

Introducción

......



Lugar / Fecha: Trujillo, 11 y 18 de Mayo 2010



MINISTERIO DE CIENCIA E INNOVACIÓN







- 1 Introducción a MPI
- 2 Nociones básicas
- **3** Funciones de MPI
- 4 Casos prácticos

#### **Tutorial MPI – Primeros pasos**



GOBIERNO DE ESPAÑA MINISTERIO DE CIENCIA E INNOVACIÓN





FEDER

Fondo Europeo de Desarrollo Regional





#### ¿Qué es MPI?

 MPI es una interfaz de paso de mensaje que representa un esfuerzo prometedor de mejorar la disponibilidad de un software altamente eficiente y portable para satisfacer las necesidades actuales en la computación de alto rendimiento a través de la definición de un estándar de paso de mensajes universal.

William D. Gropp et al.

- Protocolo de comunicación entre máquinas
- Estándar para la comunicación entre los nodos que ejecutan un programa en un sistema de memoria distribuida
- Conjunto de librerías de funciones que pueden ser utilizadas en programas escritos en C, C++, Fortran y Ada









#### Ventajas de utilizar MPI

- Portabilidad: MPI ha sido implementado para casi toda arquitectura de memoria distribuida
- Velocidad: Cada implementación de la biblioteca ha sido optimizada para el hardware en la cual se ejecuta
  - Ejemplo: Clusters de Bull utilizan MPIBull2 optimizada para IB
- Estándar: 60 personas de 40 organizaciones diferentes principalmente de U.S.A. y Europa trabajaron en estandarizarlo (IBM, INTEL, NX/, Express, nCUBE's Vernex, p4, PARMACS, Zipcode, Chimp, PVM, Chameleon, PICL...)
- Amplia funcionalidad
- Gran variedad de implementaciones libres (MPICH, OpenMPI, LAM-MPI...)







## Implementaciones de MPI ("MPI flavors/flavours")

- LAM/MPI (Local Area Multicomputer).
- **OpenMPI**
- **MPICH:** MPICH2 (estándar MPI 2.1)
- Diferentes ramas en cada implementación:
  - Para GRID: MPICH-G2, MPICH-GM, MPICH-VMI2...
  - Para Infiniband: MVAPICH, MPIBull2...
- Otros lenguajes: Python (PyMPI), Java (mpiJava), ...
- Más info sobre el estándar MPI 2.0: http://www.mcs.anl.gov/research/projects/mpi/mpi-standard/mpi-report-2.0/mpi2report.htm







**CETA-Ciemat/ Mayo 2010** 

#### Nociones básicas

- La unidad básica en MPI son los procesos.
- Tienen espacios de memoria independientes.
- El Intercambio de información es mediante paso de mensajes.
- Introduce en concepto de comunicadores (grupo de procesos más contexto).
- A cada proceso se le asigna un identificador interno propio de MPI (rank).



#### Primer programa en MPI

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

```
int main(int argc, char *argv[]) {
     int myrank, size;
     MPI_Init(&argc, &argv);
     MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
     MPI_Comm_size(MPI_COMM_WORLD, &size);
     printf("Hola soy el proceso %d de %d\n", myrank, size);
     MPI_Finalize();
     exit(0);
```







#### ¿Qué son los comunicadores?

- MPI agrupa los procesos implicados en una ejecución paralela en comunicadores
- Un comunicador agrupa a procesos que pueden intercambiarse mensajes.
- El comunicador MPI\_COMM\_WORLD está creado por defecto y engloba a todos los procesos.
- Pueden definirse varios comunicadores.
  - MPI Comm create, MPI Comm free, MPI Comm group, MPI Group include, MPI Group exclude...





### Inicio y Finalización

- int MPI\_Init(int \*argc, char \*\*argv)
  - Primera llamada OBLIGATORIA de cada uno de los procesos MPI
  - Establece entorno
  - Un proceso solo puede hacer una llamada MPI\_INIT
- int MPI\_Finalize(void)
  - Termina la ejecución en MPI
  - Libera los recursos utilizados por MPI







#### Identificación de procesos

- MPI\_Comm\_rank (comm, &identificador);
  - Devuelve el identificador del proceso dentro del comunicador comm especificado (ej.: MPI COMM WORLD).
- MPI\_Comm\_size (comm, &num\_proc);
  - Devuelve en **num\_proc** el número de procesos del comunicador especificado.





## Primer programa en MPI (segundo vistazo)

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

```
int main(int argc, char *argv[]) {
     int myrank, size;
     MPI_Init(&argc, &argv);
     MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
     MPI_Comm_size(MPI_COMM_WORLD, &size);
     printf("Hola soy el proceso %d de %d\n", myrank, size);
     MPI Finalize();
```





**exit(0)**;



#### Tipos de datos

#### Todo tipo de datos:

**MPI CHAR** signed char

**MPI SHORT** signed short int

**MPI INT** signed int

**MPI LONG** signed long int

MPI UNSIGNED CHAR unsigned char

MPI UNSIGNED SHOT unsigned short int

MPI UNSIGNED unsigned int

MPI UNSIGNED LONG unsigned long int

**MPI FLOAT** float

MPI DOUBLE double

MPI LONG DOUBLE long double

MPI también permite definir tipos de datos propios (MPI TYPE VECTOR, MPI TYPE STRUCT...)



Fondo Europeo de Desarrollo Regional





**CETA-Ciemat/ Mayo 2010** 

#### **Funciones de MPI**

- MPI 1.2 tiene 129 funciones
- Las **funciones principales** de MPI son:
  - MPI\_Init
  - MPI\_Finalize
  - MPI\_Comm\_size
  - MPI\_Comm\_rank
  - MPI\_Send
  - MPI Recv



Fondo Europeo de Desarrollo Regional





#### Tipos de mensajes (1)

- Punto a punto: 1 emisor 1 receptor
  - Síncronos: El emisor se espera a que la comunicación se produzca (o que el buffer del receptor este disponible). BLOQUEANTE
    - Send, Recv, SendRecv, Bsend, Ssend, ...
  - Asíncronos: El emisor NO se espera a que la comunicación se produzca y se comprueba más tarde si se ha producido. NO BLOQUEANTES
    - Isend, Irecv, ...







#### Tipos de mensajes (2)

- **MPI\_Send**(buf, count, datatype, dest, tag, comm)
- **MPI** Recv(buf, count, datatype, source, tag, comm, status)
- **buf**: Dirección donde comienza el mensaje.
- count: número de elementos del mensaje.
- datatype: tipo del mensaje.
- dest/source: posición relativa del proceso fuente/destino dentro del comunicador.
  - MPI ANY SOURCE: permite recibir mensaje de cualquier fuente
- status: estado de la recepción

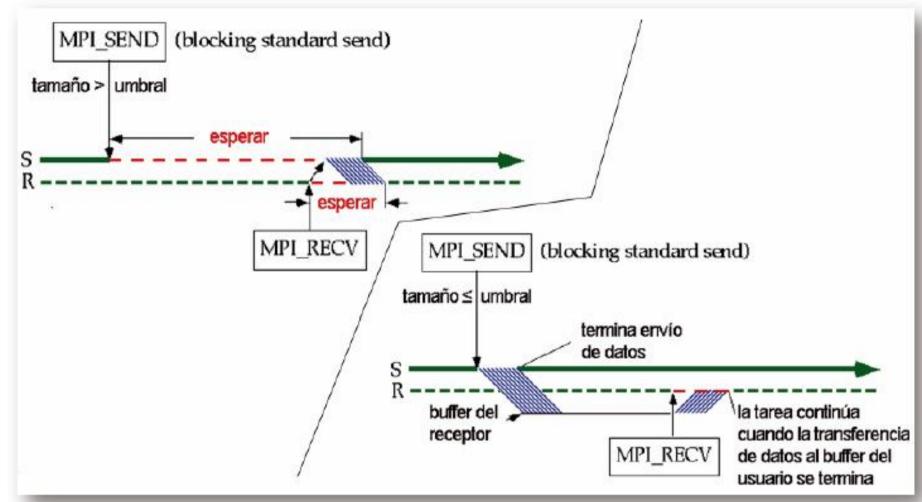






**CETA-Ciemat/ Mayo 2010** 

## Tipos de mensajes (3)





Fondo Europeo de Desarrollo Regional





CENTRO EXTREMEÑO DE TECNOLOGÍAS AVANZADAS

#### Tipos de mensajes (4)

- Variantes de MPI Send para casos especiales
- Sincrono: MPI\_Ssend(.....)
  - La operación se da por terminada sólo cuando el mensaje ha sido recibido en **destino**. (MPI Send termina cuando encuentra y receptor y envía su mensaje)
- Buffered: MPI Bsend(.....)
  - Cuando se hace un envío con buffer se guarda inmediatamente, en un buffer al efecto en el emisor, una copia del mensaje. La operación se da por completa en cuanto se ha efectuado esta copia. Si no hay espacio en el buffer, el envío fracasa.
- Ready: MPI\_Rsend(.....)
  - Sólo se puede hacer si antes el otro extremo está preparado para una recepción inmediata. No hay copias adicionales del mensaje (como en el caso del modo con buffer), y tampoco podemos confiar en bloquearnos hasta que el receptor esté preparado.

(para más info mirar referencias al final)

**CETA-Ciemat/ Mayo 2010** 







### Tipos de mensajes (5)

- Multipunto: 1 o N emisores 1 o N receptores
  - MPI\_Barrier() Sincronización global entre los procesadores del comunicador.
  - MPI\_Bcast()
  - MPI\_Gather()
  - MPI Scatter()
  - MPI Alltoall()
  - MPI\_Reduce()
  - MPI\_Reduce\_scatter( )
  - MPI\_Scan()...

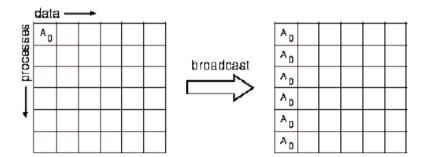


Fondo Europeo de Desarrollo Regional

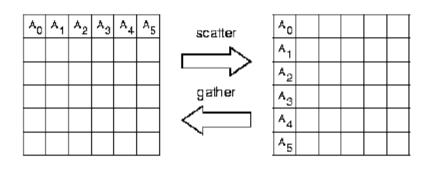




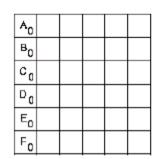
## Funciones de MPI



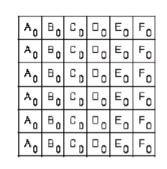
A <sub>0</sub>	A <sub>1</sub>	A <sub>2</sub>	$A_3$	A <sub>4</sub>	A <sub>5</sub>		A <sub>0</sub>	Во	Co	D a	Eo	FD
Ва	В <sub>1</sub>	В2	В3	В4	В <sub>5</sub>		A 1	В <sub>1</sub>	с <sub>1</sub>	D <sub>1</sub>	E <sub>1</sub>	F <sub>1</sub>
Co	c,	$c_2$	$c^3$	04	c 5		A <sub>2</sub>	В2	$c_2$	02	${\sf E_2}$	F <sub>2</sub>
D <sub>0</sub>	01	D <sub>2</sub>	Dg	04	D <sub>5</sub>		Ag	Вз	c <sub>3</sub>	Dg	E <sub>3</sub>	F3
E	E <sub>1</sub>	E2	Е3	E <sub>4</sub>	E <sub>5</sub>		A <sub>4</sub>	В4	c <sub>4</sub>	D <sub>4</sub>	E <sub>4</sub>	F <sub>4</sub>
Fa	F <sub>1</sub>	F <sub>2</sub>	F <sub>3</sub>	F <sub>4</sub>	F <sub>5</sub>		A <sub>5</sub>	В5	c <sub>5</sub>	D <sub>5</sub>	E <sub>5</sub>	F <sub>5</sub>

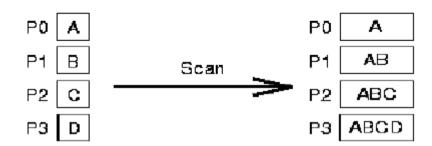


P0 <b>A</b>		P0	ABCD
P1 B	Reduce	P1	
P2 G		P2	
P3 D		P3	









Fondo Europeo de Desarrollo Regional

**FEDER** 





## Tipos de reducción

- **MPI MAX** maximum
- MPI MIN minimum
- MPI\_SUM sum
- MPI PROD product
- **MPI LAND** logical and
- **MPI BAND** bit-wise and
- MPI\_LOR logical or
- MPI BOR bit-wise or
- MPI LXOR logical xor
- **MPI BXOR** bit-wise xor
- MPI MAXLOC max value and location
- MPI\_MINLOC min value and location



Fondo Europeo de Desarrollo Regional







#### Caso práctico 1

- Crear un programa sencillo en C utilizando MPI, en el que un nodo (maestro) salude a otros nodos (esclavos).
- Condiciones:
  - El saludo tiene que realizarlo antes el maestro que los esclavos.
  - El maestro tiene que enviar un mensaje que muestren por pantalla los nodos esclavos.

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 )





#### Caso práctico 1

```
#include <stdio.h>
#include <mpi.h>
#include <unistd.h>
int main(int argc, char** argv) {
 char host[20], Recv[20]:
 MPI Status status:
 int rank, size, tag=0, i;
 gethostname(host,20):
 MPI Init(&argc, &argv):
 MPI Comm rank(MPI COMM_WORLD, &rank);
 MPI Comm size(MPI COMM WORLD. &size):
 if (rank == 0) {
                 /* proc 0 => send message */
    for (i = 1; i < size; ++i) {
    MPI Send("Hola desde P0!", 20, MPI CHAR, i, tag, MPI_COMM_WORLD);
    printf ("Hola, soy el nodo % de %i nodos. Estoy en la máquina %s y he enviado un mensaje a %i\n",rank,size, host,i);
                /* proc 1 => receive message */
 if (rank != 0) {
    MPI Recv(Recv, 20, MPI CHAR, 0, tag, MPI COMM WORLD, & status);
    printf("Hola! Sov el nodo %i de un total de %i nodos. Estov en la máguina %s v he recibido el mensaje: '%s'\n".
        rank, size, host, Recv);
 MPI Finalize();
```







CENTRO EXTREMEÑO DE TECNOLOGÍAS AVANZADAS

# Casos prácticos

#### Caso práctico 2

Crear un programa sencillo que utilice la función Broadcast para que el nodo maestro mande un entero al resto de nodos.



Int MPI\_Bcast (void \*buffer, int count, MPI\_Datatype datatype, source, MPI\_Comm comm)





#### Caso práctico 2

```
#include <stdio.h>
#include <mpi.h>
#include <string.h>
int main(int argc, char** argv) {
 int rank;
 int n:
 MPI Init(&argc, &argv):
 MPI Comm rank(MPI COMM WORLD, &rank);
                  /* proc 0 => send message */
 if (rank == 0) {
    printf("Hola, soy el nodo maestro, vamos a mandarle el numero 2010 al todos los nodos. ");
    n = 2010;
 MPI_Bcast(&n,1,MPI_INT,0,MPI_COMM_WORLD);
 printf("Hola, soy el proceso %d y el numero es el %d\n",rank,n);
 MPI_Finalize();
```







#### Enviar trabajos al cluster de GPUs

Acceder como usuario tutorial

ssh tutorial@ce-gpu-test.ceta-ciemat.es

- Crear carpeta de usuario en /home/tutorial/minombredeusuario
- Enviar un trabajo:

qsub -q gpu -pe gpupe <numero\_nodos> script.sh

Script.sh

```
export MPICH2_ROOT=/opt/mpi/mpibull2-1.3.9-14.s
export PATH=$MPICH2_ROOT/bin:$PATH
export MPD_CON_EXT="sge_$JOB_ID.$SGE_TASK_ID"
```

echo "Got \$NSLOTS slots."

mpiexec -machinefile \$TMPDIR/machines -n \$NSLOTS /opt/exp\_soft/tutorial/mpisend

exit 0









#### Referencias

#### Referencias

- MPI: Message Passing Interface. Universidad Carlos III Madrid. http://svn.cetaciemat.es/gpus/trunk/Bibliograf%c3%ada/MPI/Curso-MPI CarlosIIIMadrid.pdf
- Interfaz de Paso de Mensajes Wikipedia. http://es.wikipedia.org/wiki/Interfaz de Paso de Mensaies
- Message Passing Interface (MPI) Blaise Barney, Lawrence Livermore National Laboratory. https://computing.llnl.gov/tutorials/mpi/
- MPI: The Complete Reference http://www.netlib.org/utk/papers/mpi-book/mpibook.html
- MPI Routines (official): http://www.mcs.anl.gov/research/projects/mpi/www/







## 1 Introducción a MPI

#### Cómo funciona MPI

- SEGUNDA PARTE DEL TUTORIAL MPI
- Mpich, demonios... MPD...

Fondo Europeo de Desarrollo Regional



