

# *Procesare Paralelă*

## *Capitolul 2: Modele de programare: SPMD*

Bogdan Dumitrescu

Facultatea de Automatică și Calculatoare

Universitatea Politehnica București

# Cuprins

---

- Modelul SPMD (Single Program Multiple Data)
- Comunicație prin mesaje
- Message Passing Interface
- Comunicație globală

# Modelul SPMD

---

- Model adecvat calculatoarelor MIMD
- Ne vom ocupa doar de cazul MIMD cu memorie distribuită
- Paradigma SPMD (Single Program Multiple Data): toate procesoarele execută același program, dar fiecare utilizează un set propriu de date
- Precizări importante:
  - Execuția instrucțiunilor nu este sincronă, adică fiecare procesor execută instrucțiunile în ritmul său propriu
  - Nu toate procesoarele execută aceleași instrucțiuni; deși programele încărcate în memoriile procesoarelor sunt identice, totuși pot exista diferențe la execuție

## Adresa unui procesor

- Un program se execută pe  $p$  procesoare ale unui calculator paralel
- In general, valoarea lui  $p$  este mai mică decât numărul total de procesoare al calculatorului
- Procesoarele sunt numerotate de la 0 la  $p - 1$
- Numerotarea nu este statică, ci se face în momentul lansării în execuție
- Procesoarele pot să-și afle adresa proprie printr-o primitivă a sistemului de operare sau de bibliotecă
- In descrierea algoritmilor, adresa procesorului care execută programul este utilizată ca o variabilă
- Procesoarele pot executa instrucțiuni diferite în funcție de adresa lor

## Diferențierea la execuție

- Notăm  $id$  adresa proprie
- Exemplu simplu de program
  1.  $id \leftarrow \text{adresa proprie}$
  2. **dacă**  $id = 3$  **atunci**
    1.  $a \leftarrow 1$
  3. **altfel**
    1.  $a \leftarrow 0$
- Doar procesorul  $P_3$  execută instrucțiunea  $a \leftarrow 1$
- Toate celelalte execută  $a \leftarrow 0$
- Atenție: variabila  $a$  este în memoria locală a fiecărui procesor!

## Variabile și indici

- Așadar, o variabilă a unui program SPMD este multiplicată: fiecare procesor deține un exemplar, asupra căruia are control complet
- Un procesor nu poate modifica variabile din memoria altui procesor
- În exemplul de mai sus, putem interpreta variabila  $a$  ca un vector cu  $p$  elemente, fiecare procesor deținând doar un element, anume cel al cărui indice este egal cu adresa procesorului
- Programul de mai sus inițializează vectorul  $a$  cu zero, mai puțin elementul  $a_3$  care devine 1
- Vom folosi indici locali în programe
- Totuși, în prezentarea unor algoritmi vom folosi și indici globali

## Problemă

---

- Vrem să inițializăm un vector  $a$  de dimensiune  $n$  cu zero, mai puțin elementul  $a_m$ , care să fie 1
- Intrebare 1: cum distribuim vectorul  $a$  celor  $p$  procesoare ?
- Intrebare 2: cum facem efectiv inițializarea ?
- Răspunsuri: la curs

## Model de comunicație prin mesaje

- Model exclusiv pentru calculatoare MIMD cu memorie distribuită
- Deoarece fiecare procesor are memoria sa proprie, singura modalitate de comunicare între procesoare este transmiterea de mesaje
- Operația de bază este cea în care un procesor sursă  $P_s$  transmite un mesaj  $M$  conținând date din memoria sa  $M_s$  unui procesor destinație  $P_a$ , care stochează datele în memoria sa  $M_a$
- Când vom discuta algoritmi de comunicație vom presupune că sursa și destinația sunt vecine
- Vom scrie astfel algoritmi eficienți pe anume topologii



## Primitive de bază

- Sunt suficiente doar două rutine pentru a scrie orice operație de comunicație
- Procesorul sursă transmite prin rutina **send**
- Procesorul destinație recepționează prin rutina **recv**
- Sintaxa

*send(date, dest)*  
*recv(date, sursă)*

- *date* este o variabilă locală, indicând locul în care se află datele transmise sau unde se memorează datele recepționate
- Lungimea datelor va rezulta din context
- Presupunem implicit că datele ocupă o zonă contiguă în memorie

## Primitive de bază

- Al doilea parametru identifică vecinul cu care se comunică
- De obicei vom nota vecinul prin adresă (ca în MPI)
- Alternativ, vom preciza direcția în care se găsește vecinul
- Pe inel, grilă sau tor: stânga, dreapta, sus, jos, vest, est, nord, sud
- Pentru hipercub: dimensiunea pe care se comunică, un număr între 0 și  $d - 1$
- De exemplu, transmiterea liniei  $i$  a matricei  $A$  către vecinul din dreapta al procesorului ce execută operația se scrie

`send( $A(i, :)$ , dreapta)`

## Corectitudinea comunicației

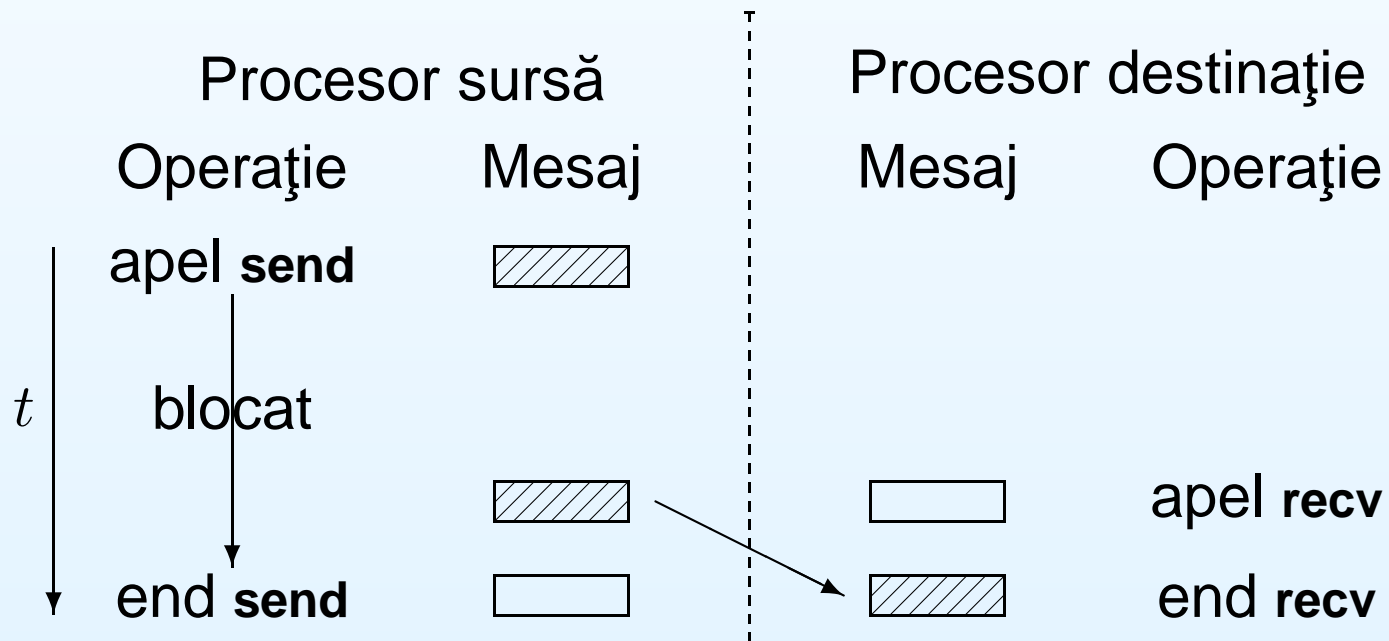
- Orice operație de transmisie a unui procesor este însoțită de una de recepție a unui vecin al său
- Primitivele **send** și **recv** trebuie să apară în perechi, pe ansamblul procesoarelor
- Exemplu: să presupunem că procesorul  $P_k$  trebuie să transmită un același mesaj  $M$  vecinilor săi pe un inel,  $P_{k-1}$  și  $P_{k+1}$
- Algoritmul are forma
  1. **dacă**  $id = k$  **atunci**
    1. **send**( $M$ , dreapta)
    2. **send**( $M$ , stânga)
  2. **altfel dacă**  $id = (k - 1) \bmod p$  **atunci** **recv**( $M$ , dreapta)
  3. **altfel dacă**  $id = (k + 1) \bmod p$  **atunci** **recv**( $M$ , stânga)

## Terminarea locală a comunicației

- Să presupunem că procesorul  $P_s$  transmite un mesaj procesorului  $P_a$
- Momentele în care  $P_s$  apelează primitiva **send**, iar  $P_a$  **recv**, sunt în general diferite
- Transmiterea efectivă a mesajului se face doar după ce ambele procesoare au apelat primitivele respective
- Se pune întrebarea: ce face procesorul care a apelat primul primitiva sa de comunicație, din momentul apelului și până când celălalt procesor apelează primitiva pereche ?
- Răspunsuri posibile:
  - așteaptă (fără să facă nimic); este cazul comunicației *blocante*
  - poate executa alte operații; comunicația este *non-blocantă*

# Comunicație blocantă

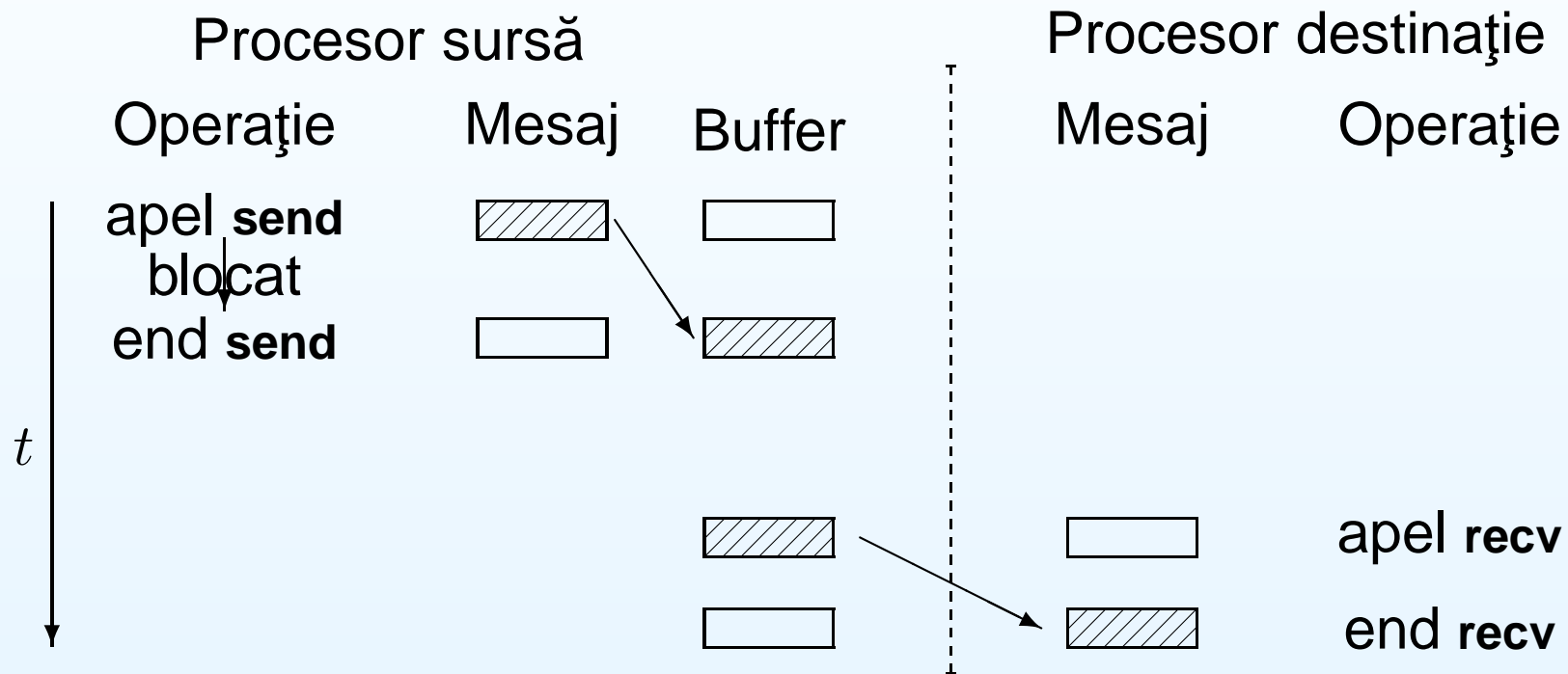
- Al doilea procesor îl așteaptă pe primul
- Cât așteaptă, nu face nimic
- Comunicația se numește blocantă și sincronă (procesoarele termină simultan)



## Comunicație blocantă prin buffer

- Comunicația blocantă are loc printr-un *buffer*
- Primitiva **send** mută mesajul din zona sa originală într-un buffer; după aceasta, execuția **send** se termină, zona de memorie a mesajului putând fi refolosită
- Procesorul sursă este blocat doar în timpul în care mesajul este copiat în buffer
- Terminarea rutinei **send** la sursă nu este condiționată de apelul **recv** la destinație
- Atât în varianta sincronă cât și în cea prin buffer, la terminarea execuției primitivei **send**, zona de memorie în care era memorat mesajul poate fi refolosită, iar la terminarea execuției **recv**, mesajul este recepționat în zona de memorie alocată acestui scop

# Comunicație blocantă prin buffer—ilustrație



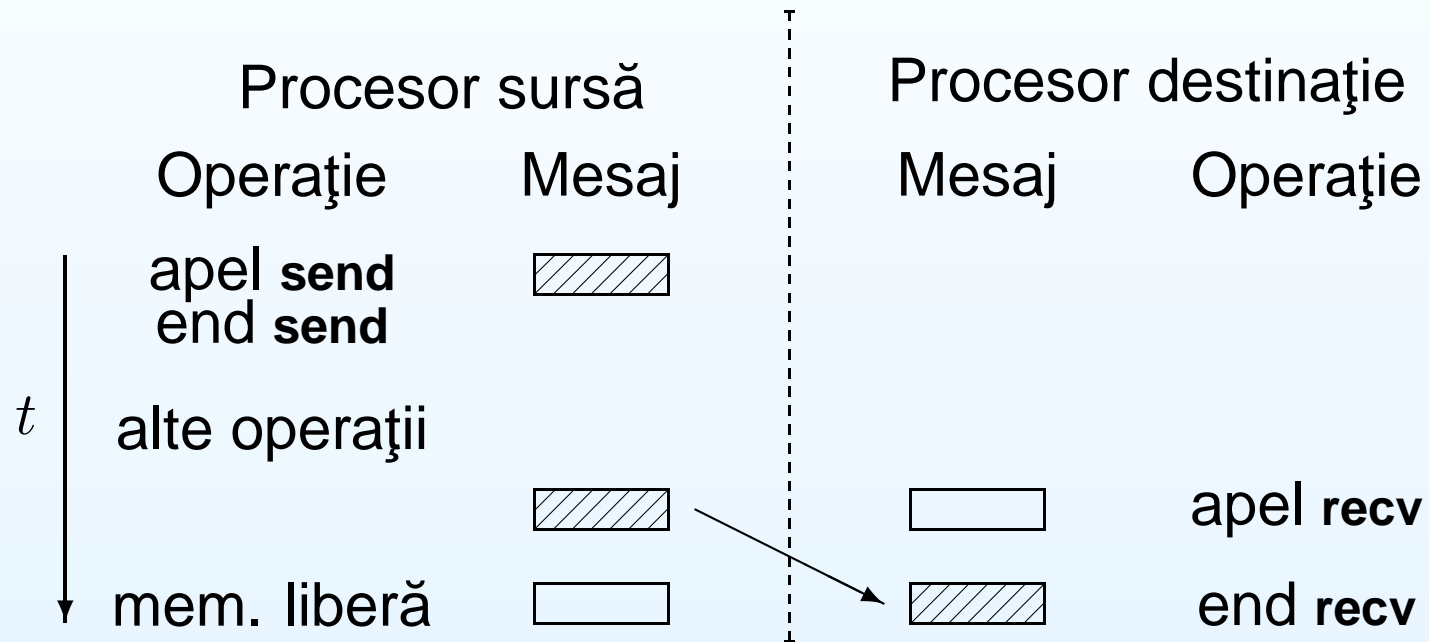
## Comunicație non-blocantă

---

- Primitivele **send** și **recv** se termină *imediat* după apel, fără ca zona de memorie a mesajului să fie liberă la **send** sau să conțină mesajul la **recv**
- Primitivele au doar rolul de a iniția comunicația, transferul efectiv fiind realizat mai târziu la un nivel inferior
- După terminarea execuției **send** sau **recv**, procesorul poate executa alte operații, în paralel sau concurent cu comunicația
- Și în modul non-blocant comunicația se poate desfășura sincron sau prin buffer



# Comunicație non-blocantă—ilustrație



## Primitive auxiliare pentru comunicația non-blocantă

- Comunicația non-blocantă trebuie să fie însoțită de posibilitatea de a afla când se termină comunicația efectivă
- În acest scop se introduce primitiva **așteaptă**, de așteptare a terminării comunicației
- Utilizare tipică
  1. **send**(*date*, dest)
  2. execută operații care nu modifică *date*
  3. **așteaptă** terminarea **send**
  4. modifică *date*
- Dacă **așteaptă** este executată imediat după **send** sau **recv**, efectul este același ca în cazul primitivelor blocante
- Altă variantă: primitivă care verifică terminarea comunicației

## Comparații

- Comunicația non-blocantă permite în general o mai bună ocupare a procesoarelor, evitând unii timpi morți
- Erorile de programare în modul blocant se manifestă de obicei prin blocarea programului
- Erorile de programare în modul non-blocant pot fi mai subtile, deoarece se poate altera informația transmisă, prin rescrierea zonei de memorie la `send` sau prin utilizarea datelor înainte de a fi recepționate efectiv
- (În general, depanarea programelor paralele este mult mai dificilă decât în cazul secvențial)

## Eroare tipică în comunicația blocantă

- Problemă: fiecare procesor transmite vecinului din dreapta un mesaj, pe un inel de procesoare
- Echivalent, fiecare procesor recepționează un mesaj de la vecinul din stânga
- Notăm  $d1$  și  $d2$  variabilele locale alocate mesajului propriu, respectiv celui recepționat
- Program incorect în modul blocant sincron:
  1. **send**( $d1$ , dreapta)
  2. **recv**( $d2$ , stânga)
- Fiecare procesor începe prin a transmite și apoi așteaptă procesorul din dreapta să execute recepția
- Acesta însă așteaptă și el, etc.
- Toate procesoarele se blochează

# Soluție

- Soluție 1: transmisia se face prin buffer
- Operația **send** se termină local la transferul mesajului în buffer (independent de acțiunile procesorului destinație), deci recepția poate să înceapă
- Soluție 2: o parte din procesoare încep prin a transmite, celelalte prin a recepționa, de exemplu
  1. **dacă** id este par
    1. **send**( $d1$ , dreapta)
    2. **recv**( $d2$ , stânga)
  2. **altfel**
    1. **recv**( $d2$ , stânga)
    2. **send**( $d1$ , dreapta)

## Soluția în cazul non-blocant

---

- În modul non-blocant, programul unui procesor poate fi
  1. **send**( $d1$ , dreapta)
  2. **recv**( $d2$ , stânga)
  3. **așteaptă** terminarea **send** și **recv**
- Eventual alte operații pot fi intercalate înainte de **așteaptă**, pentru a folosi timpul posibil mort în care are loc transferul efectiv al mesajului

# Standardul MPI

---

- Message Passing Interface (MPI) este un standard descriind primitive de comunicare—în contextul unui model de programare SPMD
- Este implementat pe toate calculatoarele MIMD și pe multe alte arhitecturi (de exemplu, rețele de calculatoare)
- Pentru Windows: Open MPI, MPICH2 și altele (WMPI)
- Avantaje
  - Portabilitate (același program MPI) se poate compila și executa pe multe calculatoare
  - Comunicația este implementată eficient pe fiecare calculator
- Se pot testa și depana programele în medii "ieftine" (pe un desktop sau o mică rețea), rulându-le pe supercalculatoare în momentul în care sunt aproape sigur funcționale

## MPI—generalități

- O rutină MPI se execută în cadrul unui grup de procesoare, numit *comunicator*
- Inițial, există doar comunicatorul `MPI_COMM_WORLD`, conținând toate procesoarele pe care se execută programul
- În cadrul unui comunicator sunt  $p$  procesoare, numerotate de la 0 la  $p - 1$
- Un procesor poate afla numărul de procesoare din comunicator și adresa proprie (numită *rang* în MPI) prin

```
int MPI_Comm_size(MPI_Comm com, int *p)
int MPI_Comm_rank(MPI_Comm com, int *my_id)
```



## MPI—generalități

- `MPI_Comm` este tipul predefinit pentru comunicatoare
- `com` este comunicatorul în care se află procesorul
- În variabilele `p` și `my_id`, rutinele returnează valorile numărului de procesoare, respectiv a adresei procesorului apelant
- Toate rutinele MPI întorc un întreg caracterizând succesul execuției

## Structura unui program MPI

---

- Încadrându-se în modelul SPMD, un program MPI are forma obișnuită a unui program secvențial
- Rutinele `MPI_Init` și `MPI_Finalize` trebuie apelate înainte, și respectiv după, orice alte rutine MPI
- `MPI_Init` primește argumentele `argc` și `argv`, acestea având sau nu o semnificație în funcție de implementarea MPI
- Aceasta este singura constrângere, în rest programatorul este liber să folosească rutinele MPI potrivite algoritmului pe care îl implementează

## Rutine de comunicație

- Oricare două procesoare își pot transmite mesaje
- Deci topologia virtuală este graful complet conectat
- Principalele rutine de comunicație sunt

```
MPI_Send(void *mes, int lung,  
          MPI_Datatype tipdate, int dest,  
          int eticheta, MPI_Comm com)
```

```
MPI_Recv(void *mes, int lung,  
          MPI_Datatype tipdate, int sursa,  
          int eticheta, MPI_Comm com,  
          MPI_Status *stare)
```

- Prin `MPI_Send` un procesor trimite procesorului `dest` mesajul aflat în memoria locală la adresa `mes`, de lungime `lung` și având tipul `tipdate`

## Rutine de comunicație

- MPI are identificatori predefiniți pentru practic toate tipurile uzuale de date; alte tipuri se pot construi cu rutine speciale
- Mesajului îi este asociată o etichetă (tag), care îl personalizează
- Transmisia mesajului se face în cadrul comunicatorului `com`, care trebuie să conțină și procesorul sursă, și cel destinație
- La `MPI_Recv` parametri similari
  - se recepționează un mesaj de lungime cel mult `lung`, trunchiind eventual mesajul transmis
  - tipul de date poate diferi la sursă și destinație
  - după recepție, variabila `stare` dă informații suplimentare despre mesajul primit, ca de exemplu lungimea efectivă a mesajului

## Moduri de comunicație

- MPI prevede rutine pentru toate tipurile de comunicație descrise anterior: blocantă sau non-blocantă, prin buffer sau sincron
- Modul de implementare a rutinelor de bază (standard) `MPI_Send`, `MPI_Recv` nu este precizat
- De exemplu, este posibil ca mesajele scurte să fie transmise prin buffer, iar cele lungi sincron; de aceea programatorul trebuie să-și ia precauțiile necesare
- Pentru comunicația non-blocantă
  - rutina `MPI_Wait` așteaptă terminarea transmisiei sau recepției
  - `MPI_Test` verifică dacă transmisia sau recepția s-au terminat sau nu

## Rutine MPI de comunicație

---

Comunicație	blocantă	non-blocantă
standard	MPI_Send MPI_Recv	MPI_Isend MPI_Irecv
prin buffer	MPI_Bsend	MPI_Ibsend
sincronă	MPI_Ssend	MPI_Issend

## Exemple de programe MPI

---

- Vezi pag.41–43 din carte
- Programele vor fi discutate la calculator

## Alte funcții MPI

---

- Standardul MPI conține numeroase alte rutine:
  - Comunicație globală: operații de comunicație implicând toate procesoarele dintr-un grup (comunicator): sincronizare, difuzare, distribuție, etc.
  - Alte operații globale: calculul maximumului unor valori distribuite tuturor procesoarelor unui grup, calculul sumei, etc.
  - Topologii virtuale. Pentru creșterea performanțelor unui program MPI pe un anumit calculator, programatorul poate descrie o topologie virtuală care să corespundă topologiei reale a calculatorului
  - Procese: se pot crea dinamic, deci se pot îmbina paralelismul și concurența



## Operații de comunicație globală

---

- Comunicații la care participă toate procesoarele (toate sau cele dintr-un grup bine precizat)
  - sincronizare (barieră)
  - difuzare (broadcast, one-to-all)
  - difuzare generală (all-to-all broadcast)
  - distribuție (difuzare personalizată, scattering)
  - colectare (gathering)
  - schimb complet (multidistribuție, transpunere)

# Sincronizare

---

- Nu este o operație de comunicație propriu-zisă
- Dacă fiecare procesor dintr-un grup apelează primitiva de sincronizare, atunci toate procesoarele așteaptă până la apelul ultimului dintre ele
- Apoi toate continuă execuția cu instrucțiunea următoare
- Cu alte cuvinte, procesoarele așteaptă la barieră până când pot trece toate (aproximativ) simultan
- In MPI funcția este `MPI_Barrier`

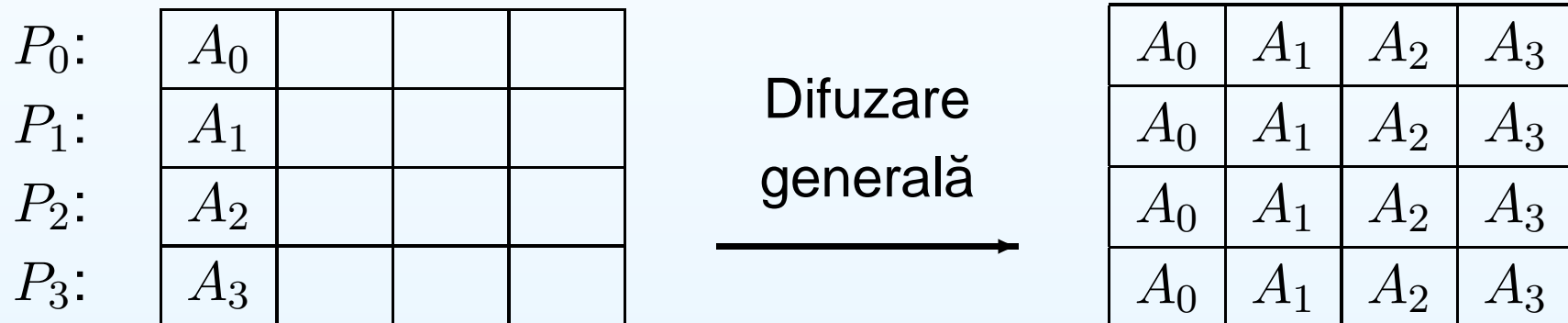
# Difuzare

- Un procesor trimite un (singur) mesaj tuturor celorlalte procesoare



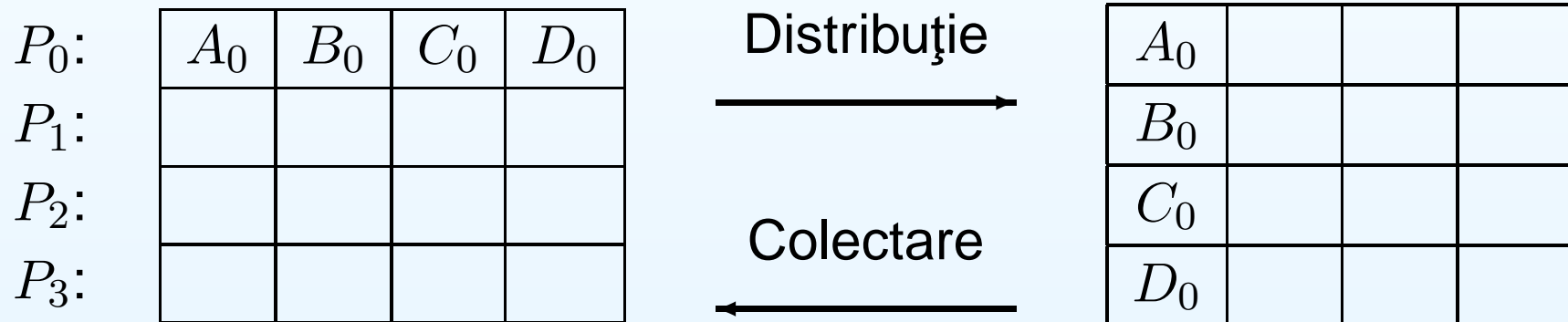
# Difuzare generală

- Fiecare procesor trimite un același mesaj tuturor celorlalte procesoare
- Echivalent, fiecare procesor efectuează o difuzare



## Distribuție și colectare

- Distribuție: un procesor trimite câte un mesaj fiecărui alt procesor
- Colectare: operația inversă, prin care un procesor primește câte un mesaj de la toate celelalte procesoare



# Schimb complet

- Fiecare procesor trimite câte un mesaj fiecărui alt procesor
- Echivalent, fiecare procesor efectuează o distribuție (sau o colectare)
- Dacă fiecare procesor deține inițial o linie a unei matrice  $p \times p$ , în final are o coloană, de unde numele de transpunere

$P_0$ :	$A_0$	$B_0$	$C_0$	$D_0$
$P_1$ :	$A_1$	$B_1$	$C_1$	$D_1$
$P_2$ :	$A_2$	$B_2$	$C_2$	$D_2$
$P_3$ :	$A_3$	$B_3$	$C_3$	$D_3$

Schimb complet



$A_0$	$A_1$	$A_2$	$A_3$
$B_0$	$B_1$	$B_2$	$B_3$
$C_0$	$C_1$	$C_2$	$C_3$
$D_0$	$D_1$	$D_2$	$D_3$