

Laborator nr. 3

Introducere în MPI (2)

1 Tipuri de date în MPI

Trebuie observat faptul că, pentru fiecare comunicație în parte, lungimea mesajului este specificată în număr de elemente, nu în număr de octeți. Tipurile de date de bază corespund tipurilor de date de bază ale limbajului folosit. Astfel avem:

Tabelul 1: Tipuri de date primare în MPI

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.FLOAT	float
MPI.DOUBLE	double
MPI.LONG.DOUBLE	long double
MPI.BYTE	
MPI.PACKED	

2 Comunicații colective

Conceptul cheie al **comunicațiilor colective** este acela de a avea un **grup de procese participante**. Pentru aceasta biblioteca MPI oferă următoarele funcții:

bariera – sincronizarea prin barieră pentru toții membrii grupului (MPI.Barrier);

broadcast – transmiterea de la un membru al grupului către ceilalți membri (MPI.Bcast);

gather – adunarea datelor de la toți membrii grupului la un singur membru din grup (MPI.Gather);

scatter – împrăștierea datelor de la un membru la toți membrii grupului (MPI.Scatter);

all gather – o variație a MPI.Gather în care toți membrii grupului primesc rezultatele (MPI.Allgather);

all-to-all scatter/gather – adunarea/împrăștierea datelor de la toți membrii grupului către toți membrii (MPI.Alltoall – o extensie a MPI.Allgather în care fiecare proces trimite date distincte către fiecare proces din grup);

operații de reducere – adunare, înmulțire, min, max, sau funcții definite de utilizator în care rezultatele sunt trimise la toți membrii grupului sau numai la un membru.

Observație: Toate funcțiile sunt executate colectiv (sunt apelate de toate procesele dintr-un comunicator); toate procesele folosesc aceiași parametri în apeluri.

Funcțiile nu au un identificator de grup ca argument explicit. În schimb avem un comunicator ca argument. Începând cu standardul MPI 2.0, comunicațiile colective acceptă ca parametru un inter-comunicator.

Caracteristicile comunicațiilor colective:

- Comunicațiile colective și cele de tip individuale **nu se vor influența reciproc**.
- Toate procesele trebuie să execute comunicațiile colective.
- Sincronizarea nu este garantată, excepție făcând doar bariera.
- Nu există comunicații colective neblocante.
- Nu există tag-uri.
- Bufferele de recepție trebuie să fie exact de dimensiunea potrivită.

2.1 Bariera

Bariera este un mecanism de sincronizare a unui grup de procese. Dacă un proces a realizat un apel către `MPI_Barrier(comm)`, unde `comm` este comunicatorul, atunci el va aștepta până când toate procesele au efectuat acest apel.

Apelul funcției **Barrier**:

`MPI_BARRIER(comm)`

Prototipul funcției este:

```
1 int MPI_Barrier( MPIComm comm )
```

Listing 1: Prototipul funcției MPI.Barrier (C)

2.2 Broadcast

Funcția `MPI_Bcast` trimite un mesaj de la procesul cu rangul *root* către toate procesele din grup, inclusiv procesului *root*. Schematic, acest comportament este prezentat în Figura 1.

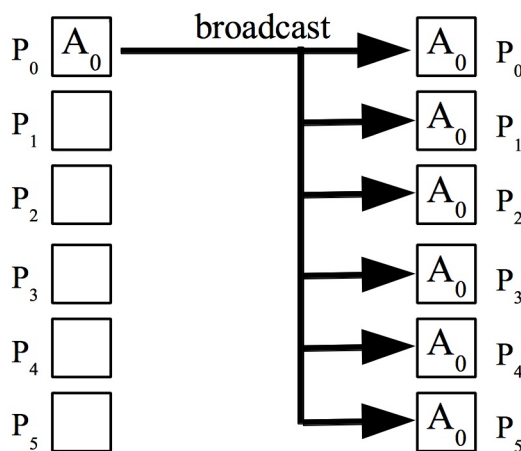


Figura 1: Operația colectivă *Broadcast*

Funcția este apelată de toți membrii grupului care folosesc același argument pentru comunicator și, respectiv, același rang pentru procesul **root**. La ieșirea din funcție, conținutul buffer-ului de comunicație al procesului **root** a fost transferat către toate procesele.

Apelul funcției **Bcast** este:

```
MPI_Bcast(buffer, count, datatype, root, comm)
```

unde:

INOUT buffer – adresa de început a buffer-ului pentru date

IN count – numărul de intrări din buffer

IN datatype – tipul datelor din buffer

IN root – rangul procesului **root** care inițiază comunicația

IN comm – comunicatorul pe care se va realiza comunicația efectivă.

Prototipul funcției este detaliat în Listing 2. În continuare, în Listing 3 este evidențiat un exemplu de utilizare.

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

Listing 2: Prototipul funcției MPI_Bcast (C)

```
1 // .....  
2 double param;  
3 MPI_Init(&argc, &argv);  
4 MPI_Comm_rank(MPI_COMM_WORLD, &rank);  
5 if(rank == 5) {  
6     param = 23.0;  
7 }  
8 MPI_Bcast(&param, 1, MPI_DOUBLE, 5, MPI_COMM_WORLD);  
9 printf("P:%d_dupa_broadcast_parametrul_este_%f\n", rank, param);  
10 MPI_Finalize();  
11 // .....
```

Listing 3: Utilizare MPI_Bcast (C)

2.3 Scatter

Procesul **root** împarte și trimite conținutul *sendbuffer*-ului către toate procesele din grup (inclusiv către procesul **root**). Schematic, acest lucru este reprezentat în Figura 2 și în Figura 3.

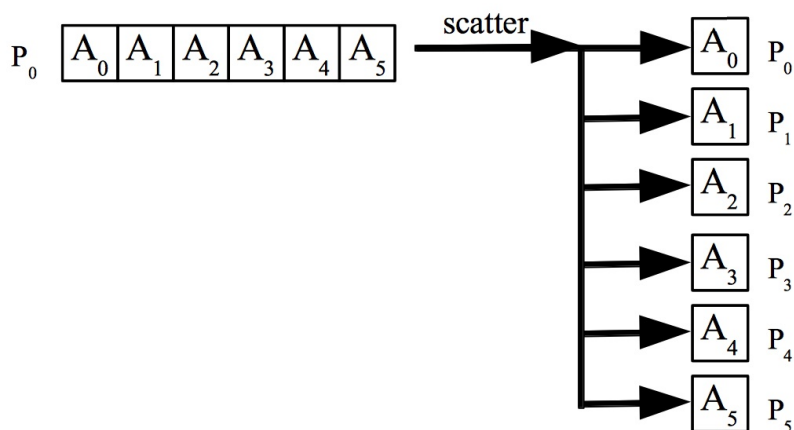


Figura 2: Operația colectivă *Scatter*

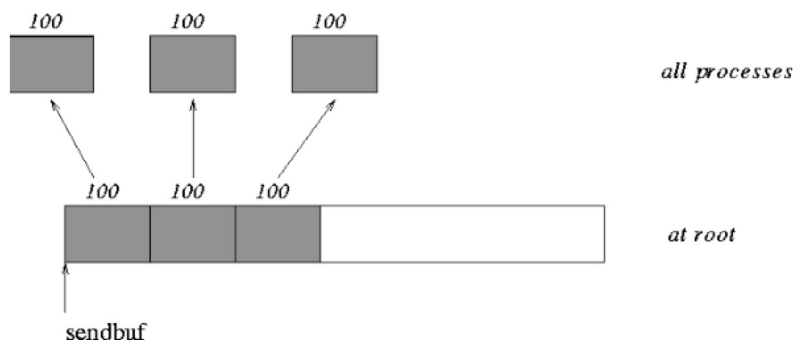


Figura 3: Comportamentul operației colective *Scatter*

Același rezultat s-ar obține dacă procesul **root** ar realiza n operații de tip *send* (unde n reprezintă numărul total de procese ce aparțin de grupul curent de lucru)

```
MPI_Send(sendbuf+i*sendcount*extent(sendtype), sendcount, sendtype, i,...)
```

și fiecare proces execută o operație *receive*

```
MPI_Recv(recvbuf, recvcount, recvtype, i,...).
```

Apelul funcției **Scatter**:

```
MPI_SCATTER(sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)
```

unde:

IN sendbuf – adresa de început a buffer-ului cu datele pentru trimitere

IN sendcount – numărul de intrări din buffer ce va fi trimis la fiecare proces

IN sendtype – tipul datelor din buffer

OUT recvbuf – adresa de început a bufferului de recepție

IN recvcount – numărul de intrări ce vor fi primite de fiecare proces

IN recvtype – tipul datelor din bufferul de recepție

IN root – rangul procesului ce trimite datele

IN comm – comunicatorul pe care se va realiza comunicația efectivă.

Prototipul funcției este expus în Listing 4.

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

Listing 4: Prototipul funcției MPI_Scatter (C)

Observații:

1. Toate argumentele funcției sunt semnificative numai în procesul *root*, în timp ce pentru celelalte procese numai argumentele *recvbuf*, *recvcount*, *recvtype*, *root* și *comm* sunt semnificative.
2. *sendbuffer*, *sendcount* și *sendtype* sunt ignorate în procesele care nu au rangul *root*.

În Listing 5 este prezentat un exemplu de utilizare a funcției MPI_Scatter.

```

1 // .....
2 MPI_Comm comm;
3 int gsize,*sendbuf;
4 int root, rbuf[100];
5 // ...
6 MPI_Comm_size(comm, &gsize);
7 sendbuf = (int *)malloc(gsize*100*sizeof(int));
8 // ...
9 MPI_Scatter(sendbuf, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm);
10 // .....

```

Listing 5: Utilizare MPI_Scatter (C)

2.4 Gather

Fiecare proces (inclusiv procesul **root**) trimite conținutul `sendbuf` er-ului către procesul **root**. Procesul **root** primește mesajele și le **stochează în ordinea rangului**. Conceptual, acest lucru este exemplificat în Figurile 4 și 5.

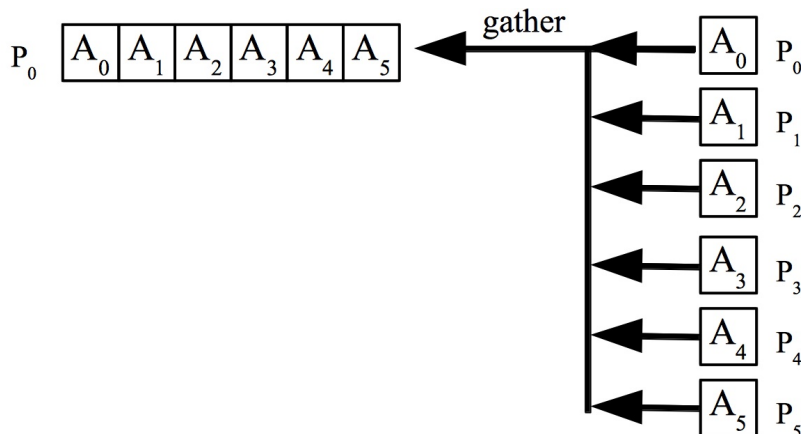


Figura 4: Operația colectivă *Gather*

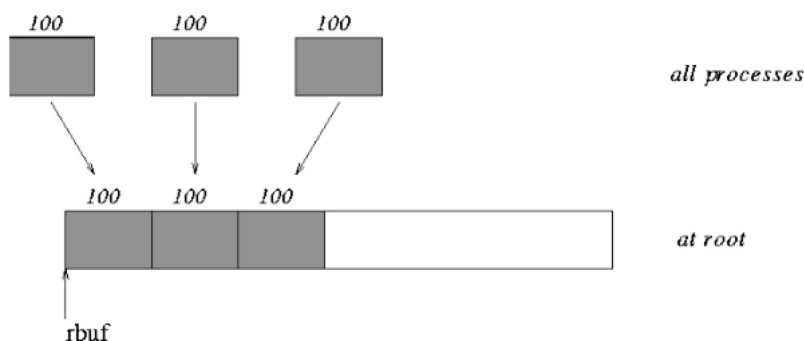


Figura 5: Comportamentul operației colective *Gather*

Același rezultat s-ar obține dacă toate procesele (inclusiv procesul **root**) ar realiza câte o operație de tip *send*

`MPI_Send(sendbuf, sendcount, sendtype, root,...)`

și procesul **root** execută n operații de tip *receive* (unde n reprezintă numărul total de procese ce aparțin de grupul curent de lucru)

```
MPI_Recv(recvbuf+i*recvcount*extent(recvtype), recvcount, recvtype, i ,...).
```

Apelul funcției **Gather**:

```
MPI_GATHER( sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm)
```

unde:

IN sendbuf – adresa de început a zonei buffer-ului de date

IN sendcount – numărul de elemente din buffer

IN sendtype – tipul datelor din buffer

OUT recvbuf – adresa buffer-ului de recepție (are semnificație numai pentru procesul root)

IN recvcount – numărul de elemente pentru fiecare recepție în parte (are semnificație numai pentru procesul root)

IN recvtype – tipul datelor din buffer-ul de recepție (are semnificație numai pentru procesul root)

IN root – rangul proceselor care recepționează datele

IN comm – comunicatorul pe care se va realiza comunicația efectivă.

Prototipul funcției este următorul:

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

Listing 6: Prototipul funcției MPI_Gather (C)

Observații:

1. *Buffer-ul de recepție este ignorat pentru procesele care nu au rangul root.*
2. *Argumentul recvcount la procesul root indică dimensiunea datelor primite de la fiecare proces și nu numărul total de date primite.*

Listing-urile 7 și 8 exemplifică utilizarea corectă a funcției Gather.

```
1 // .....
2 MPI_Comm comm;
3 int gsize, sendarray[100];
4 int root, *rbuf;
5 // ...
6 MPI_Comm_size( comm, &gsiz );
7 rbuf = (int *) malloc( gsize*100*sizeof(int) );
8 MPI_Gather( sendarray, 100, MPI_INT, rbuf, 100, MPI_INT, root, comm );
9 // .....
```

Listing 7: Utilizare MPI_Gather (C) (1)

```
1 // .....
2 MPI_Comm comm;
3 int gsize, sendarray[100];
4 int root, myrank, *rbuf;
5 // ...
6 MPI_Comm_rank( comm, &myrank );
7 if ( myrank == root ) {
8     MPI_Comm_size( comm, &gsiz );
```

```
9      rbuf = (int *) malloc ( gsize*100*sizeof (int ) );  
10 }  
11 MPI_Gather( sendarray , 100, MPI_INT, rbuf , 100, MPI_INT, root , comm );  
12 // .....
```

Listing 8: Utilizare MPI_Gather (C) (2)

3 Topologii de comunicare în MPI

O topologie oferă un mecanism convenabil pentru denumirea proceselor dintr-un grup, și în plus, poate oferi soluții sistemului, în timpul rulării, pentru maparea proceselor pe hardware-ul existent. Un grup de procese în MPI reprezintă o colecție de n procese. Fiecare proces din grup are atribuit un rang (valoare întreagă) între 0 și $n-1$. În multe aplicații paralele o atribuire liniară a acestor ranguri nu reflectă logica modelului de comunicații necesar. Adesea rangurile proceselor sunt aranjate în topologii, precum o plasă cu 2 sau 3 dimensiuni. Mai general, logica aranjării proceselor este descrisă de un graf. În continuare, această aranjare logică a proceselor o vom numi: „topologie virtuală”.

Trebuie făcută distincție între topologia virtuală a proceselor și topologia fizică (hardware). Descrierea topologiilor virtuale depinde numai de aplicație și este independentă de mașină. Structura comunicațiilor pentru un set de procese poate reprezenta un graf, în care nodurile reprezintă procesele iar muchiile conectează procesele care comunică între ele.

MPI asigură transmiterea mesajelor între oricare dintre perechile de procese din grup. Nu este nevoie ca un canal de comunicație să fie deschis explicit. **De aceea, o muchie lipsă în graful de procese definit de utilizator nu asigură faptul că procesele care nu sunt conectate nu fac schimb mesaje.**

Informațiile legate de topologie sunt asociate comunicatorilor. Astfel, următoarele informații legate de structura topologiei pot fi obținute prin interogarea comunicatorului:

1. Tipul topologiei (carteziană, graf sau nedefinită)
2. Pentru o topologie carteziană:
 - (a) ndims (numărul de dimensiuni),
 - (b) dims (numărul de procese pe direcția de coordonate),
 - (c) periods (periodicitatea informației),
 - (d) own_position (poziția din topologia creată)
3. Pentru o topologie de tip graf:
 - (a) index,
 - (b) muchii (reprezintă vectori ce definesc structura grafului).

Ca și pentru alte apeluri de funcții colective, **programul trebuie scris astfel încât să funcționeze corect indiferent dacă apelurile se sincronizează sau nu.**

3.1 Topologii carteziane

Funcția MPI_CART_CREATE poate fi folosită pentru a descrie structuri carteziane de dimensiune arbitrară. Pentru fiecare coordonată se poate specifica dacă structura este periodică sau nu. De exemplu, o topologie 1D, este liniară dacă nu este periodică și este un inel dacă este periodică. Pentru o topologie 2D, avem un dreptunghi, un cilindru sau un tor după cum o dimensiune este sau nu periodică.

Observație: Un hipercub cu n dimensiuni este un tor cu n dimensiuni cu câte 2 procese pe fiecare coordonată. De aceea nu este necesar un suport special pentru o structură de hipercub.

Apelul funcției **MPI_Cart_create**:

```
MPI_CART_CREATE(comm_old, ndims, dims, periods, reorder, comm_cart)
```

unde:

IN comm_old – comunicatorul peste care va fi creată topologia nouă

IN ndims – numărul de dimensiuni ale topologiei carteziane

IN dims – tablou de întregi de dimensiune ndims care specifică numărul de procese de pe fiecare dimensiune

IN periods – tablou de dimensiune ndims care specifică dacă „grila” este periodică (true) sau nu (false) pe fiecare dimensiune

IN reorder – rangurile pot fi reordonate (true) sau nu (false)

OUT comm_cart – comunicator cu noua topologie carteziană.

Prototipul funcției este următorul:

```
1 int MPI_Cart_create(MPIComm comm_old, int ndims, int *dims,
2     int *periods, int reorder, MPIComm *comm_cart)
```

Listing 9: Prototipul funcției MPI_Cart_create (C)

Pentru topologiile carteziane, funcția MPI_DIMS_CREATE ajută utilizatorul să realizeze o distribuire echilibrată a proceselor pe fiecare direcție de coordonate. O posibilitate ar fi împărțirea tuturor proceselor într-o topologie n-dimensională.

Apelul funcției **MPI_Dims_create**:

```
MPI_DIMS_CREATE(nnodes, ndims, dims)
```

unde:

IN nnodes – numărul de noduri din grilă

IN ndims – numărul de dimensiuni ale topologiei

INOUT dims – tablou de întregi care specifică numărul de noduri pe fiecare dimensiune

Prototipul funcției este următorul:

```
1 int MPI_Dims_create(int nnodes, int ndims, int *dims)
```

Listing 10: Prototipul funcției MPI_Dims_create (C)

Pentru fiecare `dims[i]` setat de apelul funcției MPI_DIMS_CREATE, `dims[i]` va fi ordonat descrescător. Tabloul `dims` este potrivit pentru a fi utilizat ca intrare în apelul funcției MPI_CART_CREATE. Un exemplu de rezultat oferit de funcția MPI_Dims_create poate fi observat în Tabelul 2.

Tabelul 2: MPI_Dims_create – scenarii de utilizare

dims (înainte de apel)	parametrii de apel	dims (după apel)
(0, 0)	MPI_DIMS_CREATE(6, 2, dims)	(3, 2)
(0, 0)	MPI_DIMS_CREATE(7, 2, dims)	(7, 1)
(0, 3, 0)	MPI_DIMS_CREATE(6, 3, dims)	(2, 3, 1)
(0, 3, 0)	MPI_DIMS_CREATE(7, 3, dims)	apel eronat

Odată topologia stabilită, aceasta poate fi interogată pentru a obține diferite informații.

Aflarea dimensiunilor un topologii carteziane:

```
MPI_CARTDIM_GET(comm, ndims)
```

unde:

IN comm – comunicator cu structura carteziană

OUT ndims – numărul de dimensiuni al structurii carteziane

Prototipul funcției este următorul:

```
1 int MPI_Cartdim_get(MPIComm comm, int *ndims)
```

Listing 11: Prototipul funcției MPI_Cartdim_get (C)

Aflarea rank-ului unui proces mapat la un anumit set de coordonate:

```
MPI_CART_RANK (comm, coords, rank)
```

unde:

IN comm – comunicator cu structura carteziană

IN coords – tablou de întregi care specifică coordonatele carteziane ale unui proces

OUT rank – ragul procesului specificat

Prototipul funcției este următorul:

```
1 int MPI_Cart_rank(MPIComm comm, int *coords, int *rank)
```

Listing 12: Prototipul funcției MPI_Cart_rank (C)

Pentru un grup de procese cu o structură carteziană, funcția MPI_CART_RANK transformă coordonatele logice ale procesului în rangul procesului.

Aflarea informațiilor despre topologia carteziană curentă:

```
MPI_CART_GET(comm, maxdims, dims, periods, coords)
```

unde:

IN comm – comunicator cu structura carteziană

IN maxdims – lungimea vectorilor dims, periods, și coords în program

OUT dims – numărul de procese pe fiecare dimensiune carteziană

OUT periods – periodicitatea (true/ false) pentru fiecare dimensiune

OUT coords – coordonatele procesului care apează funcția în structura carteziană

Prototipul funcției este următorul:

```
1 int MPI_Cart_get(MPIComm comm, int maxdims, int *dims, int *periods,
2 int *coords)
```

Listing 13: Prototipul funcției MPI_Cart_get (C)

Aflarea coordonatelor unui anumit proces:

```
MPI_CART_COORDS(comm, rank, maxdims, coords)
```

unde:

IN comm – comunicator cu structură carteziană

IN rank – rank-ul procesului pentru care se dorește aflarea coordonatelor

IN maxdims – lungimea vectorului dims

OUT coords – vector de ieșire ce conține coordonatele rank-ului indicat

Prototipul funcției este următorul:

```
1 int MPI_Cart_coords(MPI_Comm comm, int rank, int maxdims, int *coords)
```

Listing 14: Prototipul funcției MPI_Cart_coords (C)

Listingul 15 prezintă un model de utilizare a metodelor descrise anterior (crearea unei topologii carteziene, interogare coordonate rank, interogare rank la coordonate).

```
1 #include <stdio.h>
2 #include "mpi.h"
3
4 void main(int argc, char *argv[]) {
5     int rank;
6     MPI_Comm vu;
7     int dim[2], period[2], reorder;
8     int coord[2], id;
9
10    MPI_Init(&argc, &argv);
11    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
12    dim[0]=4; dim[1]=3;
13    period[0]=TRUE; period[1]=FALSE;
14    reorder=TRUE;
15    MPI_Cart_create(MPI_COMM_WORLD, 2, dim, period, reorder, &vu);
16
17    if(rank==5){
18        MPI_Cart_coords(vu, rank, 2, coord);
19        printf("P:%d_My_coordinates_are_%d_%d\n", rank, coord[0],
20    }
21    if(rank==0) {
22        coord[0]=3; coord[1]=1;
23        MPI_Cart_rank(vu, coord, &id);
24        printf("The_processor_at_position_(%d,%d)_has_rank%d\n",
25            coord[0], coord[1], id);
26    }
27
28    MPI_Finalize();
29    return 0;
30 }
```

Listing 15: Creare unui topologie carteziene și exemple de interogări (C)

3.2 Topologiile de tip graf

Crearea unei topologii de tip graf:

```
MPI_GRAPH_CREATE(comm_old, nnodes, index, edges, reorder, comm_graph)
```

unde:

IN comm_old – comunicator de intrare

IN nnodes – numărul de noduri din graf

IN index – tablou de întregi care descriu nodurile

IN edges – tablou de întregi care descriu muchiile grafului

IN reorder – rangurile pot fi ordonate sau nu

OUT comm_graph – comunicator cu topologia grafului

Prototipul funcției este următorul:

```
1 int MPI_Graph_create(MPIComm comm_old, int nnodes, int *index,
2 int *edges, int reorder, MPIComm *comm_graph)
```

Listing 16: Prototipul funcției MPI_Graph_create (C)

Pentru exemplificare, considerăm următorul graf:

Tabelul 3: Exemplu de graf

proces	vecini
0	1, 3
1	0
2	3
3	0, 2

În acest caz, argumentele `nnodes`, `index` și `edges` corespunzătoare vor fi:

```
nnodes=4
index={2,3,4,6}
edges={1,3,0,3,0,2}
```

Pentru o structură de tip graf, numărul de noduri este egal cu numărul de procese din grup. De aceea numărul nodurilor nu trebuie memorat separat.

Funcțiile `MPI_GRAPHDIMS_GET` și `MPI_GRAPH_GET` primesc informații privind topologia grafului care a fost asociat cu un comunicator de către funcția `MPI_GRAPH_CREATE`. Informațiile furnizate de `MPI_GRAPHDIMS_GET` pot fi utilizate pentru dimensionarea corectă a vectorilor `index` și `edges` pentru a realiza un apel `MPI_GRAPH_GET`.

```
MPI_GRAPHDIMS_GET(comm, nnodes, nedges)
```

unde:

IN comm – comunicator cu grupul care are structura de graf

OUT nnodes – numărul de noduri din graf (același cu numărul proceselor din grup)

OUT nedges – numărul de muchii din graf

Prototipul funcției este următorul:

```
1 int MPI_Graphdims_get(MPI_Comm comm, int *nnodes, int *nedges)
```

Listing 17: Prototipul funcției MPI_Graphdims_get (C)

```
MPI_GRAPH_GET(comm, nnodes, nedges)
```

unde:

IN comm – comunicator cu structura grafului

IN maxindex – lungimea vectorului index în programul apelant

IN maxedges – lungimera vectorului edges în programul apelant

OUT index – tablou de întregi ce conține structura grafului

OUT edges – tablou de întregi ce conține structura grafului

Prototipul funcției este următorul:

```
1 int MPI_Graph_get(MPI_Comm comm, int maxindex, int maxedges,  
2 int *index, int *edges)
```

Listing 18: Prototipul funcției MPI_Graph_get (C)

Funcțiile MPI_GRAPH_NEIGHBORS_COUNT și MPI_GRAPH_NEIGHBORS furnizează informații suplimentare pentru o topologie de tip graf.

```
MPI_GRAPH_NEIGHBORS_COUNT(comm, rank, nneighbors)
```

unde:

IN comm – comunicator cu topologia grafului

IN rank – rangul procesului din grupul comm

OUT nneighbors – numărul de vecini ai procesului specificat

Prototipul funcției este următorul:

```
1 int MPI_Graph_neighbors_count(MPI_Comm comm, int rank, int *nneighbors)
```

Listing 19: Prototipul funcției MPI_Graph_neighbors_count (C)

```
MPI_GRAPH_NEIGHBORS(comm, rank, maxneighbors, neighbors)
```

unde:

IN comm – comunicator cu topologia grafului

IN rank – rangul procesului din grupul comm

IN maxneighbors – dimensiunea tabloului ce va stoca rangurile vecinilor

OUT neighbors – rangurile proceselor care au ca vecini procesul specificat

Prototipul funcției este următorul:

```
1 int MPI_Graph_neighbors(MPI_Comm comm, int rank, int maxneighbors,  
2 int *neighbors)
```

Listing 20: Prototipul funcției MPI_Graph_neighbors (C)