

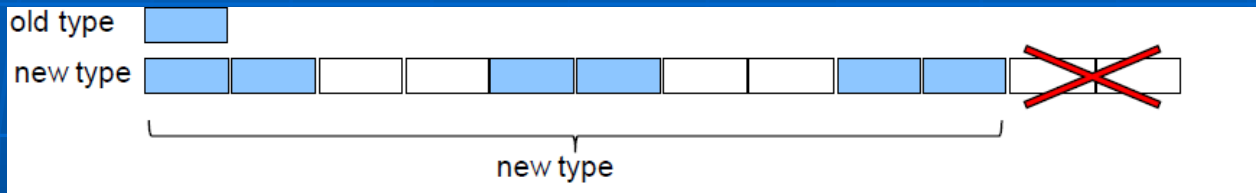
Paralelni sistemi: MPI-Izvedeni tipovi podataka i komunikatori



Prof. dr Natalija Stojanović

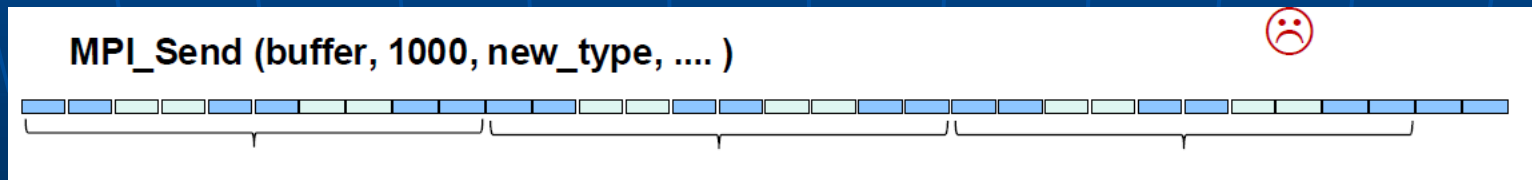
Izvedeni tipovi podataka

Ako želimo da u nekoj od funkcija za komunikaciju koristimo izvedeni tip podatka a da je pritom $\text{count} > 1$, kao na primer `MPI_Send(buff, 10, izvedeni_tip, ...)` moramo voditi računa o tome šta će u tom slučaju biti poslato i da li je to ono što želimo da pošaljemo. Npr. funkcijom **`MPI_Type_vector(3, 2, 4, oldtype, &newtype)`** kreiramo tip:



extent = 10 x extent of “old type”

Ako želimo da sa `MPI_Send` pošaljemo $\text{cnt} > 1$ podataka `newtype`, ono što će biti poslato je:



tj. slanje sledećeg podatka tipa `newtype` kreće sa adrese koja je odmah iza poslednjeg poslatog bajta podatka. Najčešće, ovo nije ono što mi želimo da bude poslato.

Izvedeni tipovi podataka

Ako je ono što je potrebno da bude poslato sa MPI_Send:

```
MPI_Send (buffer, 1000, new_type, .... )
```



Problem možemo rešiti korišćenjem funkcije

```
int MPI_Type_create_resized(MPI_Datatype oldtype, MPI_Aint  
lb, MPI_Aint extent, MPI_Datatype *newtype)
```

lb -nova donja granica tipa (namanji pomeraj u novom tipu, uglavnom je jednak onom u starom tipu, tj 0)

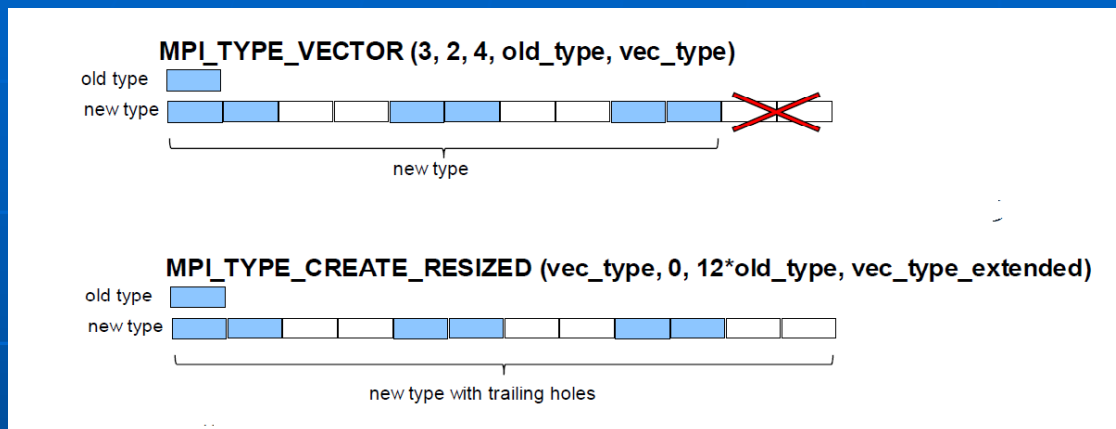
extent -nova veličina tipa koja utiče na to odakle će krenuti slanje sledeće jedinice novog tipa

oldtype-stari tip

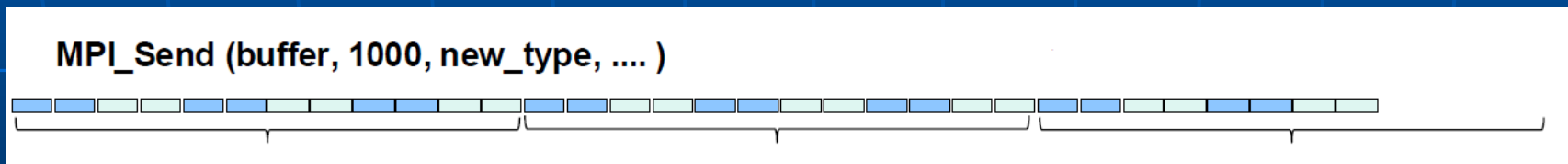
newtype-novi tip

Izvedeni tipovi podataka

Dakle, uradićemo sledeće



Ono što će biti poslato sa MPI_Send je

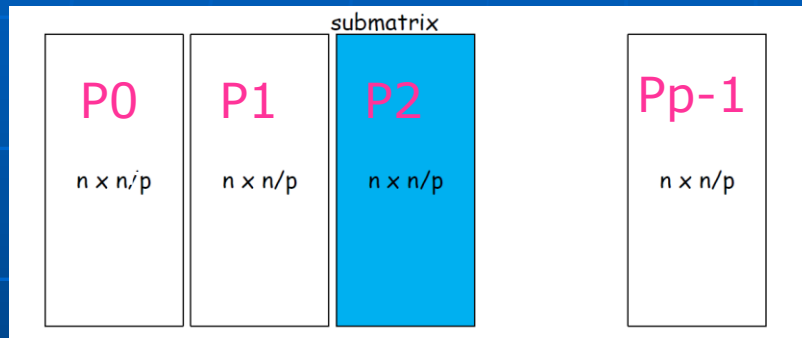


Možemo da zaključimo da MPI_Send(buff, cnt, ndt,...) praktično implementira for (i=0;i<cnt;i++) MPI_Send(buff[i*extent] , 1, ndt,...) .

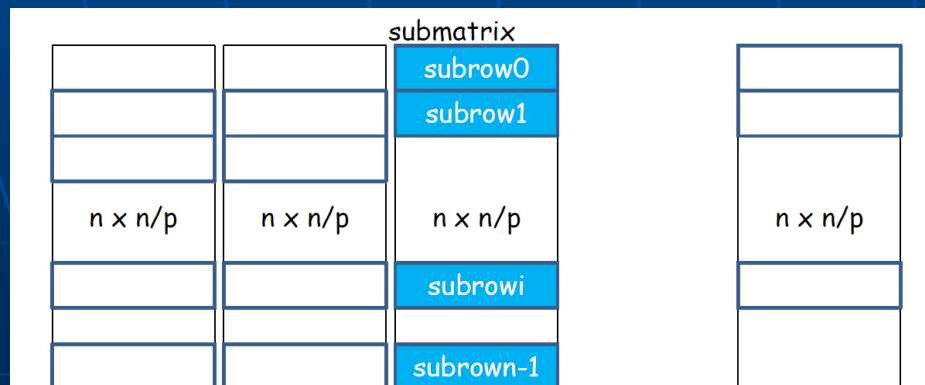
Sva ova zapažanja važe za korišćenje i drugih funkcija za komunikaciju (Scatter, Gather,...)

Izvedeni tipovi podataka

Primer. Želimo da pošaljemo iz master procesa (npr. P0) po n/p (n -dimenzija matrice, p -broj procesa, n deljivo sa p) kolona matrice A svakom procesu i da se u svakom procesu (pa i u P0) te vrednosti nakon slanja nadju u polju niz. U ovu svrhu možemo koristiti funkciju `MPI_Scatter`.

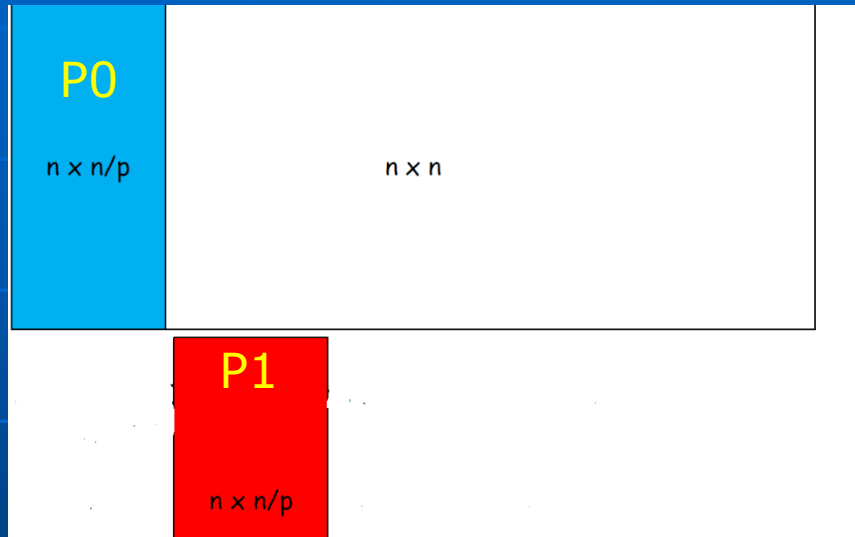


Da bi poslali n/p kolona odjednom kreiraćemo izvedeni tip `colT`:



Izvedeni tipovi podataka

Ono što bi se redom slalo iz svakog procesa sa
`MPI_Scatter(a,1,colT,niz,n*n/p,MPI_DOUBLE,0, MPI_COMM_WORLD)` je:



Pošto to ne odgovara onom što smo želeli da bude poslato, koristimo:

```
MPI_Type_vector(n,n/p,n,MPI_DOUBLE,&colT);  
MPI_Type_create_resized(colT,0,n/p*sizeof(double),&subT);  
MPI_Type_commit(&subT);
```

koji nam omogućava da sa `MPI_Scatter(a,1,subT,niz,n*n/p,MPI_DOUBLE,0, MPI_COMM_WORLD)` pošaljemo ono što smo hteli.

Zad. Napisati MPI program koji realizuje množenje matrice $A_{m \times n}$ i matrice $B_{n \times k}$, čime se dobija i prikazuje rezultujuća matrica $C_{m \times k}$. Izračunavanje se obavlja tako što master proces šalje svakom procesu po jednu kolonu matrice A i po jednu vrstu matrice B. Svi elementi kolone matrice A se šalju odjednom. Svi procesi učestvuju u izračunavanjima potrebnim za generisanje rezultata programa. Zadatak rešiti isključivo primenom grupnih operacija

$$m=2 \ n=3 \ k=4$$

$$\begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \end{bmatrix} * \begin{bmatrix} b_{00} & b_{01} & b_{02} & b_{03} \\ b_{10} & b_{11} & b_{12} & b_{13} \\ b_{20} & b_{21} & b_{22} & b_{23} \end{bmatrix} =$$

$$\begin{bmatrix} a_{00}b_{00} + a_{01}b_{10} + a_{02}b_{20} & a_{00}b_{01} + a_{01}b_{11} + a_{02}b_{21} & \dots & \dots \\ a_{10}b_{00} + a_{11}b_{10} + a_{12}b_{20} & a_{10}b_{01} + a_{11}b_{11} + a_{12}b_{21} & \dots & \dots \end{bmatrix}$$

Zad. Napisati MPI program koji realizuje množenje matrice $A_{m \times n}$ i matrice $B_{n \times k}$, čime se dobija i prikazuje rezultujuća matrica $C_{m \times k}$. Izračunavanje se obavlja tako što master proces šalje svakom procesu po jednu kolonu matrice A i po jednu vrstu matrice B. Svi elementi kolone matrice A se šalju odjednom. Svi procesi učestvuju u izračunavanjima potrebnim za generisanje rezultata programa. Zadatak rešiti isključivo primenom grupnih operacija

```
#include <stdio.h>
#include <mpi.h>
#define m 2
#define n 3
#define k 4
void main(int argc, char* argv[])
{int a[m][n],b[n][k],c[m][k],local_c[m][k],root=0,niza[m],nizb[k],rez[m],rank,i,j,p;
MPI_Datatype vector,column,dt;
MPI_Init(&argc, &argv);
MPI_Comm_rank(MPI_COMM_WORLD,&rank);
MPI_Comm_size(MPI_COMM_WORLD,&p);
```



```

MPI_Type_vector(m,1,n,MPI_INT,&vector);
MPI_Type_commit(&vector);
MPI_Type_create_resized( vector, 0, 1*sizeof(int),&column );
MPI_Type_commit(&column);
if (rank == root)
{
    for(i = 0; i < m; i++)

        for(j = 0; j < n; j++)
            a[i][j] = i+j;

    for(i = 0; i < n; i++)
        for(j = 0; j < k; j++)
            b[i][j] = 1+j-i;
}
MPI_Scatter(a,1,column,niza,m,MPI_INT,0,MPI_COMM_WORLD);
MPI_Scatter(b,k,MPI_INT,nizb,k,MPI_INT,0,MPI_COMM_WORLD);
for (i=0;i<m;i++)
    for (j=0;j<k;j++)
        local_c[i][j]=niza[i]*nizb[j];
for (i=0;i<m;i++)
    for (j=0;j<k;j++)
        c[i][j]=0;

```

```
MPI_Reduce(&local_c[0][0],&c[0][0],m*k,MPI_INT,MPI_SUM,0,MPI_COMM_WORLD;
if (rank==0)
{
    for (i = 0; i< m; i++)
    {
        for (j = 0; j< k; j++)
            printf("%d ",c[i][j]);

        printf("\n");
    }

    MPI_Finalize();
}
```

Sa koliko procesa startujemo aplikaciju?

Grupe i komunikatori

- U nekim slučajevima potrebno je da se komunikacija obavi samo između određene grupe procesa osnovne grupe. To se može izvesti kreiranjem nove grupe procesa korišćenjem odgovarajućih MPI funkcija.
- Grupa predstavlja uređen skup procesa. Grupa je uvek vezana za neki komunikator. Komunikator sadrži grupu procesa koja može da komunicira.
- Pozivi MPI funkcija za komunikaciju između procesa koriste komunikator kao argument. MPI funkcije vezane za grupe specificiraju koji će se procesi koristiti prilikom pravljenja komunikatora.
- Procesi mogu biti članovi više od jedne grupe tj. komunikatora. Unutar svake od njih imaju jedinstveni identifikator.

Identifikator procesa

```
MPI_Comm communicator;  
int myrank;  
MPI_Comm_rank(communicator, &myrank);
```

Grupe i komunikatori-funkcije

Navešćemo neke od funkcija koje se koriste u tu svrhu:

int MPI_Comm_group(MPI_Comm comm, MPI_Group *group)

–kao argument uzima komunikator i vraća odgovarajuću grupu procesa. Ova funkcija je važna za kasnije formiranje novih grupa od osnovne grupe

int MPI_Group_rank(MPI_Group group, int *rank)

–vraća identifikator procesa unutar grupe

int MPI_Group_size(MPI_Group group, int *size)

–vraća broj procesa unutar grupe

Grupe i komunikatori-funkcije

```
int MPI_Group_excl( MPI_Group group, int count, int  
*nonmembers, MPI_Group *new_group )
```

-vraća novu grupu tako što se iz stare grupe isključe procesi sa identifikatorima koji su definisani sa **nonmembers**, i kojih ima **count**. Redosled procesa u novoj grupi prati redosled procesa u staroj grupi.

```
int MPI_Group_incl( MPI_Group old_group, int count,  
int *members, MPI_Group *new_group )
```

-vraća novu grupu tako što procesi sa identifikatorima iz stare grupe koji su definisani sa **members** čine novu grupu. Procesu u novoj grupi ima **count**. Proces `members[i]` u novoj grupi ima rang `i`.

Grupe i komunikatori-funkcije

int MPI_Group_intersection(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)

-vraća novu grupu koja se sastoji od procesa preseka grupa group1 i group2, s tim što je redosled procesa u novoj grupi kao u prvoj grupi

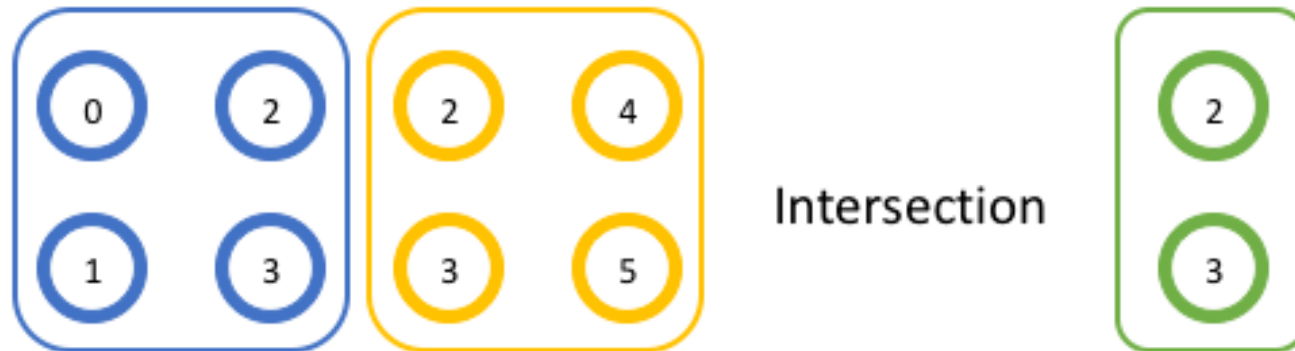
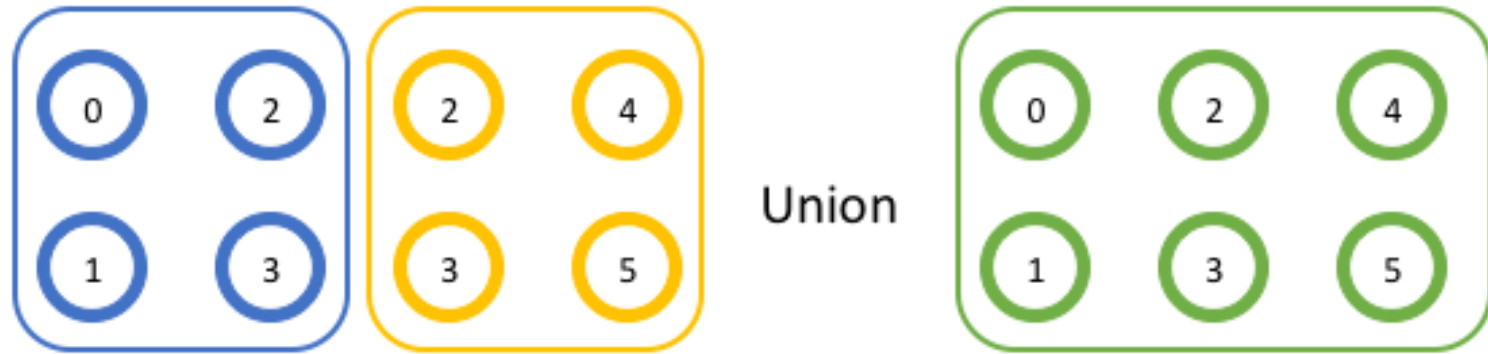
int MPI_Group_union(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)

-vraća novu grupu koja se sastoji od procesa grupe group1 na koju se nadovezuju elementi group2 koji nisu group1

int MPI_Group_difference(MPI_Group group1, MPI_Group group2, MPI_Group *newgroup)

-vraća novu grupu koja se sastoji od procesa grupe group1 koji nisu u grupi group2, i uređeni su kao u group1.

Grupe i komunikatori-funkcije



Grupe i komunikatori-funkcije

```
int MPI_Comm_create( MPI_Comm old_comm,  
MPI_Group group, MPI_Comm * new_comm)
```

-kreira novi komunikator new_comm od procesa iz grupe group koja je kreirana iz old_comm

```
int MPI_Comm_split( MPI_Comm old_comm, int color,  
int key, MPI_Comm *new_comm )
```

-razbija komunikator old_comm na više delova tako sto svi procesi koji imaju istu vrednost za color pripadaju istom podkomunikatoru. Key određuje redosled procesa u podkomunikatorima. Ako je key isti za sve procese onda se redosled procesa preuzima iz starog komunikatora

Grupe i komunikatori-funkcije

int MPI_Comm_split(MPI_Comm, int color, int key, MPI_Comm *new_comm)

-razbija komunikator old_comm na više delova tako sto svi procesi koji imaju istu vrednost za color pripadaju istom podkomunikatoru. Key određuje redosled procesa u podkomunikatorima. Ako je key isti za sve procese onda se redosled procesa preuzima iz starog komunikatora

Rank	0	1	2	3	4	5	6	7	8	9	10
Process	a	b	c	d	e	f	g	h	i	j	k
Color	U	3	1	1	3	7	3	3	1	U	3
Key	0	1	2	3	1	9	3	8	1	0	0

Grupe i komunikatori-funkcije

int MPI_Comm_split(MPI_Comm, int color, int key, MPI_Comm *new_comm)

-razbija komunikator old_comm na više delova tako sto svi procesi koji imaju istu vrednost za color pripadaju istom podkomunikatoru. Key određuje redosled procesa u podkomunikatorima. Ako je key isti za sve procese onda se redosled procesa preuzima iz starog komunikatora

Rank	0	1	2	3	4	5	6	7	8	9	10
Process	a	b	c	d	e	f	g	h	i	j	k
Color	U	3	1	1	3	7	3	3	1	U	3
Key	0	1	2	3	1	9	3	8	1	0	0

{i, c, d}
{k, b, e, g, h}
{f}

Grupe i komunikatori-zadaci

Zad. Napisati MPI program kojim se vrši podela procesa članova komunikatora MPI_COMM_WORLD u dve grupe: grupu procesa sa neparnim identifikatorima i grupu procesa sa parnim identifikatorima

```
#include <mpi.h>
#include <stdio.h>
void main(int argc, char* argv[]) {
    MPI_Group group_world, odd_group, even_group;
    int i, p, Neven, Nodd, members[8], group_rank1, group_rank2, rank;
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_group(MPI_COMM_WORLD, &group_world);
    Neven = (p+1)/2;
    Nodd = p - Neven;
    for (i=0; i<Neven; i++) {
        members[i] = 2*i;
    };
    MPI_Group_incl(group_world, Neven, members, &even_group);
    MPI_Group_rank(even_group, &group_rank1);
    MPI_Group_excl(group_world, Neven, members, &odd_group);
```

Grupe i komunikatori-zadaci

```
MPI_Group_rank(odd_group, &group_rank2);  
printf("moj rank je: %d, moj even rank je %d, moj odd rank je  
%d\n",rank,group_rank1,group_rank2);  
MPI_Finalize();  
}
```

Izlaz:

```
moj rank je: 4, moj even rank je 2, moj odd rank je -32766  
moj rank je: 6, moj even rank je 3, moj odd rank je -32766  
moj rank je: 5, moj even rank je -32766, moj odd rank je 2  
moj rank je: 2, moj even rank je 1, moj odd rank je -32766  
moj rank je: 0, moj even rank je 0, moj odd rank je -32766  
moj rank je: 3, moj even rank je -32766, moj odd rank je 1  
moj rank je: 1, moj even rank je -32766, moj odd rank je 0
```

Grupe i komunikatori-zadaci

Rešenje zadatka na drugi način:

```
#include <mpi.h>
#include <stdio.h>
void main(int argc, char* argv[]) {
int i, p,color,rank,rankc,size,key;
MPI_Comm newcomm;
MPI_Init(&argc,&argv);
MPI_Comm_size(MPI_COMM_WORLD, &p);
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
color=rank%2;key=7;
MPI_Comm_split(MPI_COMM_WORLD,color,key,&newcomm);
MPI_Comm_size(newcomm, &size);
MPI_Comm_rank(newcomm, &rankc);
printf("moj rank je: %d, moj key je %d, moj rankc je %d moj color je
%d moj size je %d\n",rank,key,rankc,color,size);
MPI_Finalize();
}
```

Grupe i komunikatori-zadaci

Izlaz:

moj rank je: 2, moj key je 7, moj rankc je 1 moj color je 0 moj size je 4
moj rank je: 3, moj key je 7, moj rankc je 1 moj color je 1 moj size je 3
moj rank je: 0, moj key je 7, moj rankc je 0 moj color je 0 moj size je 4
moj rank je: 5, moj key je 7, moj rankc je 2 moj color je 1 moj size je 3
moj rank je: 6, moj key je 7, moj rankc je 3 moj color je 0 moj size je 4
moj rank je: 4, moj key je 7, moj rankc je 2 moj color je 0 moj size je 4
moj rank je: 1, moj key je 7, moj rankc je 0 moj color je 1 moj size je 3

Grupe i komunikatori-zadaci

```
#include "stdio.h"
#include "mpi.h"
void main(int argc, char *argv[])
{
    int mcol, irow, jcol, p;
    MPI_Comm row_comm, col_comm, comm2D;
    int Iam, row_id, col_id;
    mcol=2;
    /* Starts MPI processes ... */
    MPI_Init(&argc, &argv); /* starts MPI */
    MPI_Comm_rank(MPI_COMM_WORLD, &Iam); /* get current
process id */
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    irow = Iam/mcol; /* row number: 0 0 1 1 2 2 */
    jcol = Iam%mcol; /* column number: 0 1 0 1 0 1 */

    comm2D = MPI_COMM_WORLD;
    MPI_Comm_split(comm2D, irow, jcol, &row_comm);
    MPI_Comm_split(comm2D, jcol, irow, &col_comm);
}
```


Grupe i komunikatori-zadaci

```
MPI_Comm_rank(row_comm, &row_id);  
MPI_Comm_rank(col_comm, &col_id);
```

```
    printf("%8d %8d %8d %8d %8d\n", Iam, irow, jcol, row_id, col_id);  
MPI_Finalize();  
}
```

Figure a.
2D logical Grid

(0)	(1)
(2)	(3)
(4)	(5)

Figure b.
3 Row Subgrids

(0) (0)	(1) (1)
(2) (0)	(3) (1)
(4) (0)	(5) (1)

Figure c.
2 Column Subgrids

(0) (0)	(1) (0)
(2) (1)	(3) (1)
(4) (2)	(5) (2)

<i>Iam</i>	0	1	2	3	4	5
<i>irow</i>	0	0	1	1	2	2
<i>jcol</i>	0	1	0	1	0	1

- Napisati MPI program koji realizuje množenje matrica A i B reda n, čime se dobija rezultujuća matrica $C=A*B$. Množenje se obavlja tako što master proces (sa identifikatorom 0) šalje svakom procesu radniku jednu kolonu matrice A i jednu vrstu matrice B. Master proces ne učestvuje u izračunavanju. Štampati dobijenu matricu.

$$m=2 \ n=3 \ k=4$$

$$\begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \end{bmatrix} * \begin{bmatrix} b_{00} & b_{01} & b_{02} & b_{03} \\ b_{10} & b_{11} & b_{12} & b_{13} \\ b_{20} & b_{21} & b_{22} & b_{23} \end{bmatrix} =$$

$$\begin{bmatrix} a_{00}b_{00} + a_{01}b_{10} + a_{02}b_{20} & a_{00}b_{01} + a_{01}b_{11} + a_{02}b_{21} & \dots & \dots \\ a_{10}b_{00} + a_{11}b_{10} + a_{12}b_{20} & a_{10}b_{01} + a_{11}b_{11} + a_{12}b_{21} & \dots & \dots \end{bmatrix}$$

- Napisati MPI program koji realizuje množenje matrica A i B reda n, čime se dobija rezultujuća matrica $C=A*B$. Množenje se obavlja tako što master proces (sa identifikatorom 0) šalje svakom procesu radniku jednu kolonu matrice A i jednu vrstu matrice B. Master proces ne učestvuje u izračunavanju. Štampati dobijenu matricu.

```
#include <stdio.h>
#include <mpi.h>
# define n 3
void main(int argc, char* argv[])
{
    int rank,size,m_rank;

    int a[n][n],b[n][n],c[n][n], vrsta[n], kolona[n], tmp[n][n],rez[n][n];
    int root=0;
    int no[1]={0};
    int br=1;
    int i,j;

    MPI_Status status;
    MPI_Datatype vector;
    MPI_Group mat,world;
    MPI_Comm com;
```

```
MPI_Init(&argc, &argv);
MPI_Comm_rank(MPI_COMM_WORLD,&rank);
MPI_Comm_size(MPI_COMM_WORLD,&size);

MPI_Type_vector(n,1,n,MPI_INT,&vector);
MPI_Type_commit(&vector);

MPI_Comm_group(MPI_COMM_WORLD,&world);

MPI_Group_excl(world,br,no,&mat);
MPI_Comm_create(MPI_COMM_WORLD,mat,&com);
MPI_Group_rank(mat,&m_rank);
```

```
if (rank == root)
{
    for(i = 0; i < n; i++)
    {
        for(j = 0; j < n; j++)
        {
            a[i][j] = 3*i+j+1;
            b[i][j] = i+1;
        }
    }
}
```

```

for (i= 0;i < n ; i++)
{
    MPI_Send(&a[0][i],1,vector,i+1,33,MPI_COMM_WORLD);//svakom procesu
    odgovarajucu kolonu

    MPI_Send(&b[i][0],n,MPI_INT,i+1,32,MPI_COMM_WORLD);//svakom procesu
    odgovarajucu vrstu
}
else
{
    MPI_Recv(&kolona[0],n,MPI_INT,root,33,MPI_COMM_WORLD,&status);
    MPI_Recv(&vrsta[0],n,MPI_INT,root,32,MPI_COMM_WORLD,&status);

    for (int i= 0;i<n; i++)
    {
        for (int j= 0;j<n; j++)

            tmp[i][j] = kolona[i]*vrsta[j];
    }

    MPI_Reduce(&tmp,&rez,n*n,MPI_INT,MPI_SUM,0,com);
}

```

```
if (m_rank == 0)
{
    for (i = 0; i < n; i++)
    {
        for (j = 0; j < n; j++)
        {
            printf("%d ", rez[i][j]);
        }
        printf("\n");
    }

    MPI_Finalize();
}
```

- Izračunati standardnu devijaciju skupa od N brojeva korišćenjem formule:

$$s = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2}$$

gde je x_i element skupa od N brojeva, a \bar{x} srednja vrednost svih N brojeva. U procesu P0 se obavlja inicijalizacija N brojeva a zatim se svim procesima distribuira po jednak broj brojeva (N je deljivo sa ukupnim brojem procesa). Nakon toga se izračunava srednja vrednost svih brojeva a zatim i standardna devijacija brojeva. Svi procesi učestvuju u izračunavanjima i rezultat se prikazuje u procesu P0. Zadatak rešiti isključivo korišćenjem grupnih operacija.

```
void main(int argc, char** argv) {
int world_rank,I, world_size;
float global_sum, local_sum, mean, local_sq_diff, global_sq_diff, stddev ;
MPI_Init(&argc,&argv);
MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
MPI_Comm_size(MPI_COMM_WORLD, &world_size);
num_elements_per_proc=num_elements/world_size;
float *rand_numsp = (float *)malloc(sizeof(float) * num_elements);
float *rand_nums = (float *)malloc(sizeof(float) * num_elements_per_proc);
    if(world_rank==0)
        for (i=0;i<num_elements-1;i++)
            rand_numsp[i]=i+1;

MPI_Scatter(rand_numsp,num_elements_per_proc,MPI_FLOAT,rand_nums,num
_elements_per_proc, MPI_FLOAT,0, MPI_COMM_WORLD);
local_sum = 0;

for (i = 0; i < num_elements_per_proc; i++)
    local_sum += rand_nums[i];
MPI_Reduce(&local_sum, &global_sum, 1, MPI_FLOAT,
MPI_SUM,0,MPI_COMM_WORLD);
MPI_Bcast(&global_sum, 1, MPI_FLOAT, 0,MPI_COMM_WORLD);
mean = global_sum / (num_elements);
```



```
local_sq_diff = 0;
for (i = 0; i < num_elements_per_proc; i++)
    local_sq_diff += (rand_nums[i] - mean) * (rand_nums[i] - mean);

MPI_Reduce(&local_sq_diff, &global_sq_diff, 1, MPI_FLOAT, MPI_SUM, 0,
MPI_COMM_WORLD);
if (world_rank == 0) {
    stddev = sqrt(global_sq_diff /
(num_elements_per_proc * world_size));
    printf("Mean - %f, Standard deviation = %f\n", mean, stddev);
}

MPI_Finalize();
}
```