

## Laborator nr. 2

### Introducere în MPI (1)

## 1 Standardul MPI

**Acronim:** *Message Passing Interface*

**MPI (Message Passing Interface)** reprezintă o bibliotecă pentru programarea paralelă. MPI a fost creat în cadrul grupului **MPI Forum** între 1993 și 1994, când a fost elaborat standardul 1.0. Acest standard definea, la momentul respectiv, mai mult de 120 de funcții care oferă suport pentru:

**Comunicații punct-la-punct** – operațiile de bază în recepția și trimiterea mesajelor de către procesele unui grup de lucru sunt `send` și `receive` cu variantele de implementare a acestora.

**Operații colective** – sunt operații definite ca și comunicații care implică un grup de procese.

**Grupuri de procese și contexte de comunicație**

**Topologia proceselor**

În MPI un grup de procese este o **colecție de  $n$  procese** în care fiecare proces are atribuit un **rang între 0 și  $n-1$** . În multe aplicații paralele o atribuire liniară a rangurilor nu reflectă logica necesară comunicațiilor între procese. Adesea, procesele sunt aranjate în modele topologice (de exemplu o plasă cu 2 sau 3 dimensiuni – aplicații în calcul matriceal). Într-un caz și mai general procesele pot fi descrise de un graf și vom face referire la aranjarea acestor procese într-o **topologie virtuală**.

O distincție clară trebuie făcută între o **topologie virtuală** și **topologia de bază** a suportului (hardware-ul propriu-zis). **Topologia virtuală** poate fi exploatată de sistem în atribuirea proceselor către procesoarele fizice, dacă acest lucru ajută la creșterea performanțelor de comunicație pentru o anumită mașină, dar realizarea acestui lucru nu face obiectul MPI. Descrierea unei topologii virtuale depinde numai de aplicație și este independentă de sistemul de calcul.

**Legătura cu Fortran 77 și C** – aplicațiile paralele folosind MPI pot fi scrise în limbaje ca Fortran 77 și C.

**Obținerea informațiilor despre mediu și gestiunea acestuia** – sunt introduse rutine pentru obținerea și setarea unor parametri pentru execuția programelor corespunzătoare diferitelor implementări ale MPI.

Standardul MPI a ajuns la versiunea 3.0 și oferă în plus suport pentru:

**Extinderea operațiilor colective** – în MPI-1 intercomunicarea reprezenta comunicare între două grupuri de procese disjuncte. În multe funcții argumentul de comunicare poate fi unul de **intercomunicare**, excepție făcând operațiile colective. Acest lucru este implementat în standardul MPI-2 și extinde utilizarea funcțiilor pentru operațiile colective.

**Crearea proceselor și gestiunea dinamică a acestora** – obiectivele principale sunt oferirea posibilității de a crea aplicații cu un număr variat de task-uri (*task farm*) și cuplarea aplicațiilor după modelul client/server.

**One-sided Communication** – ideea de bază este aceea ca un task să aibă acces direct la memoria altui task (asemănător modelului shared memory). În plus, se oferă și metode de sincronizare.

**Parallel File I/O** – mecanismele care oferă operații de I/O paralele, permit task-urilor să acceseze fișierele într-un mod controlat, dar sunt dependente de arhitectură. MPI-2 definește rutine de nivel șalt pentru a furniza un model de programare unic pentru operațiile de I/O paralele.

**Legătura cu alte limbaje** – este oferit suport pentru C++. De asemenea, există și o variantă pentru Java.

## 2 Modelul de execuție

Execuția unui program scris cu apeluri MPI constă în **rularea simultană a mai multor procese**, definite în mod **static** (MPI-1) sau **dinamic** (începând MPI-2), identice sau diferite și care **comunică punct-la-punct** sau **colectiv**. Procesele comunicante fac parte din același **comunicator**, care definește **domeniul de comunicație**. În cadrul domeniului de comunicație **fiecare proces are un rang**. De fapt comunicatorul asigură un mecanism de identificare a grupurilor de procese și de protecție a comunicației. Un comunicator poate fi de două feluri: **intracomunicator** pentru procese din interiorul aceluiași grup și **intercomunicator** pentru comunicațiile dintre grupuri.

Există șase apeluri de funcții mai importante care ar putea fi suficiente pentru scrierea unui program MPI:

**MPI.Init** – inițiază o aplicație MPI;

**MPI.Finalize** – termină o aplicație MPI;

**MPI.Comm\_size** – determină numărul proceselor;

**MPI.Comm\_rank** – determină identificatorul procesului;

**MPI.Send** – trimite un mesaj;

**MPI.Receive** – primește un mesaj.

Înainte de apelul MPI.Init și după apelul MPI.Finalize se pot apela doar funcțiile de interogare a mediului MPI. În Listing 1 este prezentată o primă aplicație MPI.

```
1 #include <stdio.h>
2 #include "mpi.h"
3
4 //!!! ATENTIE: PENTRU MAIN SE VA UTILIZA TOT TIMPUL
5 //          FORMA COMPLETA A FUNCTIEI !!!
6 int main ( int argc, char ** argv ) {
7     int count, myrank;
8
9     MPI_Init (&argc, &argv);
10    MPI_Comm_size (MPI_COMM_WORLD, &count);
11    MPI_Comm_rank (MPI_COMM_WORLD, &myrank);
12
13    printf ("Eu_sunt_%d_din_%d\n", myrank, count);
14
15    MPI_Finalize ();
16
17    //!!! OBLIGATORIU IN FINAL DE APLICATIE !!!
18    return 0;
19 }
```

Listing 1: Hello World MPI

### Compilarea unei aplicații MPI – GNU C/C++ compiler:

- sursă C

```
mpicc nume_sursa.c -o nume_executabil.bin
```

- sursă C++

```
mpic++ nume_sursa.cpp -o nume_executabil.bin
```

**Rularea aplicațiilor MPI – cazul generic simplist** (## reprezintă numărul dorit de procese – valoare întreagă strict pozitivă):

```
mpirun -np ## nume_executabil.bin
```

Parametrul `MPI_COMM_WORLD` precizează că toate procesele aparțin aceluiași grup și reprezintă comunicatorul inițial. Procesele pot executa instrucțiuni diferite, dar pentru aceasta trebuie să se asigure un mecanism de diferențiere (Listing 2).

```
1 #include <stdio.h>
2 #include "mpi.h"
3
4 !!!! ATENTIE: PENTRU MAIN SE VA UTILIZA TOT TIMPUL
5 //          FORMA COMPLETA A FUNCTIEI !!!
6 int main ( int argc , char ** argv ) {
7     int count , myrank ;
8
9     MPI_Init ( &argc , &argv );
10    MPI_Comm_size ( MPI_COMM_WORLD , &count );
11    MPI_Comm_rank ( MPI_COMM_WORLD , &myrank );
12
13    /*
14     .....
15    */
16
17    if ( myrank == 0 )
18        //cod program master
19    else
20        //cod program slave
21
22    /*
23     .....
24    */
25
26    MPI_Finalize ();
27
28    !!!! OBLIGATORIU IN FINAL DE APLICATIE !!!
29    return 0;
30 }
```

Listing 2: Exemplu de diferențiere a execuției proceselor MPI ce aparțin de același grup

## 3 Comunicații

O mare parte a funcțiilor MPI sunt funcții pentru comunicații. Comunicațiile pot fi **punct-la-punct** (sau **individuale**) sau **colective**. În continuare, **comunicațiile individuale** pot fi **blocante** sau **neblocante**. **Comunicațiile colective** sunt **întotdeauna blocante**, dar pot fi împărțite în comunicații colective **explicite** (*sunt caracterizate de un mesaj propriu*) sau **implicite** (*nu au un mesaj propriu*).

### 3.1 Comunicații blocante

După cum sugerează și tipul acestor comunicații, procesul implicat într-un astfel de transfer de date se blochează până la finalizarea acestui schimb. O funcție de comunicare cu blocare se consideră a fi finalizată odată cu trimiterea, respectiv primirea, mesajului. Astfel, la apelul funcției `MPI_Send` execuția poate continua după revenirea din apel, chiar dacă mesajul nu a ajuns încă la destinație. Asemănător, se revine din apelul funcției `MPI_Recv` imediat ce mesajul a ajuns în buffer-ul de destinație și poate fi citit.

Apelul unei **funcții send**:

```
MPI_SEND(buf, count, datatype, dest, tag, comm)
```

unde:

**IN buf** adresa zonei de memorie ce conține datele de trimis

**IN count** numărul de elemente de transmis

**IN datatype** tipul datelor pentru fiecare element din buffer

**IN dest** identificatorul procesului destinație

**IN tag** identificator de mesaj

**IN comm** comunicatorul

Prototipul funcției este prezentat în Listing 3.

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

Listing 3: Prototipul funcției `MPI_Send` (C)

Apelul unei **funcții receive**:

```
MPI_RECV (buf, count, datatype, source, tag, comm, status)
```

unde:

**OUT buf** adresa zonei de memorie în care va fi stocat mesajul recepționat în cazul în care este primit corespunzător parametrilor funcției

**IN count** numărul maxim de elemente care pot fi primite

**IN datatype** tipul datelor pentru fiecare element din buffer primit

**IN source** identificatorul procesului sursă

**IN tag** identificator de mesaj

**IN comm** comunicatorul

**OUT status** variabilă utilă pentru aflarea dimensiunii, identificatorului de mesaj și a sursei

Prototipul funcției este prezentat în Listing 4.

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

Listing 4: Prototipul funcției `MPI_Recv` (C)

Un exemplu de utilizare a celor două funcții este prezentat în Listing 5.

```
1 #include <stdio.h>
2 #include "mpi.h"
3
4 int main( int argc, char **argv ) {
5     char message[20];
6     int myrank;
7
8     MPI_Status status;
9
10    MPI_Init( &argc, &argv );
11    MPI_Comm_rank( MPI_COMM_WORLD, &myrank );
12
13    if (myrank == 0) { /* codul pentru procesul 0 */
14        strcpy(message, "Hello");
15        MPI_Send(message, strlen(message), MPI_CHAR,
16                1, 99, MPI_COMM_WORLD);
17    } else { /* codul pentru procesul 1 */
18        MPI_Recv(message, 20, MPI_CHAR,
19                0, 99, MPI_COMM_WORLD, &status);
20        printf("received :%s:\n", message);
21    }
22
23    MPI_Finalize();
24    return 0;
25 }
```

Listing 5: Prototipul funcției MPI\_Send (C)

### 3.2 Moduri de comunicație

Există mai multe moduri de comunicație, indicate în numele funcției printr-o literă astfel: B pentru comunicații bufferate, S pentru comunicație sincronă și R pentru comunicații de tip ready.

O operație de tip send în mod **buffered** poate fi inițiată indiferent dacă operația receive corespunzătoare a fost inițiată sau nu. În acest caz terminarea unei operații de tip send nu depinde de apariția unei operații receive. Pentru a termina operația este necesar ca datele ce trebuie recepționate să fie într-un buffer local. În acest scop, bufferul trebuie furnizat de către aplicație. Spațiul (de memorie) ocupat de buffer este eliberat atunci când mesajul este transferat către destinație sau când operația este anulată.

O operație de tip send în mod **sincron** poate fi inițiată indiferent de starea operației receive corespunzătoare (a fost sau nu inițiată). Oricum, operația va fi încheiată cu succes numai dacă o operație receive corespunzătoare a fost inițiată și a început recepția mesajului trimis prin send. Astfel, încheierea operației send sincrone nu indică numai că buffer-ul poate fi reutilizat, dar mai indică și faptul că receptorul a ajuns într-un anumit punct din execuția sa și anume faptul că a inițiat o operație receive corespunzătoare.

**Observație:** Modul sincron asigură faptul că o comunicare nu se termină la fiecare sfârșit de proces, înainte ca ambele procese să ajungă la aceeași operație comunicare.

În modul sincron operația de comunicare este încheiată simultan de către ambele funcții indiferent de primul apel.

O operație de tip send în mod **ready** poate fi inițiată numai dacă o operație receive a fost deja inițiată, altfel operația este o sursă de erori. Pe unele sisteme acest lucru poate permite înlăturarea operațiilor de tip handshake și au ca rezultat creșterea performanțelor. De aceea, într-un program scris corect o operație ready-send poate fi înlocuită de una standard fără a avea alt efect asupra comportării programului decât performanța.

### 3.3 Comunicații neblocaante

Un mecanism care adesea conferă performanțe mai bune este utilizarea comunicațiilor neblocaante. Comunicațiile neblocaante pot fi folosite în cele patru moduri ca și cele blocaante: standard, bufferate, sincrone și ready. Modul de folosire este același cu cel de la comunicațiile blocaante.

Cu excepția modului ready operația send neblocaantă poate fi inițiată chiar dacă o operație receive a fost sau nu inițiată, iar una de tip ready numai dacă a fost inițiată o operație receive. În toate cazurile operația send este inițiată local și returnează imediat indiferent de starea altor procese.

Apelul unei **funcții Isend**:

```
MPI_ISEND(buf, count, datatype, dest, tag, comm, request)
```

unde:

**IN buf** adresa zonei de memorie ce conține datele de trimis

**IN count** numărul de elemente de trimis

**IN datatype** tipul datelor pentru fiecare element din buffer

**IN dest** rangul destinației

**IN tag** identificator de mesaj

**IN comm** comunicatorul

**OUT request** handler-ul comunicației

Prototipul funcției este prezentat în Listing 6.

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

Listing 6: Prototipul funcției MPI\_Isend (C)

Apelul unei **funcții Irecv**:

```
MPI_IRecv (buf, count, datatype, source, tag, comm, request)
```

unde:

**OUT buf** adresa zonei de memorie în care va fi stocat mesajul recepționat în cazul în care este primit corespunzător parametrilor funcției

**IN count** numărul de elemente de trimis din buffer

**IN datatype** tipul datelor pentru fiecare element din buffer

**IN source** identificatorul procesului sursă

**IN tag** identificator de mesaj

**IN comm** comunicatorul

**OUT request** handler-ul comunicației

Prototipul funcției este prezentat în Listing 7.

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

Listing 7: Prototipul funcției MPI\_Irecv (C)

Din apelul de recepție se revine chiar dacă nu există nici un mesaj primit. Acest mod de lucru poate conduce la erori și în acest caz se folosesc funcții de detectare a terminării operației: `MPI_Wait()` și `MPI_Test()`. O altă observație importantă este legată de faptul că operațiile send/receive blocante pot fi completate de perechi nebloccante (exemplu în Listing 8).

```
1 // ...
2 MPI_Comm_rank (MPI_COMM_WORLD,&myrank );
3 MPI_Status status; MPI_Request req1;
4 int a,b;
5
6 if (myrank == 0) {
7     a = somevalue();
8     MPI_Isend(&a,1,MPI_INT,msgtag,MPI_COMM_WORLD,&req1);
9     calcul();
10    MPI_Wait(&req1, &status);
11 } else if (myrank == 1) {
12    MPI_Irecv(&b,1,MPI_INT,0,msgtag,MPI_COMM_WORLD,&req1);
13    MPI_Wait(&req1, &status);
14 }
15 // ...
```

Listing 8: Exemplu de utilizare – send nebloccant/recv blocant

Funcția `MPI_Wait()` se termină după finalizarea operației, iar `MPI_Test()` se termină imediat și returnează un indicator a cărui stare arată dacă operația s-a încheiat sau nu.

În afară de aceste funcții mai există `MPI_Waitany()`, `MPI_Testany()`, `MPI_Waitall()`, `MPI_Testall()`, `MPI_Waitsome()`, `MPI_Testsome()`.

**Pentru mai multe detalii despre aceste tipuri de funcții consultați help-ul distribuției MPI cu care lucrați (ex: `man MPI_Waitsome`).**

**Observație:** *Funcțiile de comunicație fără blocare oferă posibilitatea suprapunerii în timp a operațiilor de comunicație cu cele de calcul, ceea ce poate reduce timpul de execuție al aplicației.*

**Observație:** *Sistemele de comunicații cu transfer de mesaje sunt în general nedeterministe, adică nu garantează că ordinea în care se transmit mesajele este aceeași cu ordinea în care mesajele sunt recepționate, programatorului revenindu-i sarcina găsirii mecanismelor necesare pentru obținerea unui calcul determinist.*

Listing-urile 9 și 10 prezintă o variantă posibilă de implementare a problemei ”producători–consumatori” (cazul unui singur consumator) utilizând comunicații simple blocante și nebloccante.

```
1 // ...
2 typedef struct {
3     char data[MAXSIZE];
4     int datasize;
5     MPI_Request req;
6 } Buffer;
7
8 Buffer *buffer;
9 MPI_Status status;
10
11 // ...
12
13 MPI_Comm_rank(comm, &rank);
14 MPI_Comm_size(comm, &size);
15 if (rank != size - 1) { //codul pentru producator
16     //initializare producatorul alocă un buffer
17     buffer = (Buffer *)malloc(sizeof(Buffer));
18     //producatorul umple bufferul cu date si intoarce numarul
```

```

19      //de bytes din buffer
20      produce( buffer->data , &buffer->datasize );
21      //trimite datele
22      MPI_Send(buffer->data , buffer->datasize , MPI_CHAR,
23              size-1, tag , comm);
24 } else { // rank==size-1; codul pentru consumator
25      //initializare consumatorul alocă un buffer pentru fiecare
26      //producator
27      buffer=(Buffer *)malloc(sizeof(Buffer)*(size-1));
28      // se aplează un receive pentru fiecare producator
29      for(i=0; i<size-1; i++) {
30          MPI_Irecv(buffer[i].data , MAXSIZE, MPI_CHAR, i ,
31                  tag , comm, &(buffer[i].req));
32      }
33      for(i=0; i<size-1; i++) { //bucla principala
34          MPI_Wait(&(buffer[i].req), &status);
35          // găsește numărul de bytes primiți
36          MPI_Get_count(&status , MPI_CHAR, &(buffer[i].datasize));
37          // consumatorul eliberează bufferul de date
38          consume(buffer[i].data , buffer[i].datasize);
39      } //end for
40 } //end if-else

```

Listing 9: Problema producători–consumatori – cu utilizare MPI\_Wait

```

1 // ...
2 typedef struct{
3     char data[MAXSIZE];
4     int datasize;
5     MPI_Request req;
6 } Buffer;
7
8 Buffer *buffer;
9 MPI_Status status;
10
11 // ...
12
13 MPI_Comm_rank(comm, &rank);
14 MPI_Comm_size(comm, &size);
15
16 if (rank!=size-1) { //codul pentru producator
17     //initializare producatorul alocă un buffer
18     buffer= (Buffer *)malloc(sizeof(Buffer));
19     while(1) { //bucla principala
20         //producatorul umple bufferul cu date și întoarce numărul
21         //de bytes din buffer
22         produce(buffer->data , &buffer->datasize);
23         //trimite datele
24         MPI_Send(buffer->data , buffer->datasize , MPI_CHAR,
25                 size-1, tag , comm);
26     } //end while
27 } else { // rank==size-1; codul pentru consumator
28     //initializare consumatorul alocă un buffer pentru fiecare
29     //producator
30     buffer=(Buffer *)malloc(sizeof(Buffer)*(size-1));
31     // se aplează un receive pentru fiecare producator
32     for(i=0; i<size-1; i++)
33         MPI_Irecv(buffer[i].data , MAXSIZE, MPI_CHAR, i ,
34                 tag , comm, &(buffer[i].req));

```



```
35
36     i=0;
37     while(1){ //bucla principala
38         for(flag=0;!flag; i=(i+1)%(size-1)) {
39             MPI_Test(&(buffer[i].req),&flag, &status);
40         }//end for
41         MPI_Get_count(&status, MPI_CHAR, &(buffer[i-1].datasize));
42         // consumatorul elibereaza bufferul de date
43         consume(buffer[i-1].data, buffer[i-1].datasize);
44         //reinitializare
45         //(in cazul in care mai sunt mesaje de la acelasi producator)
46         MPI_Irecv(buffer[i].data, MAXSIZE, MPI_CHAR, i,
47                 tag, comm, &(buffer[i].req));
48     }//end while
49 }//end else
```

Listing 10: Problema producători–consumatori – cu utilizare MPI\_Test