

T3: Principios de Programación Paralela

T3.2: Modelo de Paso de Mensajes

Departamento de Electrónica y Sistemas

Primavera 2016



- 1 Conceptos Básicos
- 2 Operaciones Punto a Punto
- 3 Operaciones Colectivas

Conceptos Básicos

Modelo de Paso de Mensajes

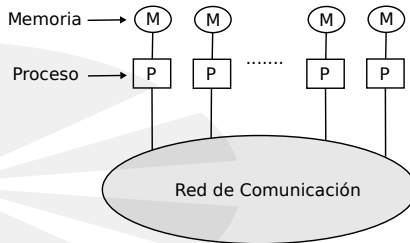
- Paradigma muy extendido en programación paralela
- MPI (Message Passing Interface) es la solución más popular (desde MPI1, 1992)
- Mínimos requerimientos al HW para su implementación
- Soporta un gran número de entornos paralelos, especialmente de memoria distribuida
- En este modelo uno o más procesos se comunican llamando a rutinas de una biblioteca para recibir y enviar mensajes entre procesos
- Control del paralelismo por el programador, que ha de evitar dependencias de datos, interbloqueos y *race conditions*
- Llamadas a MPI (u otra librería) desde programas C o Fortran
- Implementaciones de MPI: MPICH2, OpenMPI, Intel MPI, ...

Conceptos Básicos

Modelo de ejecución de un programa en paso de mensajes

- Programa paralelo compuesto de múltiples procesos/tareas que utilizan su propia memoria local durante la computación
- Generalmente un proceso/tarea por elemento de procesamiento (e.g., CPU core)
- Comunicación entre procesos mediante envío y recepción de mensajes *two-sided*, un envío se corresponde con una recepción

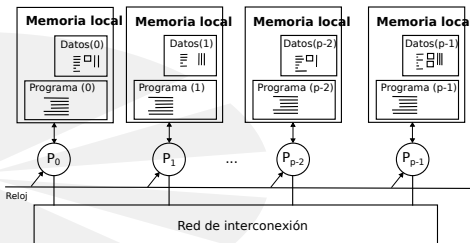
Arquitectura del modelo de paso de mensajes:



Conceptos Básicos

Estructura de un programa en paso de mensajes

- **MPMD** (Multiple Program Multiple Data): cada proceso/tarea tiene su propio programa con comunicaciones asíncronas entre ellos (máxima flexibilidad y complejidad)
- **SPMD** (Single Program Multiple Data): todos los procesos/tareas comparten un mismo programa/binario aunque en su lógica interna las tareas se pueden ejecutar de forma condicional dependiendo del proceso. Se suele hacer uso de comunicaciones síncronas con lo que suele resultar más sencillo programar pero con menor escalabilidad.



Conceptos Básicos

Características de un programa MPI C

- Incluye la librería de MPI (mpi.h)
- Las funciones MPI tienen la forma MPI_Nombre(parámetros)
- Devuelven un valor de éxito (MPI_SUCCESS) o error (MPI_ERR_{*}). Consultar man.
- Los procesos son independientes hasta que se inicializa MPI (MPI_Init), pudiendo colaborar intercambiando datos, sincronizándose tras ese punto
- Clave que los procesos conozcan el número de procesos (numprocs, obtenido con MPI_Comm_size) que se han puesto en marcha así como su identificador (entre 0 y numprocs - 1, obtenible con MPI_Comm_rank)
- MPI_Finalize se llama cuando ya no es necesario que los procesos colaboren entre sí. Libera todos los recursos reservados por MPI
- MPI_COMM_WORLD: comunicador global, incluye a todos los procesos

Conceptos Básicos

Hello World MPI C

```
#include <stdio.h>
#include <mpi.h>

int main(int argc, char *argv[]) {
    int numprocs, rank, namelen;
    char processor_name[MPI_MAX_PROCESSOR_NAME];

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Get_processor_name(processor_name, &namelen);

    printf("Process %d on %s out of %d\n", rank, processor_name, numprocs);

    MPI_Finalize();
}
```

Conceptos Básicos

Compilación

```
mpicc mpi-hello.c -o mpi-hello
```

Ejecución

```
mpirun -np 4 ./mpi-hello
```

Output

```
user@server:~$ Process 0 on localhost out of 4  
user@server:~$ Process 1 on localhost out of 4  
user@server:~$ Process 3 on localhost out of 4  
user@server:~$ Process 2 on localhost out of 4
```


Operaciones Punto a Punto

Punto a punto MPI

- Bloqueantes: MPI_Send y MPI_Recv
- Variantes (e.g., no bloqueantes, buffered, síncronas)

MPI_Send

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

- Envía un mensaje al proceso *dest* en el comunicador *comm*
- El mensaje está almacenado en *buff* y consta de al menos *count* items del tipo *datatype*
- El mensaje está etiquetado con un *tag*
- La llamada a MPI_Send finaliza cuando *buff* puede ser reusado (generalmente cuando el mensaje ha sido recibido en el destino)

Operaciones Punto a Punto

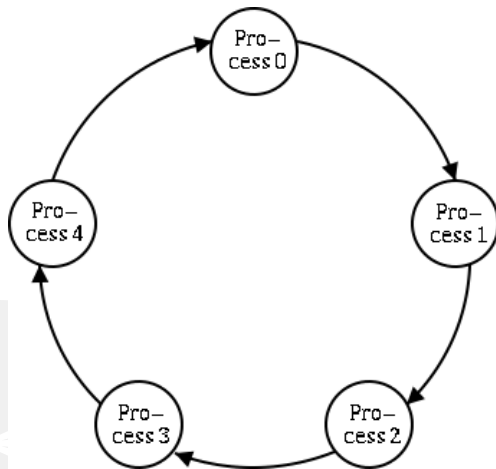
MPI_Recv

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

- Recibe un mensaje del proceso *source* del comunicador *comm* con la etiqueta *tag*
 - También se puede recibir de cualquier proceso del comunicador con *MPI_ANYSOURCE*
 - También se puede recibir mensajes con cualquier etiqueta con *MPI_ANYTAG*
- En los dos casos anteriores se recupera el *source* o *tag* recibidos accediendo a *status.MPI_SOURCE* y/o a *status.MPI_TAG*
- El mensaje se recibe en *buff* y consta de un máximo de *count* items del tipo *datatype*
- La llamada a *MPI_Recv* finaliza cuando se ha recibido el mensaje en *buff*

Operaciones Punto a Punto

Ejemplo MPI: **Ring**



Operaciones Punto a Punto

Ring MPI C (parte I)

```
MPI_Init(&argc, &argv);
MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
MPI_Comm_rank(MPI_COMM_WORLD, &my_id);

printf("my_id %d numprocs %d\n", my_id, numprocs);

if (my_id == 0) {
    passed_num = 1;

    printf("Root: before sending num=%d to dest=%d\n", passed_num, 1);
    MPI_Send(&passed_num, 1, MPI_INT, 1, 0, MPI_COMM_WORLD);

    printf("Root: before receiving from source=%d\n", numprocs-1);
    MPI_Recv(&passed_num, 1, MPI_INT, numprocs-1, 0, MPI_COMM_WORLD, &status);

    printf("Root: after receiving passed_num=%d from source=%d\n",
           passed_num, numprocs-1);
}
```

Operaciones Punto a Punto

Ring MPI C (parte II)

```
else {  
    printf("Process %d: before receiving from source=%d\n",  
           my_id, my_id-1);  
    MPI_Recv(&passed_num,1,MPI_INT,my_id-1,0,MPI_COMM_WORLD,&status);  
  
    printf("Process %d: after receiving passed_num=%d from source=%d\n",  
           my_id, passed_num, my_id-1);  
  
    passed_num++;  
  
    printf("Process %d: before sending passed_num=%d to dest=%d\n",  
           my_id, passed_num, (my_id+1)%numprocs);  
    MPI_Send(&passed_num,1,MPI_INT,(my_id+1)%numprocs,0,MPI_COMM_WORLD);  
  
    printf("Process %d: after send to dest=%d\n",  
           my_id, (my_id+1)%numprocs);  
}  
MPI_Finalize();
```

Operaciones Punto a Punto

Ring MPI C (Salida)

```
user@localhost:~/ $ mpirun -n 4 ./a.out
my_id 0 numprocs 4
Root: before sending num=1 to dest=1
Root: before receiving from source=3
my_id 1 numprocs 4
Process 1: before receiving from source=0
Process 1: after receiving passed_num=1 from source=0
Process 1: before sending passed_num=2 to dest=2
Process 1: after send to dest=2
my_id 3 numprocs 4
Process 3: before receiving from source=2
my_id 2 numprocs 4
Process 2: before receiving from source=1
Process 2: after receiving passed_num=2 from source=1
Process 2: before sending passed_num=3 to dest=3
Process 2: after send to dest=3
Process 3: after receiving passed_num=3 from source=2
Process 3: before sending passed_num=4 to dest=0
Process 3: after send to dest=0
Root: after receiving passed_num=4 from source=3
```

Operaciones Colectivas

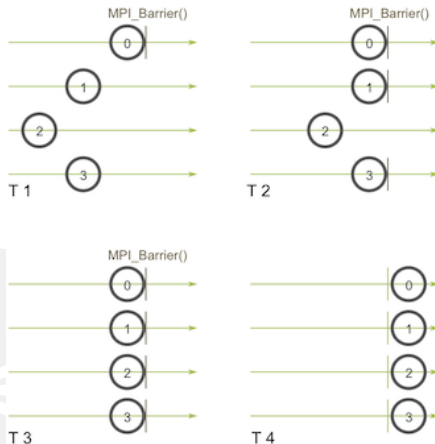
Colectivas MPI

- Operaciones típicas en las que intervienen todos los procesos de un comunicador
 - Barrier o barrera
 - Broadcast o difusión
 - Scatter o reparto
 - Gather o recolección
 - Reduce o reducción
 - Otras (e.g., Scan)
 - Combinaciones de las previas (e.g., Allreduce o Allgather)
- Uso recomendable al incrementar productividad:
 - Mayor rendimiento (optimizadas para cada librería, sistema, etc...)
 - Reducción de errores
 - Codificación a más alto nivel

Operaciones Colectivas

MPI_Barrier: Establece una barrera que bloquea el programa hasta que todos los procesos han alcanzado esta rutina.

```
int MPI_Barrier(MPI_Comm comm);
```

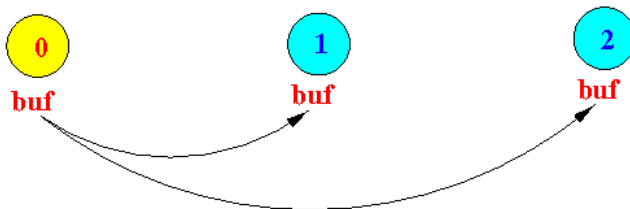


Operaciones Colectivas

MPI_Bcast: comunicación uno a todos de *count* datos del tipo *datatype* desde el proceso raíz (*root*) al resto de procesos del comunicador *comm*.

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

MPI_Bcast(buf, 10, MPI_INT, 0, MPI_COMM_WORLD)

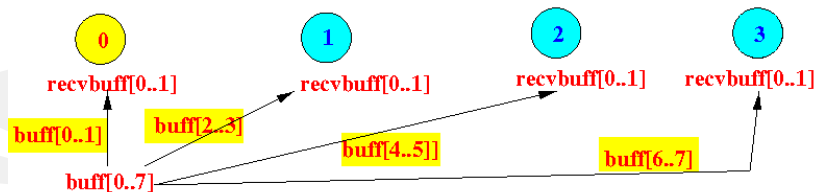


Operaciones Colectivas

MPI_Scatter: distribuye *sendcnt* elementos de *buff* de tipo *sendtype* desde el proceso *root* a todos los procesos del comunicador *comm*.

```
MPI_Scatter(void *buff, int sendcnt, MPI_Datatype sendtype,
            void *recvbuff, int recvcnt, MPI_Datatype recvtype, int root,
            MPI_Comm comm);
```

MPI_Scatter (buff, 2, MPI_INT, recvbuff, 2, MPI_INT, 0, WORLD);

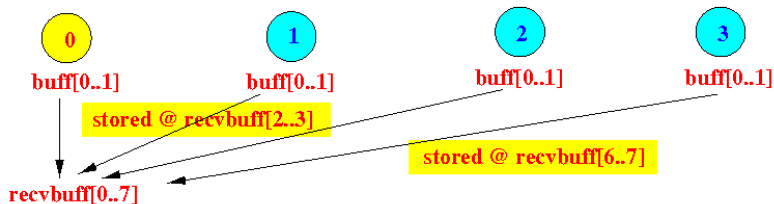


Operaciones Colectivas

MPI_Gather: recibe en el proceso *root*, en *recvbuff*, *recvcnt* elementos de tipo *recvtype* desde todos los procesos del comunicador *comm*.

```
MPI_Gather(void *buff, int sendcnt, MPI_Datatype sendtype,
          void *recvbuff, int recvcnt, MPI_Datatype recvtype, int root,
          MPI_Comm comm);
```

MPI_Gather (buff, 2, MPI_INT, recvbuff, 2, MPI_INT, 0, WORLD);



Operaciones Colectivas

MPI_Reduce: realiza una reducción todos a uno, reduciendo los datos de *buff*, *count* elementos de tipo *datatype*, y guardando el resultado en *recvbuff* del proceso *root*. Operaciones *op* disponibles:

MPI_{MAX,MIN,SUM,PROD}, MPI_{LAND,LOR,LEXOR},
MPI_{BAND,BOR,BXOR} o MPI_{MAXLOC,MINLOC}

```
int MPI_Reduce(void *buff, void *recvbuff, int count,
               MPI_Datatype datatype, MPI_Op op, int root, MPI_Comm comm);
```

MPI_Reduce (buff, recvbuff, 1, MPI_INT, MPI_SUM, 0, WORLD);

