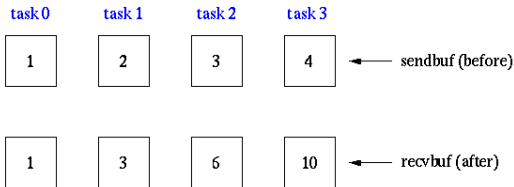
- ▶ Polegają na wykonaniu **częściowej redukcji**
- ▶ Proces i otrzymuje wynik redukcji dla danych procesów o indeksach (rangach) od 0 do i
- ▶ Proces o najwyższej randze ma wynik dla wszystkich danych
- ▶ Funkcja obsługuje takie operatory, jak MPI_Reduce

MPI_SCAN – PRZYKŁAD

MPI_Scan

Computes the scan (partial reductions) of data on a collection of processes

```
count = 1;  
MPI_Scan(sendbuf, recvbuf, count, MPI_INT, MPI_SUM,  
        MPI_COMM_WORLD);
```



TRWAŁA KOMUNIKACJA

- ▶ W przypadku częstej komunikacji, np. w pętli, można uniknąć narzutu czasowego związanego z inicjowaniem operacji komunikacji za pomocą utworzenia **trwałego** żądania komunikacyjnego (ang. persistent communication request) – „półkanał” komunikacyjny
- ▶ Oszczędność wynika z przypisania parametrów operacji na stałe do żądania
- ▶ `MPI_Send_init(buf, count, datatype, dest, tag, comm, request)` tworzy trwałe żądanie wysłania
- ▶ `MPI_Recv_init(buf, count, datatype, source, tag, comm, request)` tworzy trwałe żądanie odbierania
- ▶ Parameter request pełni podobną rolę jak w operacjach `MPI_Irecv` oraz `MPI_Isend`

TRWAŁA KOMUNIKACJA – SCENARIUSZ ZASTOSOWANIA

1. Utworzenie trwałego żądania za pomocą `MPI_Send_init` lub `MPI_Recv_init`
2. Zainicjowanie operacji wysłania / odbioru za pomocą `MPI_Start`
3. Czekanie na zakończenie zainicjowanej operacji za pomocą `MPI_Wait` lub `MPI_Waitall`
4. Usunięcie żądania za pomocą `MPI_Request_free`

Kroki 2. i 3. powtarzane są tyle razy, ile komunikatów ma zostać wysłanych / odebranych

TRWAŁA KOMUNIKACJA – PRZYKŁAD

```
MPI_Request req[2];
MPI_Recv_init(buf1, count, type, src, tag, comm, &req[0]);
MPI_Send_init(buf2, count, type, src, tag, comm, &req[1]);
const int N = 100000;
for (i=1; i < N; i++)
{
    MPI_Start(&req[0]);
    MPI_Start(&req[1]);
    MPI_Waitall(2, req, status);
    obliczenia(buf1, buf2);
}
MPI_Request_free(&req[0]);
MPI_Request_free(&req[1]);
```

DEFINIOWANIE NOWYCH TYPÓW DANYCH

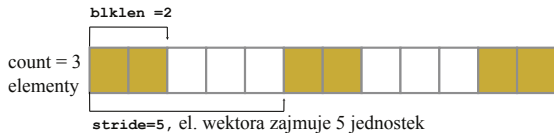
- ▶ Typy danych zdefiniowane w MPI pozwalają na łatwą wymianę tablic złożonych z elementów jednego z podstawowych typów, np. MPI_INT
- ▶ W programach często zachodzi konieczność wymiany złożonych danych różnego typu
 1. Można wysyłać kolejno składowe w osobnych operacjach – niewygodne i wolne
 2. Można zdefiniować własny typ dla przesyłanych danych – znacznie wygodniejsze
 3. Można stosować „pakowanie” i „rozpakowywanie” do i z bufora

DEFINIOWANIE NOWYCH TYPÓW DANYCH C.D.

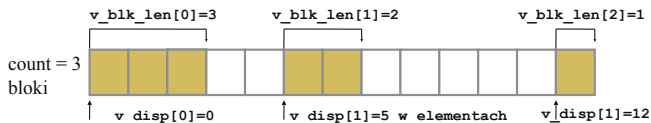
MPI udostępnia funkcje umożliwiające zdefiniowanie nowych typów danych należących do jednej z kategorii:

- ▶ Elementarne – typy danych zdefiniowane w języku programowania
- ▶ Ciągłe – tablice
- ▶ Wektory – tablice, w których kolejne elementy oddzielone są separatorem o określonym rozmiarze
- ▶ Indeksowane – podobnie jak wektor, ale adresy kolejnych elementów podawane są jawnie
- ▶ Strukturalne – najbardziej ogólne, umożliwia przesyłanie struktur (`struct`) z języka C

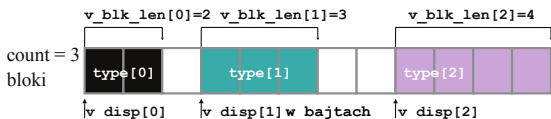
DEFINIOWANIE NOWYCH TYPÓW DANYCH C.D.



Wektor
(z separatorami)



T. indeksowany



Strukturalny

DEFINIOWANIE NOWYCH TYPÓW – PRZYKŁAD

```
typedef struct {
    float wartosc;
    char id;
} Pomiar;

int v_blk_len[2] = { 1, 1 }; // dł. bloków
MPI_Datatype v_types[2] = { MPI_FLOAT, MPI_CHAR };
Pomiar pom;
// ustal adresy kolejnych pól struktury
MPI_Aint v_addr[3];
MPI_Address(&pom, &v_addr[0]);
MPI_Address(&(pom.wartosc), &v_addr[1]);
MPI_Address(&(pom.id), &v_addr[2]);
// na ich podstawie wyznacz przesunięcia pól
MPI_Aint v_disp[2] = { v_addr[1] - v_addr[0],
                      v_addr[2] - v_addr[1] };

// zdefiniuj typ
MPI_Datatype pomiar_type;
MPI_Type_struct(2, v_blk_len, v_disp, v_types, &pomiar_type);
MPI_Type_commit(&pomiar_type);
```

DEFINIOWANIE NOWYCH TYPÓW – PRZYKŁAD, C.D.

```
const int N = 10;
Pomiar pomiary[N];
if (proc_id == 0) {
    for (int i = 0; i < N; ++i) {
        pomiary[i].wartosc = 1 + i;
        pomiary[i].id = 'A' + i;
    }
}
MPI_Bcast(pomiary, N, pomiar_type, 0, MPI_COMM_WORLD);
if (proc_id != 0) {
    for (int i = 0; i < N; i++) {
        cout << pomiary[i].wartosc
              << ", " << pomiary[i].id << endl;
    }
}
```

GRUPOWANIE DANYCH

- ▶ Definiowanie własnych typów danych wymaga dodatkowego nakładu pracy, który warto ponieść w przypadku częstej wymiany komunikatów o tej samej strukturze
- ▶ W przypadku, gdy wymiana danych złożonych pojawia się w jednym miejscu programu można skorzystać z możliwości przesyłania danych zgrupowanych (spakowanych)
- ▶ Wszystkie dane muszą zostać zapisane w ciągłym obszarze pamięci

GRUPOWANIE DANYCH – FUNKCJE

- ▶ `int MPI_Pack(void *pack_data, int in_count, MPI_Datatype datatype, void *buffer, int size, int *position_ptr, MPI_Comm comm)`
 - ▶ Bufor `pack_data` powinien zawierać `in_count` elementów typu `datatype`.
 - ▶ Parametr `position_ptr` zawiera adres początkowy w buforze wyjściowym `buffer` i po wykonaniu funkcji jego wartość jest uaktualniana (dodawana jest długość zapisanych danych)
 - ▶ Parametr `size` określa rozmiar obszaru pamięci wskazywanego przez `buffer`.
- ▶ `int MPI_Unpack(void buffer, int size, int position_ptr, void *unpack_data, int count, MPI_Datatype datatype, MPI_Comm comm)`
 - ▶ Funkcja rozpakowuje `count` elementów typu `datatype` ze zmiennej `buffer` począwszy od pozycji `position` do obszaru pamięci wskazywanego przez `unpack_data` i odpowiednio uaktualnia zmienną `position`.

GRUPOWANIE DANYCH – PRZYKŁAD

```
float a, b;
int l;
char buffer[100];
int pos;
// Proces wysyłający dane
pos=0;
MPI_Pack(&a, 1, MPI_FLOAT, buffer, 100, &pos, MPI_COMM_WORLD);
MPI_Pack(&b, 1, MPI_FLOAT, buffer, 100, &pos, MPI_COMM_WORLD);
MPI_Pack(&l, 1, MPI_INT, buffer, 100, &pos, MPI_COMM_WORLD);
MPI_Bcast(buffer, 100, MPI_PACKED, root, MPI_COMM_WORLD);
...
// Proces odbierający dane

MPI_Bcast(buffer, 100, MPI_PACKED, root, MPI_COMM_WORLD);
pos=0;
MPI_Unpack(buffer, 100, &pos, &a, 1, MPI_FLOAT, MPI_COMM_WORLD);
MPI_Unpack(buffer, 100, &pos, &b, 1, MPI_FLOAT, MPI_COMM_WORLD);
MPI_Unpack(buffer, 100, &pos, &l, 1, MPI_INT, MPI_COMM_WORLD);
```

RÓWNOLEGŁE I/O

- ▶ MPI udostępnia również szereg funkcji do wykonywania operacji wejścia/wyjścia w sposób równoległy (i rozproszony)
 - ▶ Program może jednocześnie zapisywać i odczytywać dany plik za pośrednictwem funkcji MPI-IO
 - ▶ Wymaga **równoległego systemu plików** (ang. parallel file system) – modułowy, rozproszony; Przykłady: Lustre, IBM GPFS, PanFS
- ▶ Alternatywy dla MPI-IO:
 - ▶ Wyznaczony proces zajmuje się odczytem / zapisem do pliku i rozsyłaniem oraz gromadzeniem danych
 - ▶ Każdy proces korzysta z odrębnego pliku

RÓWNOLEGŁE I/O C.D.

- ▶ Funkcje do obsługi I/O zostały zaimplementowane dopiero w MPI 2.0
- ▶ Zalety równoległego I/O, to:
 - ▶ zwiększona wydajność
 - ▶ zunifikowany dostęp do plików, podobnie jak w przypadku programu sekwencyjnego – ułatwia przetwarzanie plików z wynikami obliczeń
- ▶ Stosowanie sekwencyjnego I/O wymaga dodatkowo kopiowania plików między węzłami obliczeniowymi – każdy proces zapisuje do lokalnego systemu plików
- ▶ MPI-IO udostępnia mechanizmy do:
 - ▶ synchronizacji plików między serwerami oraz ich kopiowania
 - ▶ „nieciągłego” sposobu wykonywania operacji na pliku za pomocą widoków (ang. file views)

RÓWNOLEGŁE I/O – PRZYKŁAD

```
MPI_File fh;
MPI_Status status;
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
MPI_Comm_size(MPI_COMM_WORLD, &nprocs);
bufsize = FILESIZE/nprocs;
nints = bufsize/sizeof(int);
MPI_File_open(MPI_COMM_WORLD, "/home/datafile",
              MPI_MODE_RDONLY, MPI_INFO_NULL, &fh);
MPI_File_seek( fh, rank*bufsize, MPI_SEEK_SET);
MPI_File_read( fh, buf, nints, MPI_INT, &status);
MPI_File_close(&fh);
```