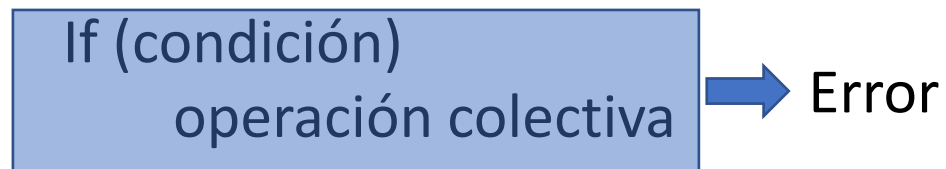


Transparencias MPI

Comunicaciones Colectivas

Operaciones de Comunicación Colectiva

- Involucran a todos los procesos del comunicador: todos ellos deben ejecutar la operación
- Operaciones disponibles:
 - Sincronización (Barrier)
 - Difusión (Bcast)
 - Reparto (Scatter)
 - Recogida (Gather)
 - Multi-recogida (Allgather)
 - Todos a todos (Alltoall)
 - Reducción (Reduce)
 - Prefijación (Scan)
- Algunas de estas operaciones tienen un argumento (root) el cual tiene un papel especial
- Algunas de estas operaciones tienen el sufijo “v”, el cual indica que cantidad de datos de una variable es distinta para los diferentes procesos
- Un error muy común consiste en utilizar una operación colectiva dentro de una sentencia del tipo **if** cuya condición dependa del identificador del proceso, ya que todos los procesos del comunicador deben realizar esa operación esa operación colectiva:



Sincronización

MPI_Barrier(MPI Comm comm)

- Todos los procesos del comunicador se detienen hasta que todos ellos han invocado esta operación
- Ejemplo: medición de tiempos

```
MPI_Barrier(comm);
```

```
t1 = MPI_Wtime();
```

```
/*
```

```
...
```

```
*/
```

```
MPI_Barrier(comm);
```

```
t2 = MPI_Wtime();
```

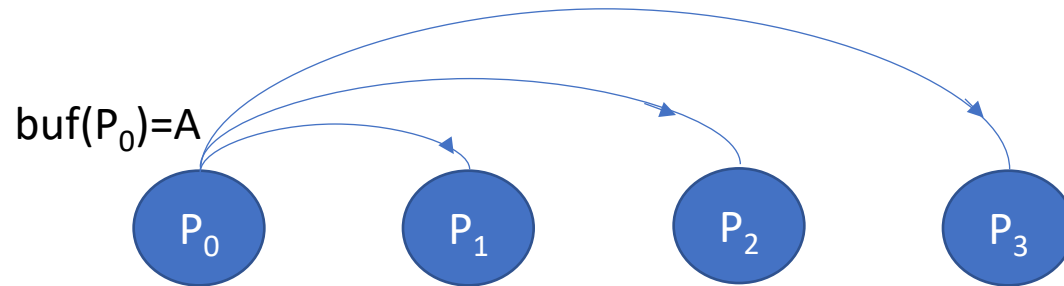
```
if (!rank) printf("Tiempo transcurrido: %f s.\n", t2-t1);
```

Difusión

MPI_Bcast(void *buf, int count, MPI_Datatype datatype, int root, MPIComm comm)

- buf=variable local que contiene el dato aportado por el proceso root
- count=número de elementos de la variable buf
- datatype=tipo de datos de buf
- root=proceso que realiza la difusión. Puede ser cualquiera, por ejemplo, el proceso P_0
- Comm= comunicador (habitualmente el comunicador universal MPI_COMM_WORLD)

Ejemplo:



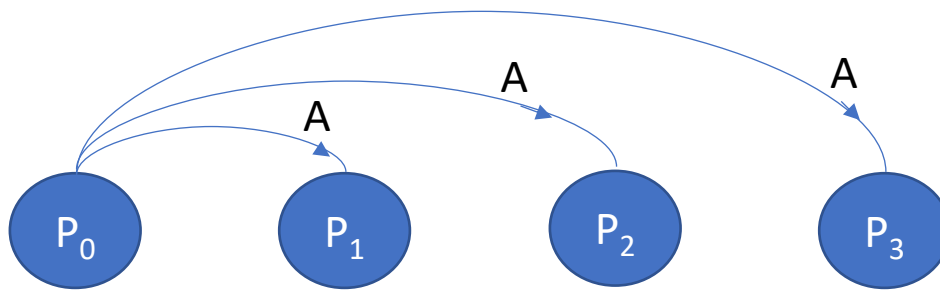
count = $|A| = n^\circ$ datos de A
root = 0

	Antes		Después
buf(P_0)	A	buf(P_0)	A
buf(P_1)		buf(P_1)	A
buf(P_2)		buf(P_2)	A
buf(P_3)		buf(P_3)	A

Ejercicio 1

Supongamos que en MPI se dispone de p procesos y que el proceso P_0 contiene una matriz A de dimensión $M \times N$ de tipo double que debe ser difundida al resto de procesos. Usa la llamada a la función de comunicación colectiva adecuada ¿Cómo podrías sustituir esa comunicación colectiva mediante comunicaciones punto a punto? Calcula el tiempo de comunicaciones.

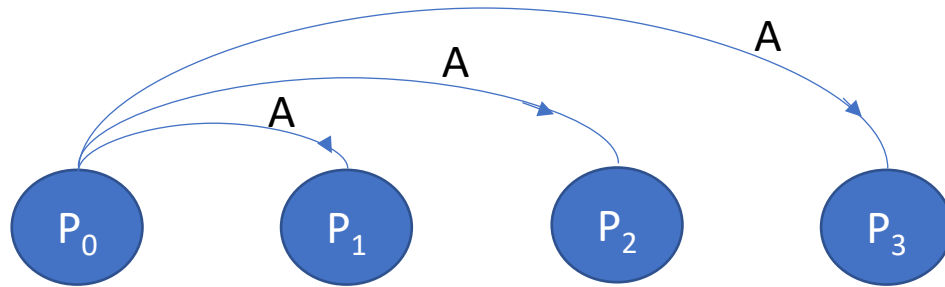
$$A(P_0) \xrightarrow[\text{MPI_Bcast}]{\text{Difusión}} A(P_1), A(P_2), A(P_3)$$



```
int MPI_Bcast(void *buf, int count,  
MPI_Datatype datatype, int root,  
MPI_Comm comm)
```

MPI_Bcast(A, M*N, MPI_DOUBLE, 0, MPI_COMM_WORLD)

Ejercicio 1



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

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

Implementación usando comunicaciones punto a punto:

```
double A[M][N];
```

```
MPI_Init(&argc, &argv);
```

```
int id, p;
```

```
MPI_Status stat;
```

```
MPI_Comm_size(MPI_COMM_WORLD, &p);
```

```
MPI_Comm_rank(MPI_COMM_WORLD, &id);
```

```
if (id==0) /*P0 envía la matriz A al resto de procesos*/
```

```
    for(i=1; i<p; i++)/*P0 envía la matriz A al resto de procesos*/
```

```
        MPI_Send(A, M*N, MPI_DOUBLE, i, 100, MPI_COMM_WORLD);
```

```
else /* El resto de los procesos reciben la matriz A */
```

```
    MPI_Recv(A, M*N, MPI_DOUBLE, 0, 100, MPI_COMM_WORLD, &stat);
```

```
MPI_Finalize();
```

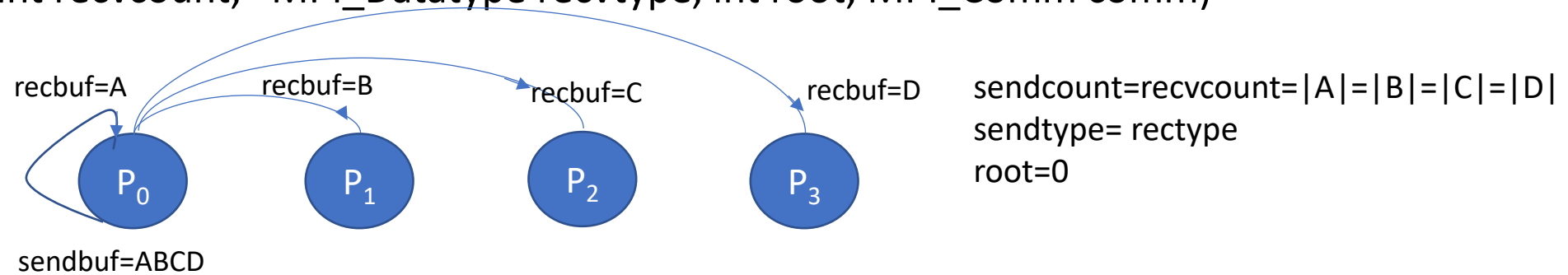
Tiempo de comunicaciones:

$$t_c = (p - 1)(t_s + MNt_w)$$

Reparto: El proceso root reparte un vector/matriz entre todos los procesos

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

Ejemplo:



	Antes		Después
sendbuf(P ₀)	ABCD	recvbuf(P ₀)	A
sendbuf(P ₁)		recvbuf(P ₁)	B
sendbuf(P ₂)		recvbuf(P ₂)	C
sendbuf(P ₃)		recvbuf(P ₃)	D

Ejercicio 2

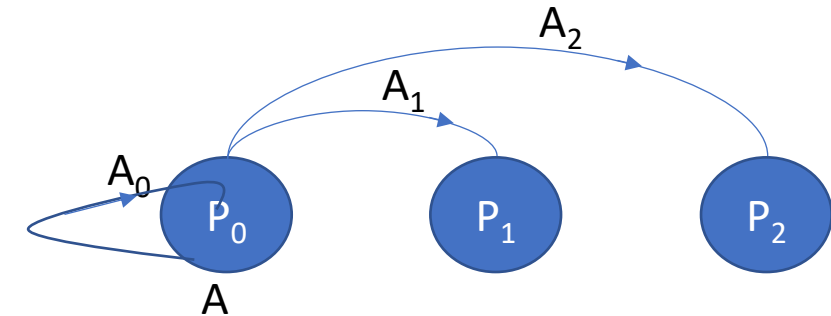
Supongamos que en MPI se dispone de p procesos y que el proceso P_0 contiene a una matriz A de dimensión $M \times N$ de tipo double, la cual debe ser repartida entre todos los procesos. Usa una operación de comunicación colectiva para realizar este reparto ¿Cómo podrías sustituir esa operación mediante comunicaciones punto a punto, suponiendo que el número de procesos p divide a M ? Calcula el tiempo de comunicaciones.

Nota: Las matrices locales quedarán almacenadas en las matrices A_i . Por comodidad, supondremos que esas matrices son también de dimensión $M \times N$.

$$A(P_0) = \begin{matrix} N \\ \left[\begin{array}{c} A_0 \\ A_1 \\ A_2 \end{array} \right] k \\ k \end{matrix} \xrightarrow{\text{Reparto MPI_Scatter}} \begin{matrix} N \\ \left[\begin{array}{c} A_i(P_0) = A_0 \\ A_i(P_1) = A_1 \\ A_i(P_2) = A_2 \end{array} \right] k \\ k \end{matrix}$$

$k = M / p$

($k = n^\circ$ de filas de $A_i = n^\circ$ de filas de $A_i(P_i)$, $p = \text{número de procesos}$)



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

`MPI_Scatter(A, k*N, MPI_DOUBLE, A_i, k*N, MPI_DOUBLE, 0, MPI_COMM_WORLD)`

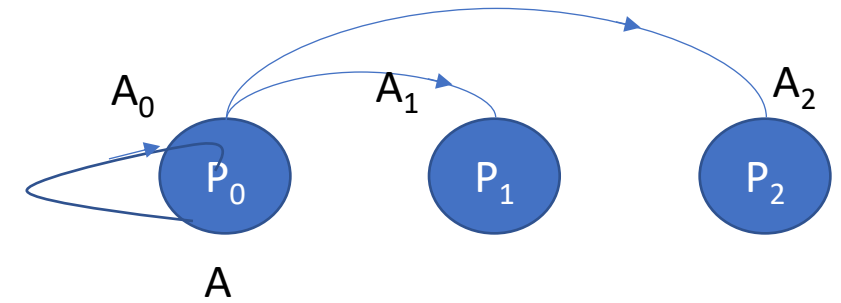
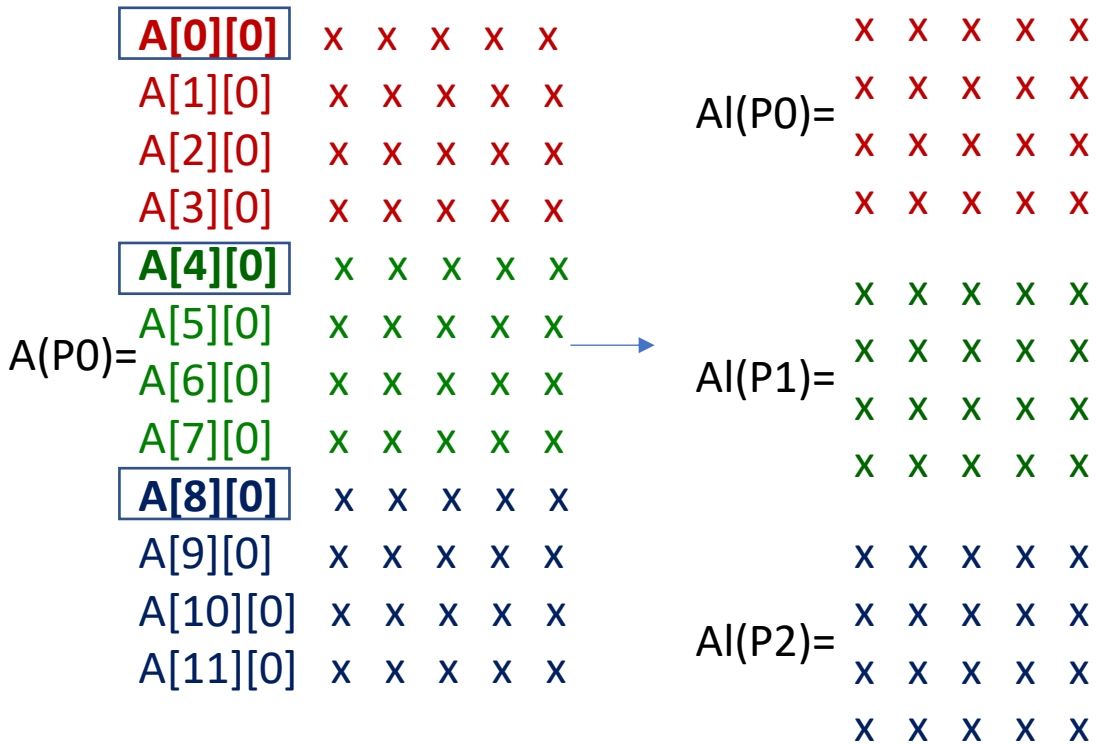
Ejercicio 2

$$A(P_0) = \begin{bmatrix} A_0 \\ A_1 \\ A_2 \end{bmatrix} \begin{matrix} k \\ k \\ k \end{matrix} \xrightarrow{\text{MPI_Scatter}} \begin{bmatrix} Al(P_0) = A_0 \\ Al(P_1) = A_1 \\ Al(P_2) = A_2 \end{bmatrix} \begin{matrix} k \\ k \\ k \end{matrix}$$

$$k = M/p$$

($k = \text{n}^\circ$ de filas de $A_i = \text{n}^\circ$ de filas de $Al(P_j)$, $p = \text{n}^\circ$ de procesos)

Ejemplo: $M=12, N=5, p=3 \rightarrow k=12/3=4$



Bloque	1ª fila del bloque
0	0
1	k
2	2k
3	3k
.....
i	i*k

P0:

for(i=1; i<p; i++)

 Enviar bloque de k*N elementos con origen A[i*k][0] a Pi

 Copiar sus 1ª k filas de A en Al (k*N elementos consecutivos)

Resto:

 Recibir en Al k*N elementos enviados por P0

Ejercicio 2

$K=M/p$

P0:

```
for(i=1; i<p; i++)
```

Enviar el bloque de $k*N$ elementos con origen $A[k*i][0]$ a P_i

Copiar sus 1ª k filas de A en A_i ($k*N$ elementos consecutivos)

Resto:

Recibir en A_i $k*N$ elementos de P_0



```
MPI_Init(&argc, &argv);
double A[M][N], A_i[M][N];
int id, p;
MPI_Status stat;
MPI_Comm_size(MPI_COMM_WORLD, &p);
MPI_Comm_rank(MPI_COMM_WORLD, &id);
int k=M/p; /* k= nº de filas de los bloques que se envían/reciben*/
if (id==0){
    for(i=1; i<p; i++)
        MPI_Send(&A[k*i][0], k*N, MPI_DOUBLE, i, 100, MPI_COMM_WORLD);
        MPI_Sendrecv(&A[0][0], k*N, MPI_DOUBLE, 0, 100, A_i, k*N, MPI_DOUBLE, 0, 100, MPI_COMM_WORLD, &stat);
}
else
    MPI_Recv(A_i, k*N, MPI_DOUBLE, 0, 100, MPI_COMM_WORLD, &stat);
MPI_Finalize();
```

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

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

int MPI_Sendrecv(void *sendbuf, int
                sendcount, MPI_Datatype sendtype, int
                dest, int sendtag, void *recvbuf, int
                recvcnt, MPI_Datatype recvtype, int
                source, int recvtag, MPI_Comm comm,
                MPI_Status *status)
```

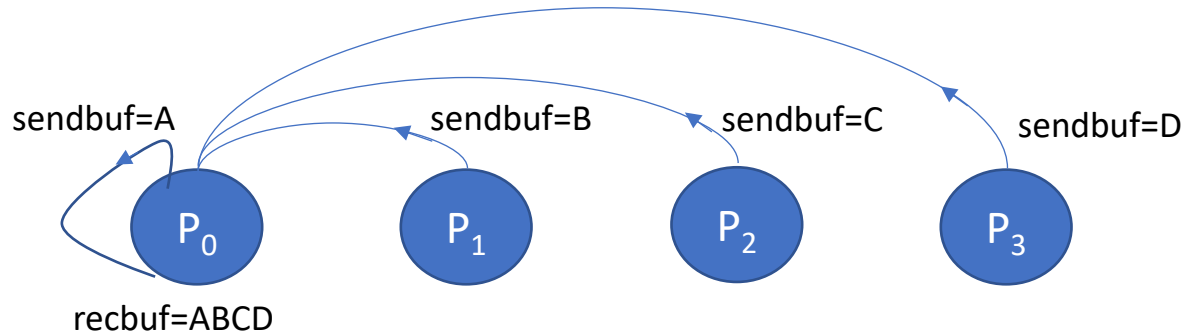
Tiempo de comunicaciones:

$$t_c = (p - 1) \left(t_s + \frac{MN}{p} t_w \right)$$

Recogida

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

Ejemplo:



sendcount=recvcount=|A|=|B|=|C|=|D|
sendtype=recvtype
root=0

	Antes
sendbuf(P ₀)	A
sendbuf(P ₁)	B
sendbuf(P ₂)	C
sendbuf(P ₃)	D

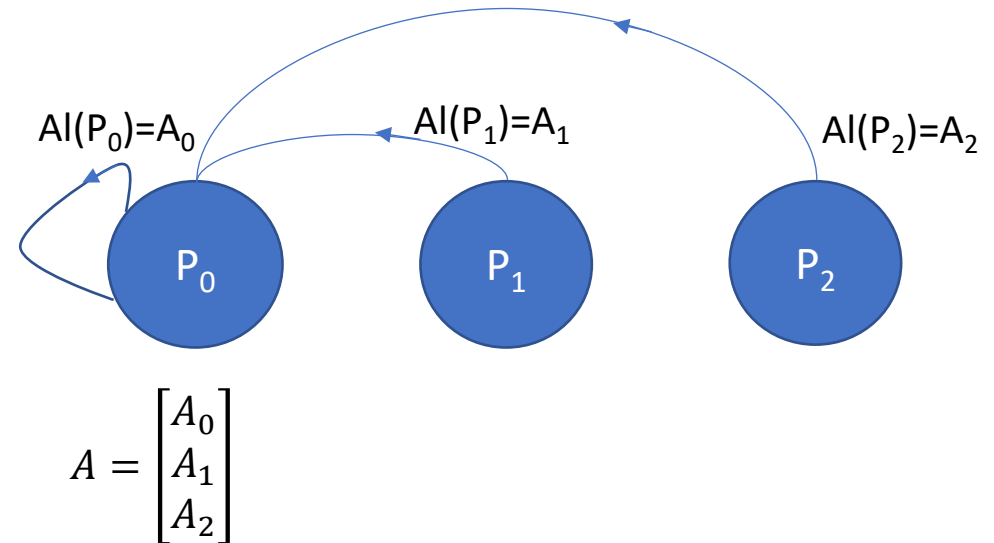
	Después
recvbuf(P ₀)	ABCD
recvbuf(P ₁)	
recvbuf(P ₂)	
recvbuf(P ₃)	

Ejercicio 3

Supongamos que en MPI se dispone de p procesos y una matriz \mathbf{A} de dimensión $M \times N$ de tipo double se encuentra repartida en matrices locales \mathbf{A}_i también de tamaños $k \times N$ ¿Qué función de comunicaciones colectivas deberían ejecutar todos los procesos para que P_0 recogiera en \mathbf{A} todas las matrices? ¿Cómo podrías sustituir esa comunicación colectiva mediante comunicaciones punto a punto, suponiendo que el número de procesos p divide a M ? Calcula el tiempo de comunicaciones.

$$\begin{bmatrix} A_l(P_0) = A_0 \\ A_l(P_1) = A_1 \\ A_l(P_2) = A_2 \end{bmatrix} \begin{matrix} N \\ k \\ k \\ k \end{matrix} \xrightarrow{\text{MPI_Gather}} A(P_0) = \begin{bmatrix} A_0 \\ A_1 \\ A_2 \end{bmatrix} \begin{matrix} k \\ k \\ k \end{matrix}$$
$$k = \frac{M}{p} \quad (p = n^\circ \text{ procesos})$$

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



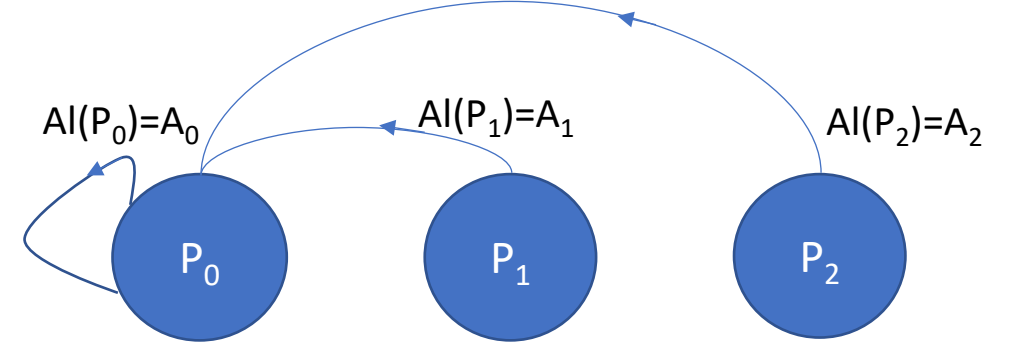
`MPI_Gather(A_l, k*N, MPI_DOUBLE, A, k*N, MPI_DOUBLE, 0, MPI_COMM_WORLD)`

Ejercicio 3

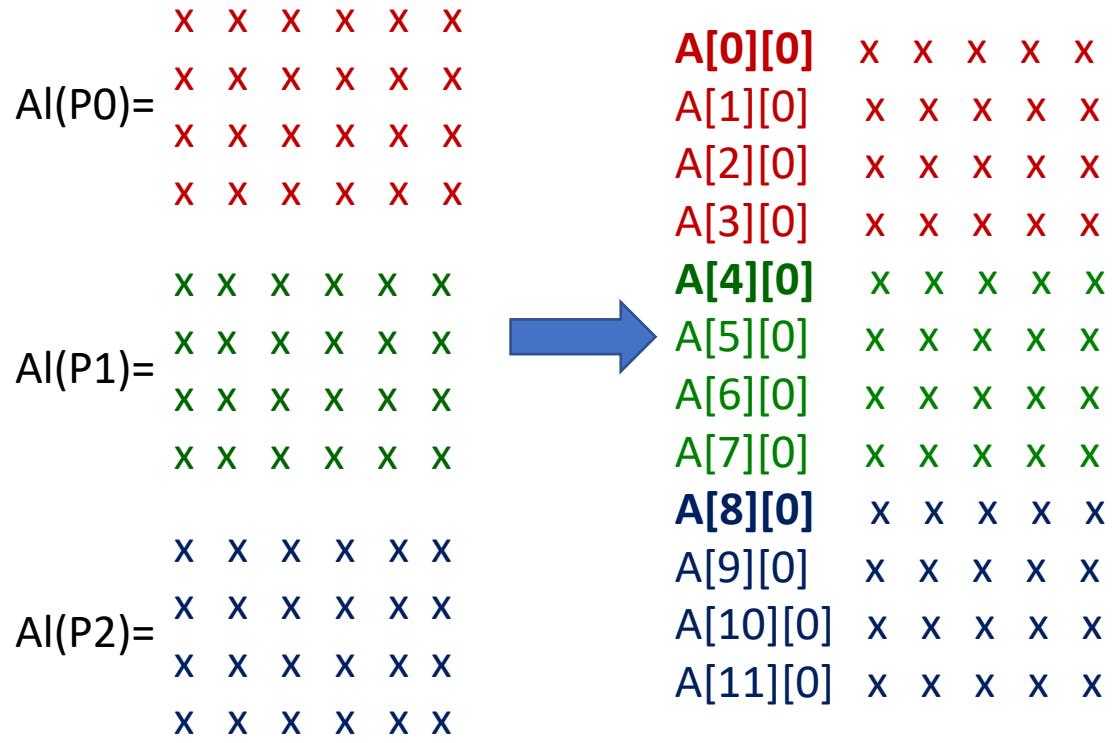
$$A = \begin{bmatrix} A_0 \\ A_1 \\ A_2 \end{bmatrix}$$

$$\begin{bmatrix} Al(P_0) = A_0 \\ Al(P_1) = A_1 \\ Al(P_2) = A_2 \end{bmatrix} k \xrightarrow{\text{Recogida MPI_Gather}} A(P_0) = \begin{bmatrix} A_0 \\ A_1 \\ A_2 \end{bmatrix} k$$

$$k = \frac{M}{p} \quad (p = n^{\circ} \text{ procesos})$$



Ejemplo: M=12, N=6, p=3 → k=12/3=4



Bloque	1ª fila del bloque
0	0
1	k
2	2k
3	3k
.....
i	i*k

P₀:

```
for (i=1; i<p; i++)
    Recibir en A[i*k][0] el bloque de k*N elementos de Pi
    Copiar Al a partir de la primera posición de A
```

Resto:

Enviar los k*N elementos de Al a P0

P0:
 for (i=1; i<p; i++)
 Recibir en A[k*i][0] bloque de kN elementos de Pi
 Copiar A1 a partir de la primera posición de A

Resto:
 Enviar los kN elementos de A1 a P0

```

MPI_Init(&argc, &argv);
double A[M][N], A1[M][N];
int id, p;
MPI_Status stat;
MPI_Comm_size(MPI_COMM_WORLD, &p);
MPI_Comm_rank(MPI_COMM_WORLD, &id);
int k=M/p; /* k= nº de filas de los bloques que se envían/reciben*/
if (id==0){
  for(i=1; i<p; i++)
    MPI_Recv(&A[i*k][0], k*N, MPI_DOUBLE, i, 100, MPI_COMM_WORLD, &stat);
  MPI_Sendrecv(A1, k*N, MPI_DOUBLE, 0, 100, A, k*N, MPI_DOUBLE, 0, 100, MPI_COMM_WORLD, &stat);
}
else
  MPI_Send(A1, k*N, MPI_DOUBLE, 0, 100, MPI_COMM_WORLD);

MPI_Finalize();

```

```

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

```

```

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

```

Ejercicio 3

```

int MPI_Sendrecv(void *sendbuf, int
                 sendcount, MPI_Datatype sendtype, int
                 dest, int sendtag, void *recvbuf, int
                 recvcount, MPI_Datatype recvtype, int
                 source, int recvtag, MPI_Comm comm,
                 MPI_Status *status)

```

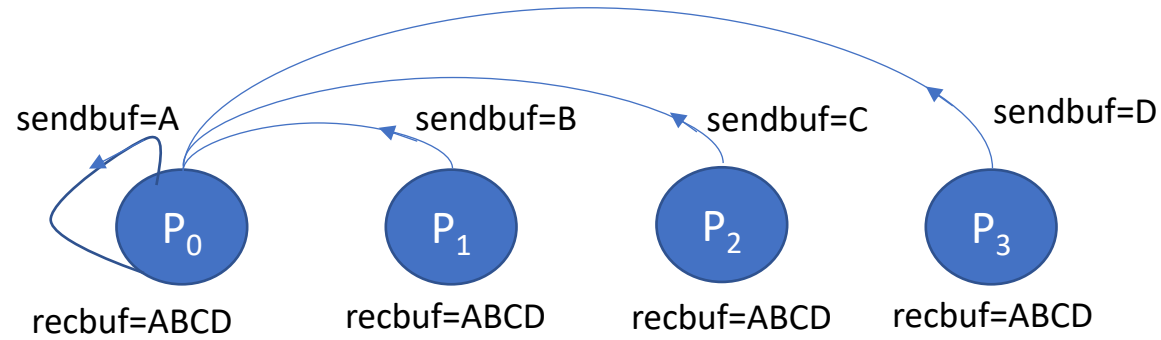
Tiempo de comunicaciones:

$$t_c = (p - 1) \left(t_s + \frac{MN}{p} t_w \right)$$

Multi-recogida

MPI_AllGather(void *sendbuf, int sendcount, MPI_Datatype sendtype, void *recvbuf, int recvcount, MPI_Datatype, recvtype, MPI_Comm comm)

Ejemplo:



sendcount=recvcount=|A|=|B|=|C|=|D|
sendtype=recvtype
root=0

	Antes		Después
sendbuf(P ₀)	A	recvbuf(P ₀)	ABCD
sendbuf(P ₁)	B	recvbuf(P ₁)	ABCD
sendbuf(P ₂)	C	recvbuf(P ₂)	ABCD
sendbuf(P ₃)	D	recvbuf(P ₃)	ABCD

Ejercicio 4

¿Qué cambiaras en el código que usa comunicaciones colectivas del ejercicio 8 para que la matriz A la tuviesen todos los procesos?

$$\begin{array}{c} N \\ \left[\begin{array}{l} Al(P_0) = A_0 \\ Al(P_1) = A_1 \\ Al(P_2) = A_2 \end{array} \right] k \\ M = kp \quad (p = n^{\circ} \text{ procesos}) \end{array} \xrightarrow{\text{MPI_Allgather}} A(P_0, P_1, P_2) = \begin{array}{c} N \\ \left[\begin{array}{l} A_0 \\ A_1 \\ A_2 \end{array} \right] k \end{array}$$

MPI_Gather(A1, k*N, MPI_DOUBLE, A, k*N, MPI_DOUBLE, 0, MPI_COMM_WORLD)



MPI_Allgather(A1, k*N, MPI_DOUBLE, A, k*N, MPI_DOUBLE, MPI_COMM_WORLD)

Todos a todos

MPI_Alltoall(void *sendbuf, int sendcount, MPI_Datatype sendtype,
void *recvbuf, int recvcount, MPI_Datatype, recvtype, MPI_Comm comm)

Ejemplo:

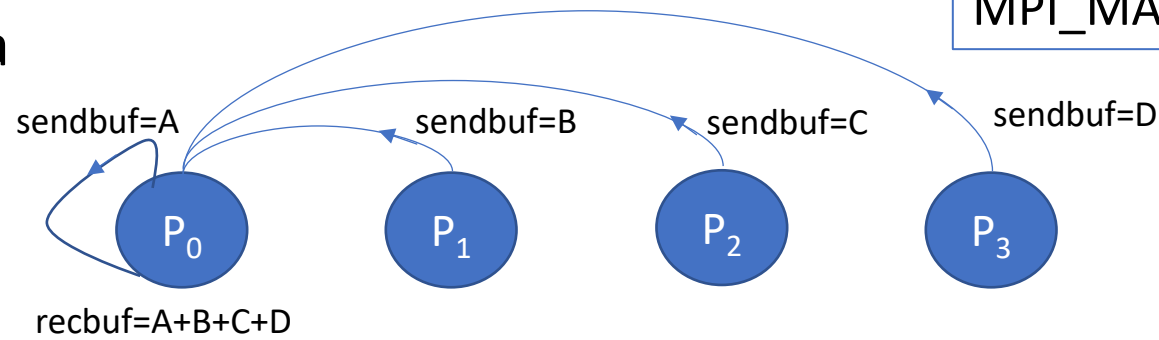
sendcount=recvcount= $|A_i|=|B_i|=|C_i|=|D_i|$
sendtype= recvtype
root=0

	Antes		Después
sendbuf(P ₀)	A ₀ A ₁ A ₂ A ₃	recvbuf(P ₀)	A ₀ B ₀ C ₀ D ₀
sendbuf(P ₁)	B ₀ B ₁ B ₂ B ₃	recvbuf(P ₁)	A ₁ B ₁ C ₁ D ₁
sendbuf(P ₂)	C ₀ C ₁ C ₂ C ₃	recvbuf(P ₂)	A ₂ B ₂ C ₂ D ₂
sendbuf(P ₃)	D ₀ D ₁ D ₂ D ₃	recvbuf(P ₃)	A ₃ B ₃ C ₃ D ₃

Reducción

`MPI_Reduce(void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm)`

Ejemplo: Suma



sendcount=|A|=|B|=|C|=|D|
op= MPI_SUM
root=0

	Antes
sendbuf(P ₀)	A
sendbuf(P ₁)	B
sendbuf(P ₂)	C
sendbuf(P ₃)	D

	Después
recvbuf(P ₀)	A+B+C+D
recvbuf(P ₁)	
recvbuf(P ₂)	
recvbuf(P ₃)	

Ejercicio 5 (Suma elementos de un vector)

Realiza una implementación paralela mediante MPI de la función **suma_vec**, la cual tiene como argumento de entrada un vector **v** de dimensión **N** y devuelve la suma de los elementos del vector **v**. Supondremos que inicialmente **v** está almacenada en P0 y que **N** es divisible entre el número de procesos **p**. Calcula el tiempo paralelo.

```
double suma_vec(double v[N]){
    int i; /*k=nº de filas de los v_i*/
    double sl=0;
    for (i=0; i<N; i++) /*calcular suma local*/
        sl+=v[i];
    return sl;
}
```

Solución:

$$v(P_0) = \begin{bmatrix} v_0 \\ v_1 \\ v_2 \end{bmatrix} k \xrightarrow{\text{MPI_Scatter}} \begin{bmatrix} vl(P_0) = v_0 \\ vl(P_1) = v_1 \\ vl(P_2) = v_2 \end{bmatrix} k \xrightarrow{\text{suma local}} \begin{bmatrix} sl(P_0) \\ sl(P_1) \\ sl(P_2) \end{bmatrix} \begin{matrix} 1 \\ 1 \\ 1 \end{matrix} \xrightarrow{\text{MPI_Reduce}} s(P_0) = sl(P_0) + sl(P_1) + sl(P_2)$$

$k = \frac{N}{p}$ (p =número de procesos)

$$v(P_0) = \begin{bmatrix} v_0 \\ v_1 \\ v_2 \end{bmatrix} \xrightarrow[k]{k \text{ MPI_Scatter}} \begin{bmatrix} vl(P_0) = v_0 \\ vl(P_1) = v_1 \\ vl(P_2) = v_2 \end{bmatrix} \xrightarrow[k]{k \text{ Suma local}} \begin{bmatrix} sl(P_0) \\ sl(P_1) \\ sl(P_2) \end{bmatrix} \begin{matrix} 1 \\ 1 \\ 1 \end{matrix} \quad s(P_0) = sl(P_0) + sl(P_1) + sl(P_2)$$

$k = N/p$ (p =número de procesos)

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

```
int MPI_Reduce(void *sendbuf, void
*recvbuf, int count, MPI_Datatype
datatype, MPI_Op op, int root, MPI_Comm
comm)
```

```
double suma_vecp(double v[N]){
    int p;
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    int i, k; /*k=nº de filas de los v_i*/
    /*En algunas ocasiones por simplicidad declaramos: double vl[N];*/
    double *vl = (double*) malloc(sizeof(double)*k);
    double sl=0, s; /* sl=suma local, s= suma total*/
    k=N/p;
    MPI_Scatter(v, k, MPI_DOUBLE, vl, k, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    for (i=0; i<k; i++) /*calcular suma local*/
        sl+=vl[i];
    MPI_Reduce(&sl, &s, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
    return s;
}
```

$$t_{\text{Scatter}} = (p - 1) \left(t_s + \frac{N}{p} t_w \right)$$

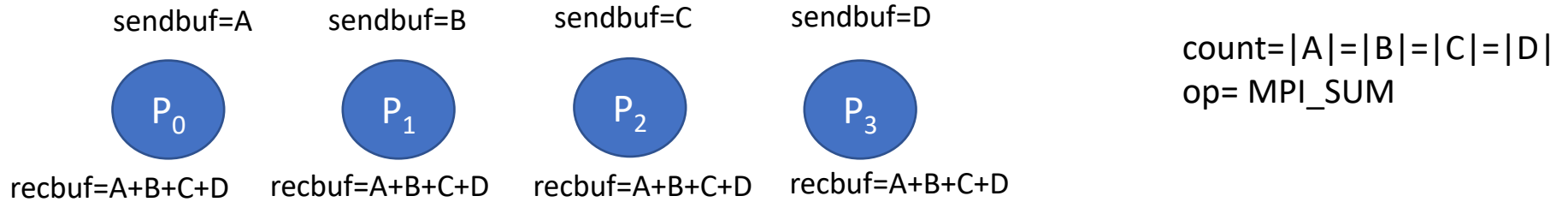
$$t_{\text{Reduce}} = (p - 1)(t_s + t_w)$$

$$t_c(N, p) = (p - 1) \left(2t_s + \left(1 + \frac{N}{p} \right) t_w \right)$$

Multi-Reducción

`MPI_Allreduce(void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)`

Ejemplo: Suma



	Antes	Después
sendbuf(P_0)	A	recvbuf(P_0) A+B+C+D
sendbuf(P_1)	B	recvbuf(P_1) A+B+C+D
sendbuf(P_2)	C	recvbuf(P_2) A+B+C+D
sendbuf(P_3)	D	recvbuf(P_3) A+B+C+D

Ejercicio 6 (Suma elementos de un vector)

¿Qué modificarías en la implementación del ejercicio 9 para que todos recibiesen la suma?

$$v(P_0) = \begin{bmatrix} v_0 \\ v_1 \\ v_2 \end{bmatrix} k \xrightarrow{\text{MPI_Scatter}} \begin{bmatrix} vl(P_0) \\ vl(P_1) \\ vl(P_2) \end{bmatrix} k \xrightarrow{\text{Suma local}} \begin{bmatrix} sl(P_0) \\ sl(P_1) \\ sl(P_2) \end{bmatrix} k \xrightarrow{\text{MPI_Allreduce}} s(P_0, P_1, P_2) = sl(P_0) + sl(P_1) + sl(P_2)$$

$k = N/p$ (p =número de procesos)

```
MPI_Reduce(&sl, &s, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
```

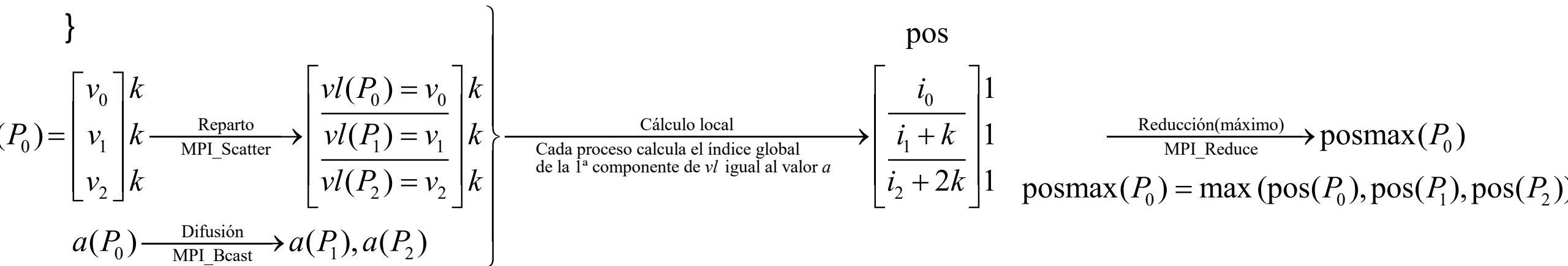
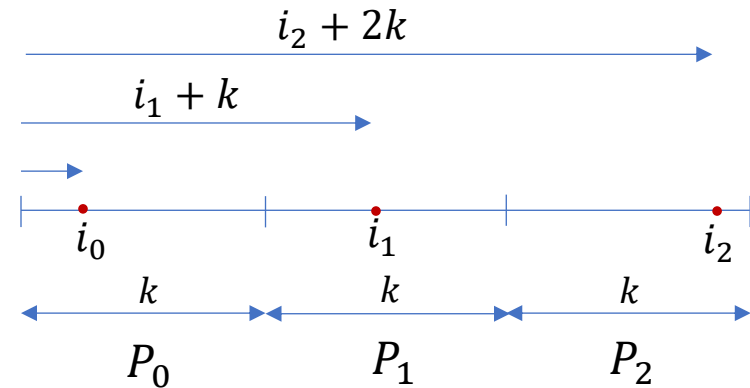


```
MPI_Allreduce(&sl, &s, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
```

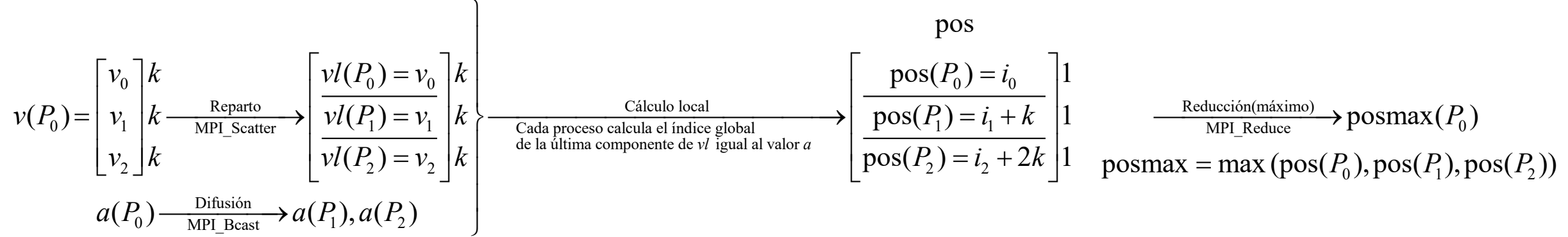
Ejercicio 7 (Búsqueda lineal)

La siguiente función devuelve la última posición del vector v coincidente con a o -1 si no la hubiese. Realiza una implementación MPI, suponiendo que inicialmente v y a están inicialmente almacenados en la memoria local de P_0 . Calcular el coste de las comunicaciones.

```
int pos(int v[N], int a){
    int pos=-1, i;
    for(i=0; i<N; i++){
        if (v[i]==a){
            pos=i;
        }
    }
    return pos;
}
```



$= N / p$ ($p =$ número de procesos)



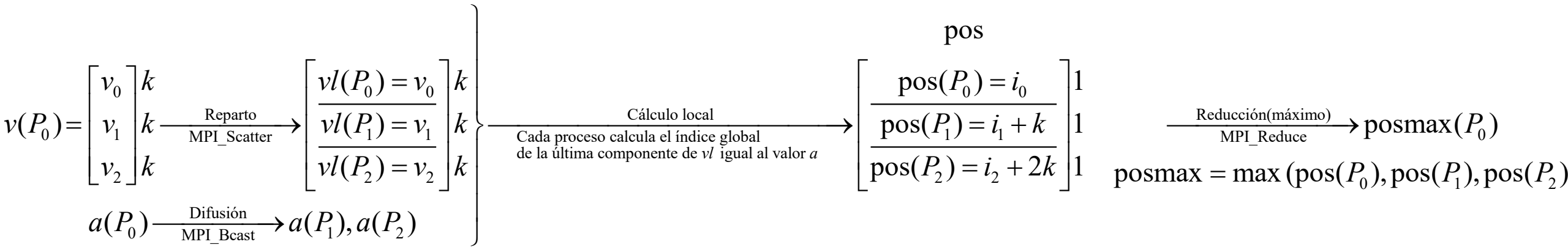
$k = N / p$ ($p =$ número de procesos)

```
int posp(int v[N], int a){
    int p;
    MPI_Comm_size(MPI_COMM_WORLD,&p);
    MPI_Comm_rank(MPI_COMM_WORLD,&id);
    int k=N/p; /* Número de componentes de vl que le corresponden a cada proceso*/
    int pos, posmax, i;
    int *vl = (int*) malloc(sizeof(int)*k);
    MPI_Scatter(v, k, MPI_INT, vl, k, MPI_INT, 0, MPI_COMM_WORLD);
    MPI_Bcast(&a, 1, MPI_INT, 0, MPI_COMM_WORLD);
    int pos=-1;
    for(i=0; i<k; i++)
        if (vl[i]==a){
            pos=i+id*k; /* posición global=posición local+k*id */
        }
    MPI_Reduce(&pos, &posmax, 1, MPI_INT, MPI_MAX, 0, MPI_COMM_WORLD);
    return posmax;
}
```

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

int MPI_Bcast(void *buf, int count,
MPI_Datatype datatype, int root,
MPI_Comm comm)

int MPI_Reduce(void *sendbuf, void
*recvbuf, int count, MPI_Datatype
datatype, MPI_Op op, int root, MPI_Comm
comm)
```



$k = N / p$ ($p =$ número de procesos)

Costes de comunicaciones:

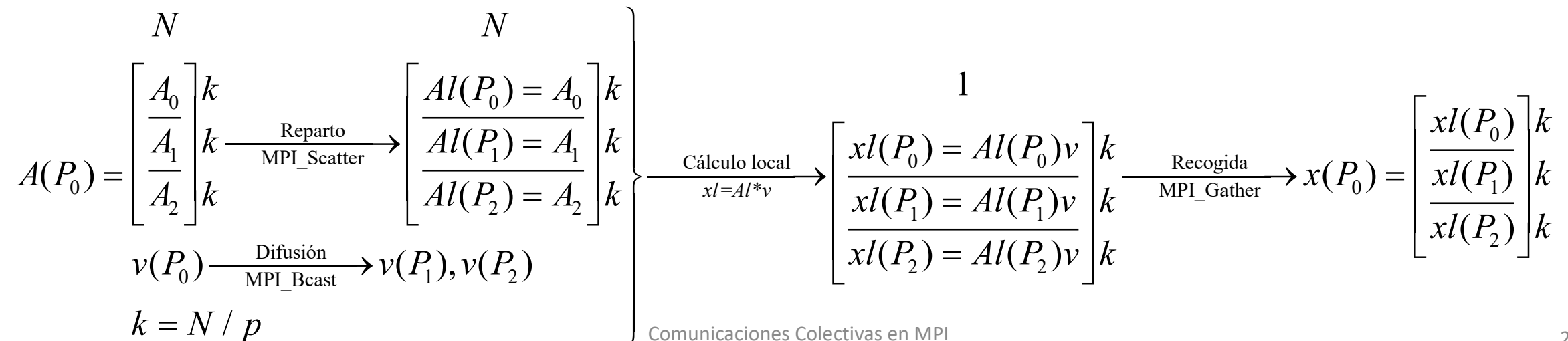
$$\left. \begin{aligned}
 &\text{MPI_Scatter}(\text{vector de dimensión } N): t_{c_s} = (p-1) \left(t_s + \frac{N}{p} t_w \right) \\
 &\text{MPI_Bcast}(\text{dato simple}): t_{c_g} = (p-1) (t_s + t_w) \\
 &\text{MPI_Reduce}(\text{dato simple}): t_{c_r} = (p-1) (t_s + t_w)
 \end{aligned} \right\} t_c = (p-1) \left(3t_s + \left(\frac{N}{p} + 2 \right) t_w \right)$$

Ejercicio 8

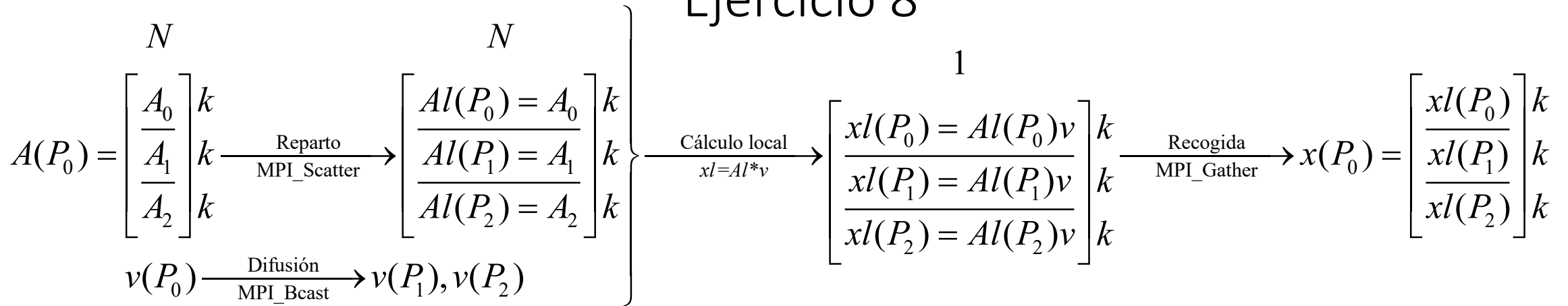
La siguiente función calcula el producto de una matriz cuadrada por un vector, ambos de dimensión N:

```
void prod_mat_vec(double A[N][N], double v[N], double x[N]){
    int i, j;
    for (i=0;i<N;i++) {
        x[i]=0;
        for (j=0; j<N; j++)
            x[i] += A[i][j]*v[j];
    }
}
```

Paraleliza la función anterior mediante MPI, teniendo en cuenta que el proceso P0 contiene a la matriz **A** y al vector **v**, y que el vector **x** debe quedar almacenado en P0. Calcula el coste paralelo.



Ejercicio 8



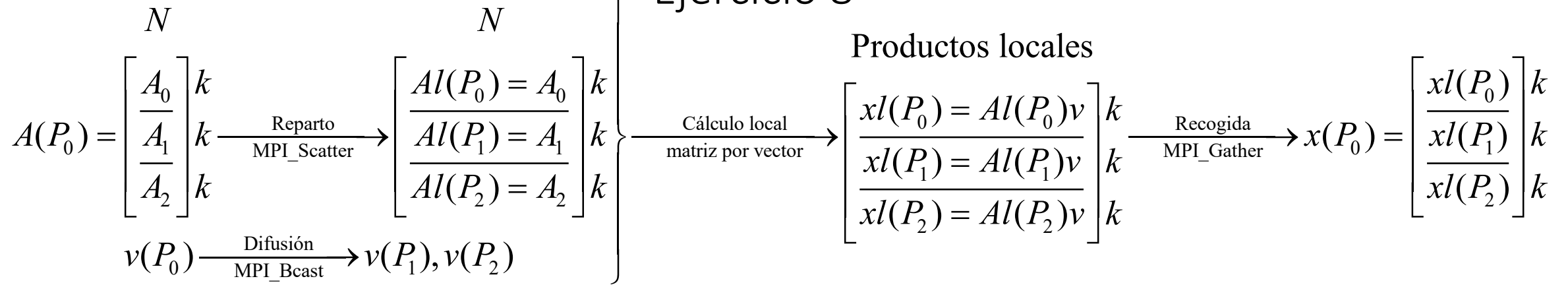
```
void prod_mat_vec(double A[N][N], double v[N], double x[N]){
int i, j, k, p;
double Al[N][N], xl[N];
MPI_Comm_size(MPI_COMM_WORLD, &p);
k=N/p;
MPI_Scatter(A, k*N, MPI_DOUBLE, Al, k*N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
MPI_Bcast(v, N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
for (i=0; i<k; i++) {
    xl[i]=0;
    for (j=0; j<N; j++)
        xl[i] += Al[i][j]*v[j];
}
MPI_Gather(xl, k, MPI_DOUBLE, x, k, MPI_DOUBLE, 0, MPI_COMM_WORLD);
```

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

```
int MPI_Bcast(void *buf, int count,
MPI_Datatype datatype, int root,
MPI_Comm comm)
```

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

Ejercicio 8



```

for (i=0;i<k;i++) {
    xl[i]=0;
    for (j=0;j<N;j++)
        xl[i] += Al[i][j]*v[j];
}
    
```

Tiempo aritmético:

$$t_a = \sum_{i=0}^{k-1} \sum_{j=0}^N 2 = 2kN = \frac{2N^2}{p} \text{ flops}$$

Tiempo de comunicaciones:

$$\text{MPI_Scatter} : t_{c_s} = (p-1) \left(t_s + \frac{N^2}{p} t_w \right)$$

$$\text{MPI_Bcast} : t_{c_b} = (p-1) (t_s + N t_w)$$

$$\text{MPI_Gather} : t_{c_g} = (p-1) \left(t_s + \frac{N}{p} t_w \right)$$

$$t_c = (p-1) \left(3t_s + \left(\frac{N^2}{p} + N + \frac{N}{p} \right) t_w \right)$$

Comunicaciones Colectivas en MPI

Tiempo paralelo:

$$t_p = t_a + t_c$$

Ejercicio 9

Realiza una implementación paralela mediante MPI de la función **suma_mat**, la cual tiene como argumento de entrada una matriz A de dimensión MxN y devuelve la suma de los elementos de la matriz A. Supondremos que inicialmente A está almacenada en P0 y que **M** es divisible entre el número de procesos **p**. Calcula el tiempo paralelo.

Solución:

```
double suma_mat (double A[M][N]){
```

```
.....
```

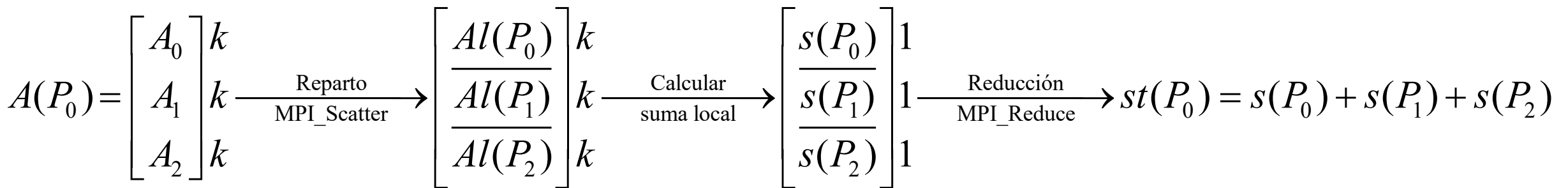
```
return 0;
```

```
}
```

N

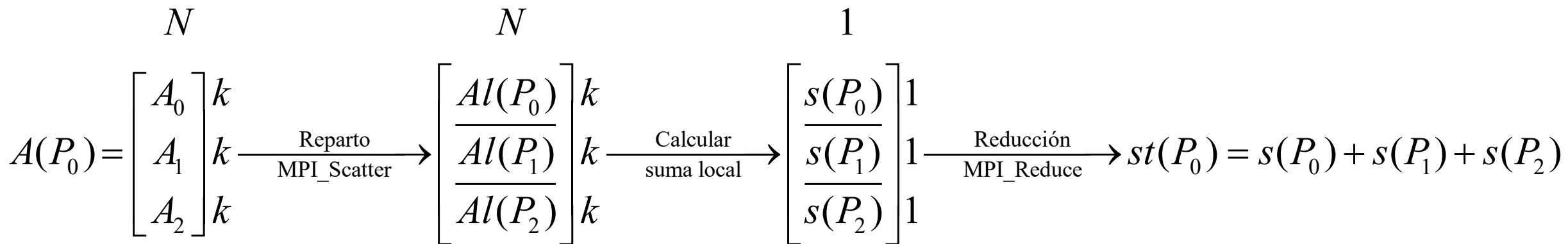
N

1



$$k = M / p \quad (p = \text{número de procesos})$$

Ejercicio 9



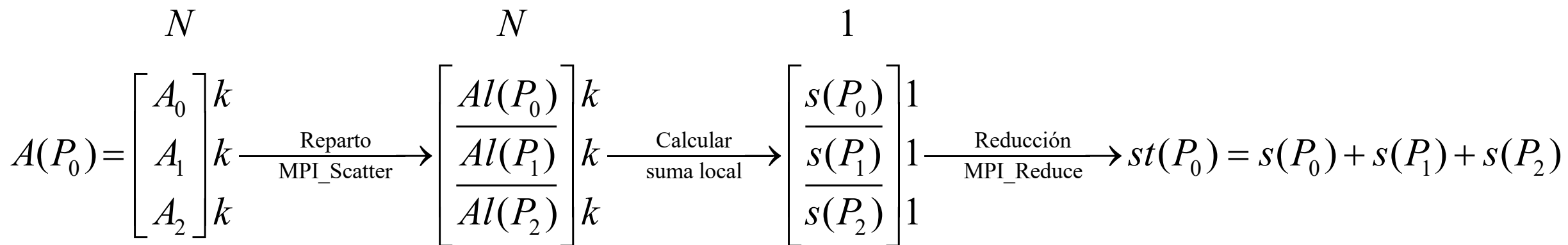
$k = M / p$ (p =número de procesos)

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

```
int MPI_Reduce(void *sendbuf, void
*recvbuf, int count, MPI_Datatype
datatype, MPI_Op op, int root, MPI_Comm
comm)
```

```
double suma_mat p(double A[M][N]){
    int p;
    MPI_Comm_size(MPI_COMM_WORLD, &p);
    int i, j, k= M/p; /*k=nº de filas de Ai = nº de filas de Al(Pi)*/
    double Al[M][N], s=0, st; /* s=suma local, st= suma total*/
    MPI_Scatter(A, k*N, MPI_DOUBLE, Al, k*N, MPI_DOUBLE, 0, MPI_COMM_WORLD);
    for (i=0; i<k; i++) /*calcular suma local*/
        for (j=0; j<N; j++)
            s+=Al[i][j];
    MPI_Reduce(&s, &st, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
    return st;
}
```

Ejercicio 9



$$k = M / p \quad (p = \text{número de procesos})$$

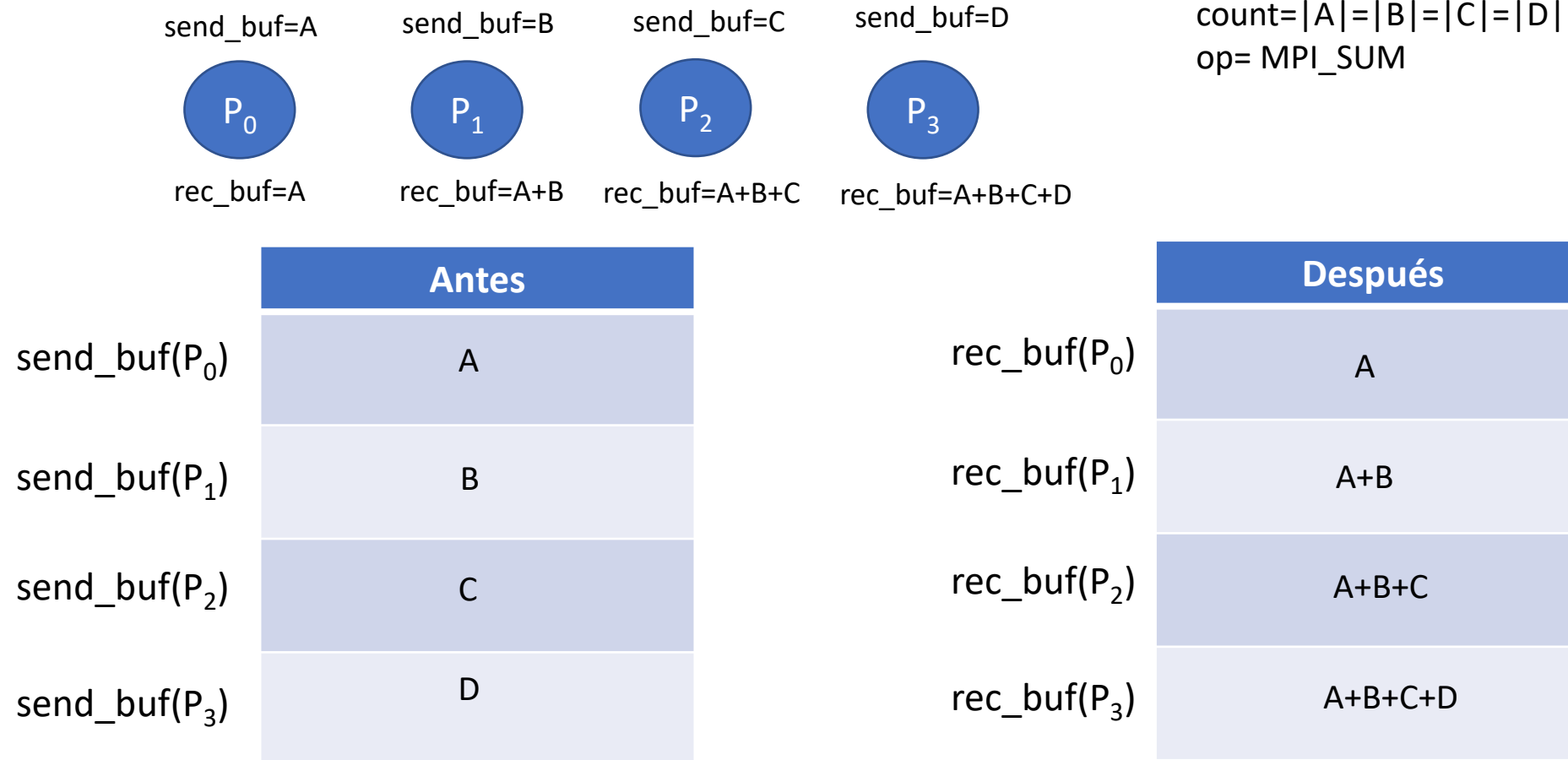
Costes de comunicaciones:

$$\begin{aligned}
 t_{\text{Scatter}} &= (p - 1) \left(t_s + \frac{MN}{p} t_w \right) \\
 t_{\text{Reduce}} &= (p - 1) (t_s + t_w) \\
 t_c(N, p) &= (p - 1) \left(2t_s + \left(1 + \frac{MN}{p} \right) t_w \right)
 \end{aligned}$$

Prefijación

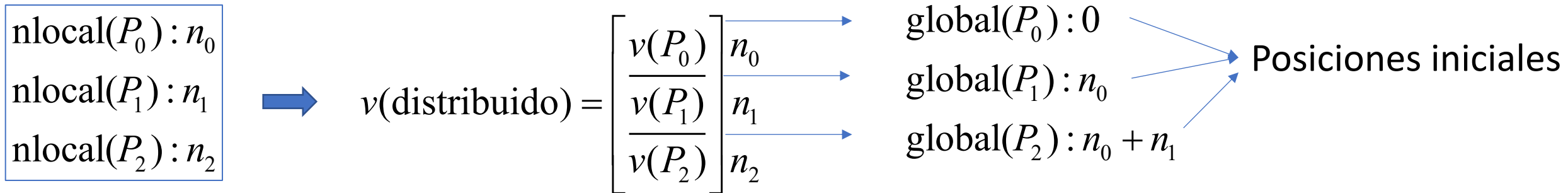
`MPI_Scan(void *sendbuf, void *recvbuf, int count, MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)`

Ejemplo: Suma



Ejercicio 10

Dado un vector v de longitud N , distribuido entre los procesos, donde cada proceso tiene n_{local} elementos consecutivos del vector, se quiere obtener la posición inicial del subvector local. Suponed que el número de elementos n_{local} se puede conocer al realizar la siguiente llamada: `calcula_nlocal(N,&nlocal)`.



Ejercicio 10

```
int global, nlocal, N;
```

```
calcula_nlocal(N,&nlocal);
```

```
MPI_Scan(&nlocal,&global,1,MPI_INT,MPI_SUM,comm);
```

```
global -= nlocal;
```

$nlocal(P_0) = n_0$

$nlocal(P_1) = n_1$

$nlocal(P_2) = n_2$

$global(P_0) = n_0$

$global(P_1) : n_0 + n_1$

$global(P_2) : n_0 + n_1 + n_2$

$global(P_0) = n_0 - n_0 = 0$

$global(P_1) = n_0 + n_1 - n_1 = n_0$

$global(P_2) : n_0 + n_1 + n_2 - n_2 = n_0 + n_1$

MPI_Scan

	Estado Inicial	Estado Final
P ₀	A	A
P ₁	B	A+B
P ₂	C	A+B+C
P ₃	D	A+B+C+D
P ₄	E	A+B+C+D+E
P ₅	F	A+B+C+D+E+F

Posiciones globales:

