

# Programação em Computação Paralela e Distribuída

Laboratório, ERBASE 2009

## MPI



Domingo Giménez

**Universidad de Murcia, Spain**

Grupo de Computación Paralela

<http://dis.um.es/~domingo>

Murilo Boratto

**UFBa**

Leandro Coelho

**UNEB**



# Noções básicas: história

---

- MPI: Message Passing Interface
  - Uma especificação para passo de mensagens
  - A primeira livreria de passagem de mensagens estándar e portable
  - Por consenso MPI Forum. Participantes de umas 40 organizações
  - Acabado e publicado em maio 1994. Actualizado em junho 1995
- Previamente PVM: Parallel Virtual Machine
- Outras versões: MPI2, HMPI



# Noções básicas: conteúdo

---

- Estandarização
- Portabilidade: multiprocesadores, multicomputadores, redes, heterogéneos...
- Boas prestações, se está disponível para o sistema
- Ampla funcionalidad. Umas 140 funções. Programas em C ou Fortran e chamadas as funções da livraria MPI
- Implementaciones gratuitas (mpich, lam, ...)



# Exemplo

---

Todos os processos executam o mesmo código

As variables estão replicadas nas memórias  
de todos os processos

Trabalham com MPI entre **Init** e **Finalize**

codigo3-5.c

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

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 desde el proceso %d de %d\n", myrank, size);
    MPI_Finalize();
}
```



# Exemplo

---

**compilação:**

**mpicc** codigo3-5.c -o codigo3-5 -O3

**execução:**

**\$> mpirun -np 4 codigo3-5**

**executa 4 processos iguais nos processadores estabelecidos por defeito**

**\$> mpirun -np 4 -machinefile maquinas codigo3-5**

**executa-os nos processadores onde se indica no arquivo de máquinas**

**\$> mpirun n0,1,2,3 codigo3-5**

**os executa nos nodos 0, 1, 2 e 3**



# Exemplo

---

**compilação:**

**mpicc** codigo3-5.c -o codigo3-5 -O3

**execução:**

**\$> mpirun -np 4 codigo3-5**

**executa 4 processos iguais nos processadores estabelecidos**

**\$> mpirun -np 4 -machinefile maquinas codigo3-5**

**executa-os nos processadores onde se indica no arquivo**

**\$> mpirun n0,1,2,3 codigo3-5**

**os executa nos nodos 0, 1, 2 e 3**

```
$> mpirun -np 4 codigo3-5
Hola desde el proceso 0 de 4
Hola desde el proceso 2 de 4
Hola desde el proceso 1 de 4
Hola desde el proceso 3 de 4
```

```
$> mpirun -np 4 codigo3-5
Hola desde el proceso 0 de 4
Hola desde el proceso 3 de 4
Hola desde el proceso 1 de 4
Hola desde el proceso 2 de 4
```



# Funções

---

- Incluir livreria: `#include <mpi.h>`
- Formato das funções: `erro=MPI_nome(parâmetros ...)`

`int MPI_Init ( int *argc , char ***argv )`

- Inicializa MPI

`int MPI_Finalize ( )`

- Finaliza MPI



# Funções

---

**MPI\_Comm\_rank** ( MPI\_Comm comm , int \*rank)

- Obtém o identificador de processo que chama à função
- É um valor entre 0 e  $p-1$

**MPI\_Comm\_size** ( MPI\_Comm comm , int \*size)

- Obtém o número total de processos,  $p$
- Comunicador: Conjunto de processos em que se fazem comunicações

MPI\_COMM\_WORD: o mundo dos processos MPI





# Comunicações

---

int **MPI\_Send** ( void \*buffer , int contador , MPI\_Datatype tipo ,  
int destino , int tag , MPI\_Comm comunicador )

int **MPI\_Recv** ( void \*buffer , int contador , MPI\_Datatype tipo ,  
int origen , int tag , MPI\_Comm comunicador , MPI\_Status  
\*estado)

- pode-se indicar qualquer tipo ou origem com as constantes

MPI\_ANY\_TAG

MPI\_ANY\_SOURCE



# Comunicações

---

usam-se tipos MPI:

Básicos:

<b>MPI_CHAR</b>	signed char
<b>MPI_SHORT</b>	signed short int
<b>MPI_INT</b>	signed int
<b>MPI_LONG</b>	signed long int
<b>MPI_UNSIGNED_CHAR</b>	unsigned char
<b>MPI_UNSIGNED_SHORT</b>	unsigned short int
<b>MPI_UNSIGNED</b>	unsigned int
<b>MPI_UNSIGNED_LONG</b>	unsigned long int
<b>MPI_FLOAT</b>	float
<b>MPI_DOUBLE</b>	double
<b>MPI_LONG_DOUBLE</b>	long double
<b>MPI_BYTE</b>	
<b>MPI_PACKED</b>	



# Comunicações síncronas

---

- Nas comunicações síncronas o processo que envia continua sua execução, e o que recebe espera a que chegue a mensagem
- O momento em que o que envia continua depende da função de envio:
  - **MPI\_Ssend**: acaba quando começa a recepção
  - **MPI\_Bsend**: envio com buffer
  - **MPI\_Send**: síncrono com buffer
  - **MPI\_Rsend**: acaba independentemente da recepção



# Comunicações não síncronas

---

- **MPI\_Isend**(buf, count, datatype, dest, tag, comm, request)
- **MPI\_Irecv**(buf, count, datatype, source, tag, comm, request)
- **MPI\_Wait**(request, status): volta se a operação completou-se, espera até que se completa
- **MPI\_Test**(request, flag, status): devolve um **flag** dizendo se a operação completou-se



# Exemplos

---

- **codigo3-6.c** mostra um exemplo de "Hello world" com funções de comunicação
- **codigo3-7.c** e **integral.c** são exemplos de aproximação de uma integral definida
- Pode-se praticar com estes exemplos mudando o tipo de dados e de comunicações, e utilizando comunicações assíncronas
- Em **integral.c**

Custo de envio:

$$3(p-1)(t_s+t_w)$$



# Comunicações colectivas

---

- **MPI\_Barrier( )** bloqueia os processos até que a chamam todos
- **MPI\_Bcast( )** broadcast do processo raiz a todos os demais
- **MPI\_Gather( )** recebe valores de um grupo de processos
- **MPI\_Scatter( )** distribui um buffer em partes a um grupo de processos
- **MPI\_Alltoall( )** envia dados de todos os processos a todos
- **MPI\_Reduce( )** combina valores de todos os processos
- **MPI\_Reduce\_scatter( )** combina valores de todos os processos e distribui
- **MPI\_Scan( )** redução prefixa (dados 0,...,i-1 em i)



# Comunicações colectivas

---

MPI Operator	Operation
<hr/>	
<b>MPI_MAX</b>	maximum
<b>MPI_MIN</b>	minimum
<b>MPI_SUM</b>	sum
<b>MPI_PROD</b>	product
<b>MPI_LAND</b>	logical and
<b>MPI_BAND</b>	bitwise and
<b>MPI_LOR</b>	logical or
<b>MPI BOR</b>	bitwise or
<b>MPI_LXOR</b>	logical exclusive or
<b>MPI_BXOR</b>	bitwise exclusive or
<b>MPI_MAXLOC</b>	max value and location
<b>MPI_MINLOC</b>	min value and location



# Exemplo

---

- **integralCC.c** calcula a integral utilizando comunicações colectivas

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

todos os processadores no comunicador executam a função, e todos indicam qual é o processo raiz

Custo de envio:

$$3(t_{sb} + t_{wb})$$

mas  $t_{sb}$  e  $t_{wb}$  são maiores que  $t_s$  e  $t_w$





# Agrupamento de dados

---

- Para reduzir o custo das comunicações podem-se agrupar os dados e realizar uma comunicação com os três dados juntos

Custo de envio:

$$(p-1)(t_s+3t_w)$$

- Agrupando os dados num array (exemplo **integralAR.c**)
- Com empacotamento
- Com tipos derivados



# Empacotamento

---

exemplo **integralPack.c**

```
int MPI_Pack( void *pack_data, int in_count, MPI_Datatype datatype,  
             void *buffer, int size, int *position_ptr, MPI_Comm comm)
```

```
int MPI_Unpack( void *buffer, int size, int *position_ptr, void  
               *unpack_data, int count, MPI_Datatype datatype, MPI_Comm comm)
```

agrupam-se os dados antes de enviá-los, e depois de recebê-los se desempacotam no mesmo ordem em que se agruparam

- entender o funcionamento do exemplo



# Tipos derivados

---

exemplo **integralTD.c**

- Criam-se em tempo de execução
- Podem-se construir tipos de maneira recursiva
- A criação de tipos requer trabalho adicional
- Especifica-se a disposição dos dados no tipo:

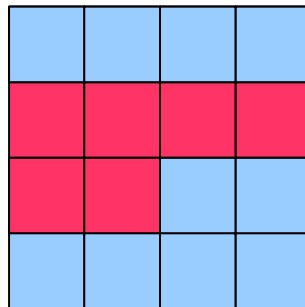
```
int MPI_Type_Struct (int count, int *array_of_block_lengths,  
MPI_Aint *array_of_displacements, MPI_Datatype *array_of_types,  
MPI_Datatype *newtype )
```

- Chamar à rotina **MPI\_Type\_commit** antes de usar o novo tipo MPI
- Não se calculam as direcções como em C, usa-se a rotina **MPI\_Address**

# Tipos derivados

- Se os dados que constituem o novo tipo são um subconjunto de entradas há mecanismos especiais:
- `int MPI_Type_contiguous( int count, MPI_Datatype oldtype, MPI_Datatype *newtype)`

cria um tipo derivado formado por **count** elementos do tipo **oldtype** contíguos em memória.

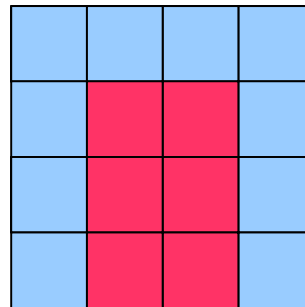


# Tipos derivados

- `int MPI_Type_vector( int count, int block_lenght, int stride, MPI_Datatype element_type, MPI_Datatype *newtype);`

formado por **count** elementos, a cada um com **block\_lenght** elementos de tipo **element\_type**. **stride** é o número de elementos de tipo **element\_type** entre elementos sucessivos de **new\_type**.

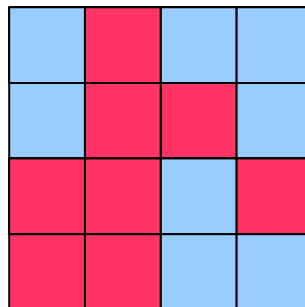
Os elementos podem ser entradas igualmente espaciadas num array



# Tipos derivados

- `int MPI_Type_indexed( int count, int *array_of_block_lengths, int *array_of_displacements, MPI_Datatype element_type, MPI_Datatype *newtype);`

cria um tipo derivado com **count** elementos; na cada elemento há **array\_of\_block\_lengths[i]** dados de tipo **element\_type**, e o deslocamento é **array\_of\_displacements[i]** unidades do tipo **element\_type** desde o começo de **newtype**





# Comunicadores

---

- **MPI\_COMM\_WORLD** inclui a todos os processos
- Pode-se definir comunicadores com um número menor de processos: para comunicar dados na cada bicha ou coluna de processos na malha...
- As comunicações em diferentes comunicadores podem ser ao mesmo tempo
- Consta de:
  - um grupo: colecção ordenada de processos numerados  $0, \dots, p-1$
  - um contexto: um identificador que associa o sistema ao grupo
  - e a um comunicador pode-se-lhe associar uma topologia virtual



# Comunicadores

---

- Criação de um comunicador cujos processos são os da primeira bicha de nossa malha virtual. `MPI_COMM_WORLD` consta de  $p=q^2$  processos agrupados em  $q$  bichas e colunas. O processo número  $x$  tem as coordenadas  $(x \text{ div } q, x \text{ mod } q)$

```
MPI_Group MPI_GROUP_WORLD;
```

```
MPI_Group first_row_group;
```

```
MPI_Comm first_row_comm;
```

```
int row_size;
```

```
int *process_ranks;
```

```
process_ranks=(int *) malloc(q*sizeof(int));
```

```
for(proc=0;proc<q;proc++)
```

```
    process_ranks[proc]=proc;
```

```
MPI_Comm_group(MPI_COMM_WORLD,&MPI_GROUP_WORLD);
```

```
MPI_Group_incl(MPI_GROUP_WORLD,q,process_ranks,&first_row_group);
```

```
MPI_Comm_create(MPI_COMM_WORLD,first_row_group,&first_row_comm);
```





# Comunicadores

---

```
int MPI_Comm_split( MPI_Comm old_comm, int split_key, int rank_key,  
                    MPI_Comm *new_comm)
```

- ▣ Cria um novo comunicador para a cada valor de **split\_key**
- ▣ Os processos com o mesmo valor de **split\_key** formam um grupo
- ▣ Se dois processos *a* e *b* têm o mesmo valor de **split\_key** e o **rank\_key** de *a* é menor que o de *b*, no novo grupo *a* tem identificador menor que *b*
- ▣ Todos os processos no comunicador devem a chamar
- ▣ Os processos que não se quer incluir em nenhum novo comunicador podem utilizar o valor **MPI\_UNDEFINED** em **rank\_key**, com o que o valor de volta de **new\_comm** é **MPI\_COMM\_NULL**



# Topologías

---

- A um grupo pode-se-lhe associar uma topología virtual para melhorar as prestações
- As topologías podem ser:
  - de grafo em general
  - de malha ou cartesiana
- Há funções de:
  - criar topologías
  - obter identificadores de processo na topología
  - particionar malhas....