

4.MPI – interfața pentru programarea paralelă a sistemelor cu memorie distribuită

Specificația *MPI (Message Passing Interface)* cunoaște diverse implementări (API-uri MPI) sub forma unor biblioteci ce conțin un set larg (peste 250) de rutine pentru schimbul de mesaje, managementul datelor și proceselor, rutine ce pot fi utilizate pentru un spectru variat de aplicații științifice implementabile eficient pe diverse sisteme paralele și/ sau distribuite.

4.1. Obiective

- Studiul bibliotecii de programare paralelă a sistemelor cu memorie distribuită bazată pe mesaje MPI
- Procese MPI și tipuri de comunicație
- Comunicatori, grupuri de procese, topologii virtuale, comunicația colectivă-mecanisme suport pentru procesare paralelă performantă
- Implementarea unor algoritmi paraleli în modelul transferului de mesaje

4.2. Concepte

Standardul MPI definește o *interfață de programare paralelă prin schimb explicit de mesaje*, ce poate fi utilizată pentru scrierea de aplicații pentru sisteme de tip MIMD (multiple instructions multiple data) cu memorie distribuită / partajată, respectiv pentru sisteme SIMD (Single instruction multiple data) ce presupun paralelizare prin decompoziția datelor. MPI este o bibliotecă de funcții, ea definește nume de funcții, secvențe de apel și rezultatele pentru un set de funcții care pot fi apelate în programe C, C++ și Fortran 90. Compilarea programelor poate fi realizată cu un compilator standard, sau cu diverse alte compilatoare performante, iar în faza de editare a legăturilor este introdusă și biblioteca MPI.

Principalele caracteristici care au determinat o largă răspândire a standardului sunt:

- oferă o interfață de programare pentru aplicații paralele/distribuite
- oferă performanță, portabilitate și ușurință în folosire, semantica sa permite implementări diverse;
- semantica interfeței este independentă de limbaj (C, C++, Fortran)
- permite comunicație eficientă *punct-la-punct* și *punct-multipunct* în medii eterogene (rețele diverse);
- conține un set de rutine pentru lucrul cu *grupuri de procese*, permite scheme de adresare complexe.
- permite *comunicație colectivă* și *operații de calcul global*, ce conferă scalabilitate sistemului
- oferă contexte de comunicare și comunicatori pentru dezvoltarea de biblioteci de funcții necesare programării paralele;
- conține rutine pentru definirea de topologii virtuale de grup;
- oferă simplitate dar și completitudine bazat pe un număr extrem de mare de funcții, însă aplicații relativ complexe pot fi programate folosind un *set minimal de 6 funcții esențiale*.

4.2.1. Modelul de comunicație

Soluția oferită de MPI constă în specificarea mesajelor la un nivel înalt de abstractizare, capabil să reflecte faptul că transmisia unui mesaj necesită o structurare mai complexă decât un simplu string de biți, astfel structura unui mesaj devine *adresa*, *contor*, *tip_date*. Termenul de *buffer de mesaj* poate fi utilizat cu semnificații multiple, astfel:

- poate referi o zonă de memorie *specificată în aplicație* pentru stocarea datelor sau utilizatorul poate seta o zonă de memorie utilizată ca o zonă intermediară pentru memorarea unor mesaje arbitrare ce apar în aplicație
- poate reprezenta o zonă de memorie creată și gestionată de *sistemul MPI*

Sincron vs asincron. Mesajul poate fi transferat direct între bufferele declarate în aplicația utilizator, sau pentru comunicație asincronă, mesajul este *stocat temporar* într-un buffer creat dinamic de sistem înaintea depozitării sale în bufferul receptor, respectiv utilizatorul poate declara inițial un bufer suficient de mare astfel încât să permită transferul oricărui mesaj ce necesită stocare intermediară și utilizare în operația de transfer. Comunicare prin transfer sincron de mesaje definește modul de comunicare ce poate fi regăsit în limbajul teoretic de comunicare prin transfer de mesaje al lui Hoare- CSP. Cele două procese corespunzătoare transmiterii-recepției de mesaje se încheie doar atunci când mesajul transmis este și recepționat. Avantajele pe care le prezintă referă semantica simplă și ușor de utilizat. Un proces care revine din subrutina send va ști cu certitudine că mesajul trimis a fost recepționat. De asemenea un avantaj important îl constituie faptul că nici sistemul și nici utilizatorul nu trebuie să mențină un buffer pentru mesaj, acesta poate fi copiat direct în spațiul de adrese al procesului receptor. Datorită simplității, un număr mare de sisteme paralele implementează diferite variante ale acestui mod de comunicare.

Comunicarea bazată pe emisie - recepție blocantă este caracterizată de transmiterea unui mesaj în mod blocant, executată *fără a aștepta recepția acestuia* de către procesul destinație, astfel este posibilă încheierea execuției unei operații de transmisie chiar înaintea începerii execuției operației de recepție corespunzătoare. Tipurile de comunicare sunt prezentate în tabelul 1.

Evenimentul de comunicație	Sincron	Blocant	Nonblocant
Revenirea din send este indicată	Mesaj recepționat	Mesaj transmis	Transmitere de mesaj inițiată
Bufferare mesaj	Nu e necesară	E necesară	E necesară
Verificarea stării	Nu e necesară	Nu e necesară	E necesară
Încărcare generată de așteptare	Mare	Medie	Mică
Suprapunere comunicație și procesare	Nu	Da	Da

Tabel 4.1. Moduri de comunicare aflate la baza primitivelor MPI

Comunicare bazată pe emisie - recepție neblocantă. Transmisia neblocantă se încheie *imediat după ce a fost notificat sistemul* de existența unui mesaj de transmis, respectiv recepționat, mesaj care poate fi în tranzit, deja sosit sau chiar netransmis încă. Pentru această variantă, sistemul trebuie să implementeze un *buffer temporar*. Implementările performante conțin funcții specializate pentru verificarea *stării bufferului* sau pentru implementarea unor mecanisme de așteptare până ce procesul poate continua (tip *wait-for*). Bufferul necesar transferului de mesaje poate fi implementat de către *sistemul de transfer al mesajelor sau mod utilizator*, la nivelul aplicației.

Avantajele oferite de comunicația neblocantă pot fi exploatate în contextul utilizării suportului hardware de comunicație dedicat, astfel încărcarea indusă de comunicație poate fi aproape integral mascată prin operații neblocante. Operațiile neblocante pot fi utilizate și cu un protocol bufferat, astfel emițătorul inițiază operația DMA și revine, însă datele sunt disponibile accesului (sunt sigure) doar după ce operația DMA s-a încheiat.

Tipuri de date MPI. O caracteristică deosebit de puternică a bibliotecii o constituie posibilitatea introducerii de argumente sub forma unor tipuri de date variate, posibil combinate, pentru toate mesajele trimise și recepționate. MPI asigură un set deosebit de complex de tipuri de date predefinite, set care include toate tipurile de bază din C, C++ și Fortran, alături de două tipuri de date specifice MPI: MPI_BYTE și MPI_PACKED.

Context de comunicație. Fiecare comunicare de mesaj se derulează într-un anumit context. Mesajele sunt întotdeauna primite în contextul în care au fost transmise, iar mesaje transmise în contexte diferite nu interferă.

Proces și grup de procese. Procesul MPI este unitatea fundamentală de calcul, fiind constituit dintr-un *fir de control independent* și un *spațiu de adrese privat*, astfel un proces nu poate accesa direct variabilele definite în spațiul de adrese al altui proces, comunicarea fiind bazată pe transfer de mesaje. Contextul este partajat de un grup de procese. Procesele MPI se execută fără a avea însă mecanisme de încărcare a codului pe procesoare, atribuire a proceselor procesoarelor sau mecanisme de creare și distrugere de procese. Aceste mecanisme sunt de obicei, în cadrul implementărilor curente, *preluate de la sistemul de operare pe care rulează* procesul MPI, fapt ce poate crea dificultăți de portare a acestor biblioteci între sisteme de operare diferite. MPI este proiectat a fi thread-safe, iar grupurile de procese definite în MPI sunt dinamice, apartenența la grup este statică însă suportă atașarea-detașarea din grup. Grupurile se pot suprapune fapt ce permite unui proces să fi membru al mai multor grupuri simultan. Procesele sunt identificate în cadrul grupului printr-un întreg (rank), care ia valori între 0 și n-1, unde n este numărul proceselor din grup. Inițializarea unei aplicații MPI folosind rutina de inițializare MPI_INIT, creează un singur grup, numit MPI_COMM_WORLD, care include toate procesele MPI ce formează aplicația.

Procesele pot utiliza *comunicație punct-la-punct* cu scopul transmisiei- recepționării de mesaje utilă implementării comunicațiilor locale nestructurate. De asemenea un grup de procese poate apela operații de tip *comunicații colective* pentru a simplifica implementarea unor operații globale, deși standardul nu oferă în mod explicit suport pentru multithreading, unele implementări includ mecanisme de acest fel.

Comunicatorul integrează informațiile de context și informațiile de identificare a grupului. MPI_INIT definește pentru fiecare proces apelat comunicatorul implicit MPI_COMM_WORLD. Toate comunicațiile MPI necesită un comunicator ca argument, iar procesele MPI nu pot comunica decât dacă *partajează* un comunicator. Fiecare comunicator identifică un grup reprezentând o listă de procese, fiecare proces este numerotat, identificatorul fiecăruia numindu-se rang (*rank*). Acest rang identifică în mod unic un proces și poate fi folosit pentru specificarea sursei sau destinației unui mesaj. Folosind MPI_COMM_WORLD, fiecare proces poate comunica cu oricare altul, iar grupul MPI_COMM_WORLD reprezintă mulțimea tuturor proceselor MPI. Informațiile specifice comunicatorului pot fi accesibile prin funcțiile:

MPI_COMM_RANK (MPI_comm comm, int *rank)

Funcția MPI_COMM_RANK returnează rangul procesului apelant în grupul comunicatorului *comm*.

MPI_COMM_SIZE (MPI_comm comm, int *size)

Funcția MPI_COMM_SIZE returnează în variabila *size* numărul de procese din grupul asociat comunicatorului *comm*. Fiecare proces aparținând unui comunicator este identificat prin rangul său.

4.2.2. Modelul computațional

Stilul de programare care poate fi folosit cu MPI se bazează fie pe modelul **MIMD** (Multiple instructions multiple data) ce presupune existența de programe diferite pentru fiecare procesor, fie pe modelul **SPMD** (Single Program Multiple Data), caz în care există un singur program, astfel procesele MPI provin din același program, dar execută porțiuni diferite de cod selectate pe baza unor instrucțiuni de condiție. Standardul MPI a fost proiectat astfel încât programatorul începător să nu simtă complexitatea de la primele programe. Programarea unei aplicații minimale MPI se poate face

cunoscând doar un set minimal de 6 funcții ce reprezintă nucleul suficient pentru scrierea unei aplicații:

MPI_Init (int *argc, char**argv) este funcția care inițializează variabilele locale MPI și înștiințează sistemul global MPI despre execuția unui nou proces MPI. Funcția acționează și ca barieră de inițializare, în sensul ca nici un proces MPI nu-și poate continua execuția până când toate celelalte procese nu au executat partea lor de inițializare, argc/argv sunt argumentele la nivel de linie de comanda pentru programul C.

MPI_Comm_rank(), returnează identificatorul procesului MPI în cadrul grupului specificat în lista de parametri.

MPI_Comm_size() întoarce numărul de procese care fac parte din grupul specificat ca parametru

MPI_Send() și **MPI_Recv()** sunt funcțiile de bază care implementează comunicația prin mesaje între procese.

MPI_Finalize() pentru anunțarea sistemului MPI despre terminarea procesului MPI cu succes (valoarea lui ierr=0), această rutină este apelată de fiecare proces MPI. După această funcție nici o altă funcție nu poate fi apelată, toate comunicațiile în așteptare trebuie încheiate înaintea apelului ei.

Apelul funcțiilor de bibliotecă necesită includerea antetului **mpi.h** ce conține constantele referite în rutinele MPI (tipuri de date, comunicatori, erori) și interfețele funcțiilor MPI. O operație este considerată *terminată local* în cadrul unui anumit proces, dacă acesta a terminat partea sa ,atât în comunicația punct-la-punct, cât și în cea colectivă. O operație de comunicare este *terminată global*, dacă toate procesele implicate în comunicare au *terminat local* părțile ce le revin din aceasta.

4.2.3. Modul de comunicație punct-la-punct

O aplicație poate reduce latența de comunicare presupunând că platformele hardware pe care este realizată implementarea includ mecanisme care permit *rutarea și gestiunea mesajelor paralel cu alte procesări*. Biblioteca MPI garantează respectarea următorului set de proprietăți (reguli semantice): pentru comunicarea punct la punct, astfel: garantarea preluării mesajelor în ordinea în care au fost trimise, perechile send și receive nu pot rămâne în permanență nerezolvate (necorelate). Mesajele transmise reprezintă o secvență de articole de același tip, identificarea mesajelor realizându-se prin identificatorul **sursei (rank)** și o **etichetă (tag)** atașată mesajului, valori interpretate relativ, în cadrul zonei de comunicație identificată de un anumit comunicator. Există situații în care este necesar ca recepția să fie efectuată indiferent de sursă (MPI_ANY_SOURCE) sau de etichetă (MPY_ANY_TAG).

Programele paralele necesită implementarea unor soluții algoritmice capabile să minimizeze timpul de răspuns, iar MPI oferă propriile mecanisme pentru instrumentarea codului astfel încât să permită analiza performanței. Funcția **MPI_Wtime**, returnează o valoare dublă precizie a numărului de secunde raportat la un anumit moment de timp trecut, garantat nemodificabil pe durata execuției programului, fapt ce permite inserarea diverselor apeluri ale rutinei în codul sursă pentru măsurarea unor intervale de timp de execuție. Funcțiile bibliotecii MPI utilizate în analiza performanțelor bazat pe elementul timp de procesare sunt **MPI_Wtime (void)**, **MPI_Wtick(void)**. Lista parametrilor utilizați în rutinele MPI și semnificația lor, respectiv lista codurilor de eroare sunt prezentate în tabelul 4.2.:

ierror	Cod de eroare MPI
Comm	Identificarea comunicatorului utilizat în transmisia mesajului
Rank	Rangul procesului apelant în grup sau comunicator
buffer	Adresa de start a buferului ce conține mesajul(transmis-receptionat)
count	Număr de articole specificate în bufer (număr maxim de articole recepționate)
datatype	Tipul de dată a fiecărui element transmis

parametru	Semnificație
Root	Rangul unui proces specific (sursa sau destinația pot fi identice)
Size	Numărul de procese din grup
Tag	Tagul de mesaj transmis-recepționat
source	Rangul procesului sursă în comunicator
dest	Rangul procesului destinație în comunicator
status	Tablou de întregi ce conține starea, referă o modul de recepție a mesajului (sursa, tagul și codul de eroare)

Tabel 4.2. Coduri de eroare și parametrii rutinelor MPI

Pentru a putea construi o aplicație paralelă performantă este necesar un control performant asupra transmiterii mesajelor. În acest sens, MPI definește patru moduri de comunicare: **standard, sincron, bufferat și ready** (vezi tabelul 4.3.). Funcția de comunicare se obține prin combinarea modului de comunicație cu un mod de blocare, moduri ce vor fi detaliate la secțiunea funcții pentru comunicație. Toate modurile de comunicație există atât în forma blocantă cât și neblocantă, în forma blocantă, întoarcerea dintr-o funcție implică terminarea ei cu succes MPI nu impune restricții asupra modului de potrivire a funcțiilor de transmisie cu cele de recepție, existând astfel mai multe posibile combinații a căror funcționare corectă intră însă în sarcina programatorului. MPI definește două moduri de comunicare:

- fără blocare (imediat), caz în care funcțiile de transmisie sau recepție se termină înainte ca informația să fie trimisă/recepționată, iar bufferul indicat într-o funcție neblocantă nu trebuie folosit până la detectarea sfârșitului funcției inițiate.
- cu blocare, caz în care după terminarea funcției, bufferul de transmisie sau recepție poate fi refolosit.

Tip	Descriere
Send sincron	operația de trimitere poate începe oricând, dar nu se poate termina decât atunci când mesajul a ajuns la receptor
Send buferat	utilizatorul poate furniza sistemului un buffer pentru a permite ca operația de trimitere să se termine întotdeauna înainte ca mesajul să fie recepționat la destinație.
Send standard	operația de trimitere poate începe chiar dacă operația corespunzătoare de recepție nu a început.
Send ready	operația send poate începe doar dacă operația corespunzătoare de recepție a început.
Receive	Se completează la sosirea unui mesaj

Tabel 4.3. Moduri de comunicare MPI

Sintaxa rutinelor pentru transmisia mesajelor este:

MPI_Send(void*buffer, int count, MPI_datatype datatype, int dest, int tag, MPI_comm comm)

MPI_Recv(void* buffer, int count, MPI_datatype datatype, int source, int tag, MPI_comm comm, MPI_status status)

După ce mesajul a fost transmis, variabila Status poate fi utilizată pentru a colecta informații referitoare la operația MPI_recv, aceasta este o structură de date cu trei câmpuri: sursa mesajului, tagul și cod de eroare, iar lungimea mesajului transmis poate fi verificată cu ajutorul funcției MPI_get_count.

Transmisie blocantă standard (Standard send). Forma *Standard send* se încheie odată ce mesajul a fost trimis, mesaj care poate să ajungă sau nu la destinație, iar ignorarea acestui fapt poate duce la comportări inconsistente ale programului. Sintaxa:

int MPI_SEND (void *buffer, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)

unde:

buffer este adresa datei ce va fi trimise;

count este numărul de elemente de tipul *datatype* conținute de *buffer*

dest reprezintă procesul destinație specificată prin intermediul rangului din cadrul comunicatorului *comm*;

tag este un marcator la dispoziția programatorului pentru a distinge între diferite tipuri de mesaje;

Transmisie sincronă (Synchronous send). *Synchronous send* se folosește atunci când se dorește confirmarea primirii mesajului, astfel destinația trimite spre sursă o confirmare iar transmisia este considerată încheiată numai după primirea confirmării. Primitiva *send* se consideră încheiată atunci când o operație *receive* a fost lansată, astfel bufferul aplicației este disponibil pentru a recepționa mesajul. Este posibilă o transmisie sincronă neblokantă ce nu presupune transmiterea mesajului, ci doar inexistența bufferării suplimentare la receptor, dar un buffer sistem la emițător poate fi necesar, iar pentru a elimina copieri multiple ale mesajului se poate utiliza emisiia blocantă. Această formă de transfer se folosește atunci când este nevoie de siguranță a transmisiei, pentru a oferi comportament determinist. Sintaxa:

int MPI_SSEND (void * buffer, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)

Transmisie buferată (Buffered send). Forma *Buffered send* folosește un buffer pentru trimiterea mesajului, dacă acesta nu poate fi trimis la momentul dorit. Avantajul față de metoda clasică îl constituie predictibilitatea recepției mesajului. Atașarea unui buffer se realizează cu ajutorul funcției:

int MPI_BUFFER_ATTACH (void *buffer, int size)

Numai un singur buffer poate fi atașat unui proces la un moment dat. Detașarea se realizează cu:

int MPI_BUFFER_DETACH (void *buffer, int size)

Observații:

- după detașare, utilizatorul poate reutiliza sau dealoca spațiul ocupat de tampon
- unele operații, de atașare și detașare, au un argument de tip *void**; ele sunt folosite diferit – un pointer la tampon este transmis la *attach*; adresa pointerului este utilizată la *detach*, astfel că acest apel întoarce valoarea pointerului. Argumentele sunt definite ambele ca *void** (și nu *void** respectiv *void***) pentru a evita conversii forțate de tip (*cast*).

Transmisie ready (Ready send). Formele *Ready send*, ca și *buffered send* se termină imediat, iar comunicarea este garantată a se efectua cu succes dacă o operație *receive* corespunzătoare a fost apelată. La expedierea mesajului, procesul expeditor pune mesajul în comunicator, sperând ca destinatarul să îl aștepte și să îl primească. Dacă acesta nu îl acceptă, mesajul poate fi abandonat sau se generează o eroare.

int MPI_RSEND (void *buffer, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)

Recepție blocantă (Blocking receive). Sintaxa:

int MPI_RECV (void *buffer, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Status status)

unde:

- *source* este rangul expeditorului în grupul asociat comunicatorului *comm*; ca și sursă se poate specifica MPI_ANY_SOURCE, ceea ce este echivalent cu orice sursă;
- *tag* specifică tipul mesajelor ce vor fi acceptate, se poate de asemenea specifica MPI_ANY_TAG ca fiind identificator universal de mesaj.

Informația este returnată de către funcția MPI_RECV în variabila *status*. Acest argument poate fi testat direct pentru a afla sursa mesajului, dimensiunea și identificatorul de mesaj (de exemplu poate fi necesară identificarea sursei unui mesaj recepționat folosind MPI_ANY_SOURCE).

4.2.4. Comunicații colective, comunicatori și topologii virtuale

Mecanismele abstracte folosite în comunicație sunt : comunicatori, context de comunicare, topologii virtuale. În multe aplicații paralele este necesar ca operațiile de comunicație să fie restricționate la un set limitat de procese. MPI definește o serie de concepte și introduce diverse mecanisme pentru lucrul cu procese structurate în grupuri și organizate folosind comunicatorii. MPI introduce noțiunea de *comunicator* pentru a defini cadrul unei comunicații colective, în particular procesele care sunt implicate în această comunicare, dar și pentru a permite construcția modulară a programelor paralele în mod: secvențial, paralel sau-și concurent.

Un *comunicator* definește un *context de comunicare alături de o colecție ordonată de procese* implicate într-un subset al calculului de efectuat. Prin introducerea comunicatorilor, comunicația între subseturi de procese ale programului poate fi structurată. Comunicația între procese din grupuri distincte este posibilă prin definirea unui *inter-comunicator*.

Comunicația inter-grupuri este, de obicei, gestionată de un conducător al grupului. Grupurile pot fi create și prin intermediul reuniunii, intersecției, diferenței între grupuri deja existente. Noțiunea de comunicator este fundamentală pentru definirea într-un mod modular a bibliotecilor de funcții care pot fi invocate de procesele utilizator. Putându-și defini propriul comunicator (alocat de către sistem, cu asigurarea unicității), procesele sistem pot schimba mesaje fără a interfera cu comunicația dintre procesele aplicației utilizatorului.

Este necesară distincția clară între diferitele tipuri de comunicatori. Astfel un *comunicator* este un obiect cu atribute și reguli pentru creare, utilizare și distrugere. Comunicatorul specifică un domeniu de comunicare ce poate fi utilizat de comunicarea punct la punct sau colectivă

- Un *intracomunicator* este utilizat pentru comunicarea în interiorul unui grup (comunicare intra-grup); el are două atribute: *grupul de procese și topologia grupului*.
- Un *intercomunicator* este utilizat pentru comunicarea punct la punct între două grupuri diferite (comunicare inter-grup); atributele lui sunt cele două grupuri, fără a specifica topologia.
- Un *domeniu de comunicare* este un set de comunicatori. Dacă acest domeniu este pentru comunicare intra-grup atunci toți comunicatorii sunt intracomunicatori și au aceleași atribute.

Unui comunicator *i* se pot atașa informații adiționale alături de cele de grup și context, informație numită **topologie**. Topologia reprezintă un mecanism ce realizează asocierea diverselor *scheme de adresare proceselor* aflate într-un grup. Grupurile și comunicatorii sunt obiecte opace, detaliile reprezentării lor interne sunt specifice implementării MPI, iar accesarea lor este realizată folosind un handle. Contextele nu sunt utilizate explicit în funcțiile MPI, ci sunt asociate în mod *implicit* grupurilor de procese la crearea lor.

Procesele unui grup au ranguri de la 0 la $n-1$, n fiind numărul de procese din grup. În multe aplicații, ordonarea liniară a proceselor, după rang, nu reflectă tiparul de comunicare între procese.

Adesea, procesele trebuie aranjate în topologii cu două sau mai multe dimensiuni. În cazul general, tiparul de comunicare între procese corespunde unui graf. Un astfel de aranjament al proceselor, care reflectă comunicările punct la punct dintre ele reprezintă topologia virtuală a grupului de procese. Trebuie făcută distincție între topologia virtuală a proceselor și topologia reală a sistemului pe care acestea sunt executate. Topologia virtuală poate fi exploatată de sistem în plasarea proceselor pe procesoare, pentru a ameliora performanțele.

Trebuie reținut faptul că topologia virtuală reflectă caracteristicile aplicației și poate fi folosită în *îmbunătățirea performanțelor*. Sunt cunoscute tehnici standard de mapare a unor topologii grilă / tor pe hipercuburi sau grile de procesoare. Pentru topologiile graf metodele sunt mai complexe și se bazează adesea pe euristici. În afara informației furnizate algoritmului de plasare a proceselor pe procesoare, topologiile virtuale au ca rol facilitarea scrierii programelor, într-o formă mai ușor de înțeles. Topologia virtuală poate fi considerată ca un atribut deosebit de puternic al unui comunicator, atribut ce permite descrierea structurii de interconectare a proceselor dintr-un grup, ea este folosită pentru a mapa procesele pe arhitectura fizică.

Topologia poate fi definită prin grupul de procese, numărul de conexiuni al fiecărui proces și partenerul pentru fiecare conexiune. Topologia virtuală este distinctă de graful programului, pentru că două procese pot comunica chiar dacă în topologie nu sunt conectate. Topologia virtuală descrie acei comunicatori care trebuie considerați de maparea grupului de procese pe arhitectura fizică. Noțiunea de topologie virtuală este deosebit de utilă, atunci când nu este disponibil un mecanism eficient de mapare a topologiei virtuale pe un calculatormasiv paralel. Astfel, dacă procesele utilizatorului vor comunica în principal cu cei mai apropiați vecini după modelul unui grid bidimensional, se va putea crea o topologie virtuală care să reflecte acest fapt. În acest mod este realizat accesul la rutinele convenabile, de exemplu, calculând rangul oricărui proces dat prin coordonatele sale în grid, luând în considerare condițiile limită și returnând MPI_NULL_PROC dacă se iese din grid. În particular există rutine care calculează rangurile celor mai apropiați vecini, rang ce poate fi folosit ca argument pentru primitivele de transmisie /recepție MPI_SEND, MPI_RECV, MPI_SENDRECV etc. Deși o topologie virtuală scoate în evidență șabloanele utilizate în comunicare într-un comunicator pentru o “conexiune”, orice proces din comunicator poate comunica cu oricare altul.

O topologie virtuală poate fi modelată ca un graf în care nodurile reprezintă procese, iar arcele perechile de procese care comunică. Nu trebuie ca o comunicare să fie precedată de o deschidere explicită de canal. Ca urmare, absența unui arc între două noduri din graf nu înseamnă că procesele respective nu vor putea comunica între ele, aci că o astfel de comunicare nu contează la maparea topologiei virtuale de procese pe o topologie reală de procesoare. Arcele grafului de comunicare nu sunt ponderate. Deși specificarea topologiei în termeni de graf este foarte generală, în multe aplicații specificarea unei topologii mai simple, regulate, este mai convenabilă. De exemplu inel, grile bi- sau tri-dimensionale, tor. Aceste topologii sunt complet definite prin numărul de dimensiuni și prin numărul de procese pe fiecare axă de coordonate. Maparea unor topologii grilă sau tor este mai simplă decât maparea unui graf, motivând astfel interesul MPI pentru tratarea separată, explicită a acestora.

MPI permite o serie largă de topologii virtuale între care cele mai uzuale sunt topologiile de tip graf virtual, topologiile carteziane, s. a. Două noduri în graf sunt conectate dacă pot comunica. Grafurile de procese pot fi utilizate pentru a specifica orice topologie, cele mai utilizate fiind gridurile. O topologie este carteziană dacă fiecare proces este conectat cu vecinii săi printr-un grid virtual. Mecanismul utilizat de MPI pentru a asigura ranguri proceselor dintr-un anumit domeniu de comunicație nu utilizează informație despre rețeaua de comunicație, ci oferă un set de funcții ce permit programatorului *aranjarea proceselor în diverse topologii fără a specifica explicit maparea lor la procesoare*, responsabilitatea unei mapări capabile să reducă costurile transmisiei de mesaje este în totalitate a bibliotecii MPI.

Comunicația colectivă permite difuzarea datelor către toate procesele unui grup, sincronizarea la bariere ca mecanism de control al execuției, distribuția colectarea datelor de la procesele unui grup, operații globale de reducere pe grupuri de procese (suma, minim, maxim) sau comunicație de tip all-to-all. Toate aceste operații se pot împărți în două mari clase: *operații de transfer de date și operații de calcul global*.

Deși operațiile de comunicație colectivă nu acționează asemeni unei bariere ele realizează o *sincronizare virtuală*, iar programul necesită o semantică corectă astfel încât să se manifeste determinist chiar dacă este necesară includerea unor operații de sincronizare globală înainte respectiv după apelul funcției colective. Deoarece operațiile colective sunt virtual sincrone, ele *nu necesită marcajul de mesaj* (tag-ul). Pentru cele mai multe operații colective există două variante astfel: transfer de date de aceeași dimensiune, respectiv transfer de date de dimensiuni diferite, funcție de ce tip și volum de date necesită transfer.

Rutinele de transfer de date referă *difuzarea datelor* (broadcast), *distribuirea* (scatter) și *colectarea rezultatelor* (gather) spre sau de la procesele unui grup. Colectarea de rezultate se poate realiza de către *toate* procesele (de ex. funcția MPI all-gather) sau de către *un singur* proces. De asemenea, există operații *all-to-all*, în care fiecare proces trimite date, respectiv recepționează date de la celelalte procese, prin utilizarea unei singure funcții cu semnătură complexă.

Operațiile de calcul global pot fi împărțite în operațiile de reducere și operații de scanare. Pentru operațiile de reducere funcția de reducere se aplică pe toate datele corespunzătoare fiecărui proces din grup, iar pentru operațiile de scanare, funcția de reducere se aplică pe datele proceselor cu rang mai mic decât cel al procesului care execută operația.

Funcția de reducere trebuie să fie asociativă și comutativă, iar rezultatul unei operații globale poate fi cunoscut de *toate procesele sau doar de un singur proces*. În operațiile globale pot fi folosite atât operațiile predefinite în MPI cât și operații definite de utilizator. În cadrul comunicației colective, nu se folosesc etichete de mesaj, coordonarea comunicației proceselor este implicită și toate rutinele de comunicație colectivă se blochează până când comunicația este terminată local.

Operații colective. Comunicarea colectivă implică un grup de procese între membrii căruia se desfășoară comunicația. Biblioteca MPI suportă comunicații colective cum ar fi difuzarea, *dispersarea/adunarea* (scatter/gather), *schimb total de date*, *agregare și bariera*. Pentru toate formele de comunicare colectivă, toate procesele unui grup trebuie să realizeze apelul corespunzător, cu argumente compatibile. Orice omisiune în acest sens constituie o eroare.

Caracteristici ale comunicațiilor colective:

- comunicațiile colective nu pot interfera cu comunicațiile punct la punct și vice-versa;
- comunicație *colectivă poate sau nu să realizeze sincronizarea* proceselor participante;
- comunicațiile colective sunt blocante;
- *toate procesele participante într-un comunicator trebuie să apeleze comunicațiile colective*;
- mesajul transmis este un vector cu elemente de un anumit tip
- tipul de date trebuie să fie același pentru transmisie și recepție

Sincronizarea la barieră este cea mai simplă operație de comunicare colectivă care nu implică transfer de date. Forma sa este: **MPI_BARRIER (MPI_comm comm)**

Unicul argument este comunicatorul *comm* ce definește grupul de procese sincronizate, astfel procesul apelant se blochează până când toți membrii grupului execută apelul funcției (ating bariera).

Transmisia broadcast. În transmisia *one-to-all*, procesul root transmite același mesaj stocat în bufferul *buffer* (un număr de *count* de intrări de date având tipul *datatype*) tuturor proceselor inclusiv lui însuși, procese incluse în comunicatorul *comm*. Toate procesele implicate în comunicație trebuie să specifice aceeași rădăcină (root, procesul cu identificatorul 0 de regulă).

int MPI_BCAST (void *buffer, int count, MPI_datatype datatype, int source, MPI_comm comm)

- *buffer* este adresa de început a buffer-ului, *count* este numărul intrărilor în buffer
- *datatype* este tipul datelor din buffer, *source* (poate fi cazul particular root) este rangul procesului sursă, iar *comm* este comunicatorul.

Conținutul mesajului este identificat de tripletul *buffer, count, datatype*, pentru procesul sursă acest triplet specifică bufferele de emisie și recepție, iar pentru celelalte procese buferul de recepție. Volumul de date transferat de procesul sursă trebuie să fie identic cu volumul de date recepționat de fiecare proces.

Funcția Gather. Procesul țintă (poate fi chiar rădăcina pentru topologii de tip arbore) recepționează mesaje personalizate de la fiecare din cele n procese (inclusiv de la el însuși). Cele n mesaje sunt *concatenate în ordinea rangurilor proceselor* și stocate în bufferul de recepție al procesului receptor. Fiecare buffer emițător este identificat de tripletul *sendbuf, sendcount, sendtype*, iar bufferul de recepție este ignorat pentru procesele nonroot, iar pentru procesul root este identificat de tripletul *recvbuf, recvcoun, recvtpe*.

Funcția implementează comunicația de tip *all-to-one*, toate procesele transferă datele aflate în *inbuf* procesului rădăcină, care plasează datele în locații contigue nesuprapuse în buferul *outbuf*.

int MPI_GATHER (void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcoun, MPI_Datatype recvtpe, int target, MPI_comm comm)

unde:

- *sendbuf* este adresa de start a send buffer
- *sendcount* este numărul elementelor din send buffer
- *sendtype* este tipul datelor din send buffer
- *recvbuf* este adresa lui receive buffer
- *recvcoun* este numărul elementelor pentru fiecare recepție
- *target* este rangul procesului receptor

Funcția Scatter descrie operația inversă operației GATHER. Sintaxa sa este:

int MPI_SCATTER (void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcoun, MPI_Datatype recvtpe, source, MPI_comm comm).

Pentru lista completă a funcțiilor de comunicație colectivă se vor consulta resursele [1][3].

Pentru *analiza de performanță* și instrumentarea codului aplicației cu scopul de a identifica timpul necesar procesării este utilă funcția *MPI_wtime()* ce returnează valoarea timpului raportat la ceasul sistem. Deasemenea este util Analizorul de concurență pentru Visual Studio care poate fi descărcat ca și componentă distinctă de la adresa <https://docs.microsoft.com/en-us/visualstudio/profiling/concurrency-visualizer?view=vs-2019> și care va fi instalat pentru integrare în Visual Studio- IDE, cu scopul de a vizualiza variația gradului de paralelism pentru fiecare aplicație și consumul de memorie la execuția codului.

4.3.Exemple

4.3.1. Realizați o transformare a programului secvențial pentru calculul lui PI bazat pe integrarea numerică a funcției $f(x)=4/(1-x^2)$ în intervalul $[0, 1]$. Să se ofere o soluție de paralelizare prin împărțirea intervalului de integrare aferent numărului de procesoare din sistem. Pentru algoritm se poate folosi aproximarea funcției cu sumarea $k=1..n$ $4/(1+((k-1)/n)^2)$ și fiecare nod va primi numărul de dreptunghiuri folosite în aproximare, va calcula aria dreptunghiului și apoi se sincronizează pentru calculul sumei globale. Implementați un program pentru calculul valorii lui PI ce poate rula pe un sistem cu un număr arbitrar de procesoare, astfel va trebui ca indiferent de numărul de procesoare, valoarea calculată să fie aceeași (sursa [5]).

```
/* Calculates the PI. It is the area of f(x)=4/(1+x^2) in interval from 0 to 1.The
interval in splitted into n intervals which are distributed equally to all
processes. The result consists in a sum of areas calculated by each processor !!! */
```

```

#include <stdio.h>
#include <math.h>
#include "mpi.h"

/* Function definition */
double f( double a ) {
    return (4.0 / (1.0 + a*a));}

int main( int argc, char *argv[]) {
    int done = 0, n, myid, numprocs, i;
    double PI25DT=3.141592653589793238462643;
    double mypi, pi, h, sum, x;
    double startwtime=0.0, endwtime;
    int namelen;
    char processor_name[MPI_MAX_PROCESSOR_NAME];

    /* Initialize MPI */
    MPI_Init(&argc,&argv);

    /* Get no. of processes */
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);

    /* Get current process ID */
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    /* Get current processor name */
    MPI_Get_processor_name(processor_name,&namelen);

    /* Print info about MPI WORLD and about current process */
    printf("MPI WORLD has %d processes\n", numprocs);
    printf("Process %d on %s\n",myid, processor_name);

    /* Initialize no. of intervals to 0 */
    n = 0;
    while (!done) {
        /* What root process does: */
        if (myid==0) {

            /* Read no. of intervals prom keyboard */
            printf("Enter the number of intervals: (0 quits) ");
            scanf("%d",&n);

            /* Next line should only be used if you are too lazy to enter no.
            of intervals */
            /* if (n==0) n=100; else n=0; */

            /* Get current time */
            startwtime = MPI_Wtime(); }

    /* Send no. of intervals to all processes. Only root sends the message
    and the others receive the message */
    MPI_Bcast(&n, 1, MPI_INT, 0, MPI_COMM_WORLD);

    /* In no intervals defined then stop program */
    if (n==0)
        done=1;
    /* If intervals defined */
    else {
        /* Calculate f(x) where x=1/n*(i-0.5)*/
        h=1.0/(double)n;
        sum=0.0;
        for (i=myid+1;i<=n;i+=numprocs) {
            x=h*((double)i-0.5);
            sum+=f(x);}

        /* Area is 1/n*f(x) */
        mypi=h*sum;
    }
}

```

```

        /* Send calculated area to root through a SUM function. I don't
        know yet if root sends it's value, but it's for sure he receives
        as result the SUM of all numbers from the other processes including
        his own. */
        MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);

        /* What root does */
        if (myid==0)
        {
            /* Print the result of PI*/
            printf("pi is approximately %.16f, Error is %.16f\n",
                pi, fabs(pi - PI25DT));

            /* Get current time and print total time needed */
            endwtime = MPI_Wtime();
            printf("wall clock time = %f\n",
                endwtime-startwtime);}}

    /* Terminate MPI */
    MPI_Finalize();
    /* Success */
    return 0;

```

4.3.2. Se vor testa exemplele ce ilustrează modul de utilizare a API-ului MPI, exemple pe care le puteți descărca de la adresa https://ftp.utcluj.ro/~civan/CPD/1_LABORATOR/04_MPI

4.4. Întrebări teoretice

4.4.1. Ce tipuri de comunicație oferă interfața MPI? Explicați-le succint semnificația și aplicabilitatea.

4.4.2. Identificați minim 3 tipuri de operații de tip SEND și ilustrați modul de utilizare.

4.4.3. Care sunt funcțiile MPI pentru inițializarea/închiderea contextului de execuție

4.4.4. Ce rol joacă un comunicator? Dar un grup de procese?

4.4.5. Ce proces joacă un rol distinct în structurarea și execuția codului? Explicați prin exemplificare cu o secvență de cod, posibila execuție paralelă pe o mașină cu memorie distribuită ce dispune de mai multe procesoare.

4.5. Probleme propuse

4.5.1. Să se implementeze o aplicație simplă folosind transmiterea de mesaje în accepțiunea în care se definesc un proces receptor și n-1 procese emițător. Acestea își transmit identificatorul de proces și numele hostului procesului receptor care le tipărește.

4.5.2. Să se implementeze un program în care sunt definite un număr impar de procese. Procesul zero inițializează un tablou de întregi și distribuie tabloul tuturor proceselor folosind primitiva Scatter. Procesele receptor își primesc părțile corespunzătoare din tabloul sursă. Fiecare proces răspunde cu un mesaj ce conține numele hostului pe care procesul se execută și partea sa din tabloul distribuit anterior, iar procesul rădăcină recepționează și tipărește mesajele.

4.5.3. Să se calculeze produsul a doi vectori (x, y de dimensiune n) distribuiți unui grup de procese. Tablourile ce conțin vectorii sunt alocate static, iar dimensiunea vectorilor este divizibilă cu p numărul de procesoare din sistem. Propuneți două soluții de implementare pentru operația de reducere și analizați modul de execuție și performanța folosind analizorul de concurență ce poate fi descărcat de la adresa <http://msdn.microsoft.com/en-us/library/dd537632.aspx> și integrat în Visual Studio.

4.5.4. Implementați un program care să determine timpul necesar execuției operației MPI_Barrier. Utilizați comunicatorul implicit MPI_COMM_WORLD. Tipăriți dimensiunea comunicatorului

asigurându-vă ca atât emițătorul cât și receptorul sunt pregătiți la începutul testului. Cum variază performanța primitivei variind dimensiunea comunicatorului?

4.5.5. Propuneți un program care să realizeze recepția corectă a mesajelor de la toate procesele emițătoare, astfel încât toate cele 100 procese să transmită procesului 0, utilizând recepția nonblocantă și primitiva `Wait_some`.

4.5.6. Modificați problema conectării procesoarelor unei mașini paralele într-o topologie de tip inel înlocuind referirea proceselor vecine folosind rangurile `rank-1` și `rank+1` folosind rutinele MPI specifice topologice (topologii catenale, liniare) utile pentru “lumi” de procese ce necesită comunicatori distincti.

4.6. Miniproiect

Să se descrie în pseudocod și să se implementeze și evalueze performanța implementării pentru soluții paralele ale algoritmilor de înmulțire matrici optimizați pentru implementare paralelă folosind transferul de mesaje. (algoritmul Fox (<https://www.cs.usfca.edu/~peter/ipp/>) și algoritmul Cannon (<https://bit.ly/2DdmHgr>)). Pentru dezvoltarea aplicației este acceptată orice implementare a bibliotecii din cele existente (MS-MPI, OpenMPI, MPICH2).

Deasemenea, se vor reprezenta grafic evaluările de performanță (Accelerarea S și Eficiența E) realizând un set de execuții consecutive și modificând corespunzător numărul de procese și/sau numărul de procesoare. Se vor analiza rezultatele obținute și se va evalua scalabilitatea.

4.7. Referințe bibliografice

1. API MPI <https://www.mpich.org/static/docs/v3.2/www3/index.htm> (similar MS-MPI)
2. Standard MPI <https://www.mpi-forum.org/docs/>
3. Tutorial MPI <https://computing.llnl.gov/tutorials/mpi/>
4. Implementări diverse :
 - Microsoft MPI** v6 <https://www.microsoft.com/en-us/download/details.aspx?id=47259>
 - Compile și executie*
<http://blogs.technet.com/b/windowshpc/archive/2015/02/02/how-to-compile-and-run-a-simple-ms-mpi-program.aspx>
 - MPI project template (VS2015)*
<https://visualstudiogallery.msdn.microsoft.com/90fb60d4-0b8c-472f-8135-683d3c45f45a>
 - OpenMPI** : <http://www.open-mpi.org/>
 - Instrucțiuni de instalare*
<http://programmers-journal.blogspot.ro/2013/09/netbeans-ide-configuration-for-openmpi.html>
 - MPICH2**: <http://www-unix.mcs.anl.gov/mpi/mpich2>
 - Instrucțiuni de instalare*
http://yazid.blog.umd.edu.my/files/2010/09/05-2-MPICH_VS2008manual.pdf
5. Programare paralelă- manual open source Universitatea Princeton :
https://princetonuniversity.github.io/PUbootcamp/sessions/parallel-programming/Intro_PP_bootcamp_2018.pdf