



MPI (Message Passing Interface)

于策



Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- MPI并行程序的基本模式
 - 对等模式
 - 主从模式
- MPI数据及进程
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信



Outline

- **MPI概述**

- MPI基本调用
- 点到点通信与组通信

- MPI并行程序的基本模式

- 对等模式
- 主从模式

- MPI数据及进程

- 自定义数据类型
- 虚拟进程拓扑

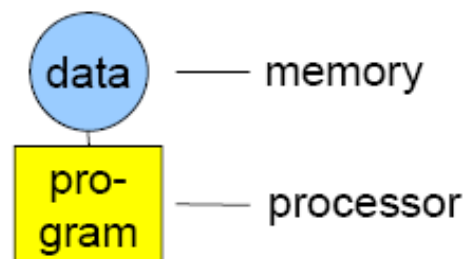
- 消息传递方式

- MPI通信模式
- 非阻塞通信

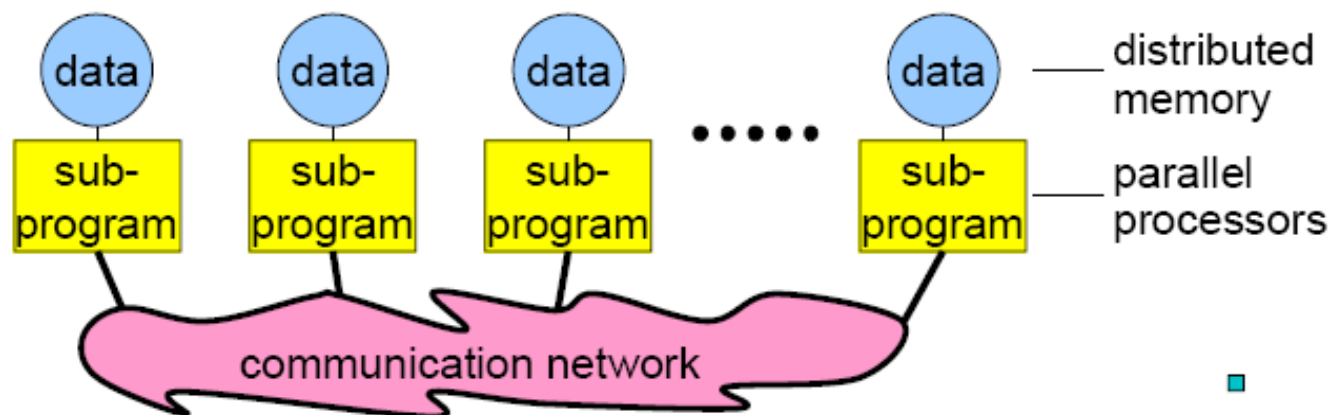


MPI概述

■ 串程序序



■ MPI并行程序





MPI (Message passing interface)

- **MPI**是一种标准或规范的代表，而不特指某一个对它的具体实现。 **MPI**同时也是一种消息传递编程模型，并成为这种编程模型的代表和事实上的标准。
 - 迄今为止所有的并行计算机制造商都提供对**MPI**的支持，可以在网上免费得到**MPI**在不同并行计算机上的实现。
- **MPI**的实现是一个库，而不是一门语言。
 - 可以把**FORTRAN+MPI**或**C+MPI** 看作是一种在原来串行语言基础之上扩展后得到的并行语言。



MPI程序示例: Hello World!

Fortran

```
PROGRAM hello
  INCLUDE 'mpif.h'
  INTEGER err
  CALL MPI_INIT(err)
  PRINT *, "hello world!"
  CALL MPI_FINALIZE(err)

END
```

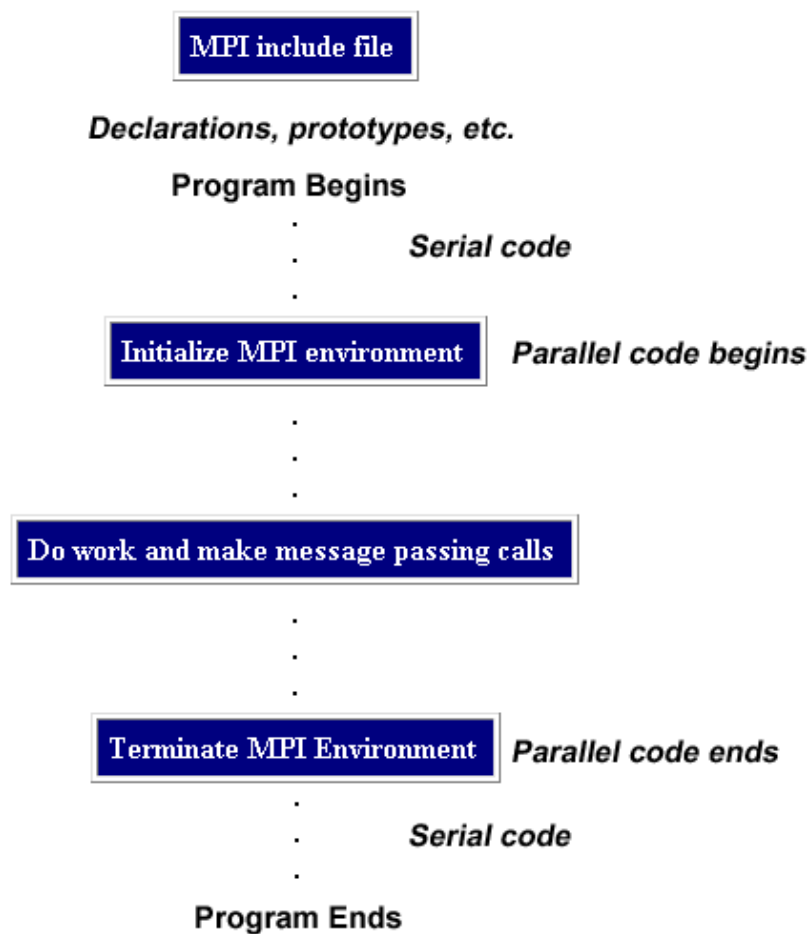
C

```
#include <stdio.h>
#include <mpi.h>
void main (int argc, char * argv[])
{
  int err;

  err = MPI_Init(&argc, &argv);
  printf( "Hello world!\n" );
  err = MPI_Finalize();
}
```



MPI程序结构





Outline

- **MPI概述**
 - **MPI基本调用**
 - 点到点通信与组通信
- MPI并行程序的基本模式
 - 对等模式
 - 主从模式
- MPI数据及进程
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信



MPI 的六个基本接口

- 开始与结束
 - MPI_INIT
 - MPI_FINALIZE
- 进程身份标识
 - MPI_COMM_SIZE
 - MPI_COMM_RANK
- 发送与接收消息
 - MPI_SEND
 - MPI_RECV



MPI 程序的开始与结束

- MPI代码开始之前必须进行如下调用:

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

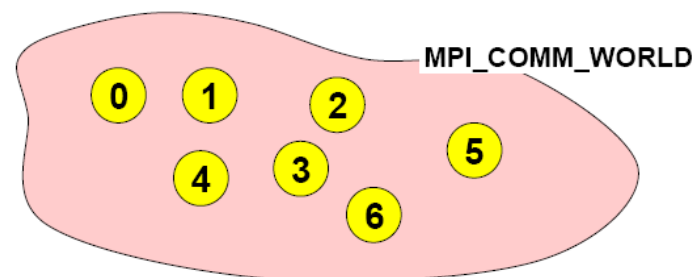
- MPI代码的最后一行必须是:

```
MPI_Finalize();
```

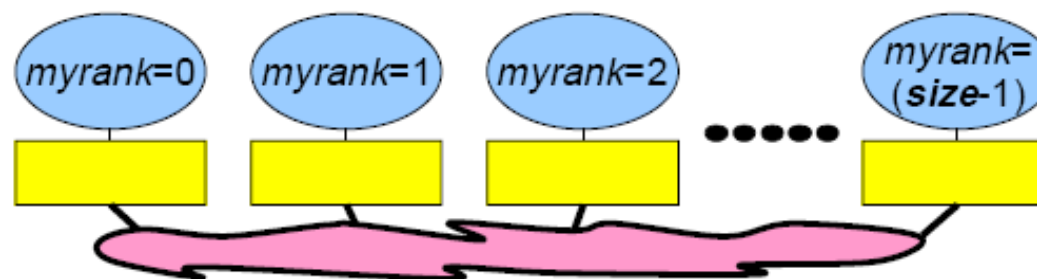
- 如果没有此行, MPI程序将不会终止。



MPI进程身份标识



- 通信域
 - 缺省的通信域为 `MPI_COMM_WORLD`
- `MPI_Comm_size(MPI_COMM_WORLD, &size)`
 - 获得缺省通信域内所有进程数目，赋值给 `size`
- `MPI_Comm_rank(MPI_COMM_WORLD, &myrank)`
 - 获得进程在缺省通信域的编号，赋值给 `myrank`





发送和接收消息

`MPI_Send(buf, count, datatype, dest, tag, comm)`

Address of
send buffer

Number of items
to send

Datatype of
each item

Rank of destination
process

Message tag

Communicator

`MPI_Recv(buf, count, datatype, src, tag, comm, status)`

Address of
receive buffer

Maximum number
of items to receive

Datatype of
each item

Rank of source
process

Message tag

Communicator

Status after operation



一个计算 $\sum \text{foo}(i)$ 的MPI SPMD消息传递程序

```
#include "mpi.h"
int foo(i)
int i;
{...}
main(argc, argv)
int argc;
char* argv[]
{
    int i, tmp, sum=0, group_size, my_rank, N;
    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &group_size);
    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
    if (my_rank==0) {
        printf("Enter N:");
        scanf("%d",&N);
        for (i=1;i<group_size;i++)
            MPI_Send(&N, 1, MPI_INT, i, i, MPI_COMM_WORLD);
        for (i=my_rank;i<N;i=i+group_size) sum=sum+foo(i);
        for (i=1;i<group_size;i++) {
            MPI_Recv(&tmp,1,MPI_INT,i,i,MPI_COMM_WORLD,&status);
            sum=sum+tmp;
        }
        printf("\n The result = %d", sum);
    }
    else {
        MPI_Recv(&N,1,MPI_INT,0,i,MPI_COMM_WORLD,&status);
        for (i-my_rank;i<N;i=i+group_size) sum=sum+foo(i);
        MPI_Send(&sum,1,MPI_INT,0,i,MPI_COMM_WORLD);
    }
    MPI_Finalize();
}
```

初始化MPI环境

获得总进程数

得到每个进程在组
中的编号

发送消息

接收消息

终止MPI环境



Outline

- **MPI概述**

- MPI基本调用
- 点到点通信与组通信

- MPI并行程序的基本模式

- 对等模式
- 主从模式

- MPI数据及进程

- 自定义数据类型
- 虚拟进程拓扑

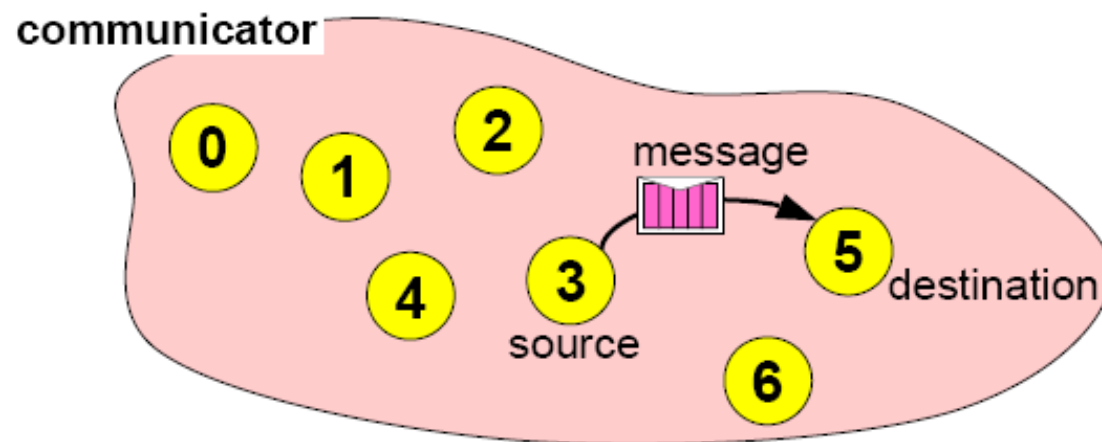
- 消息传递方式

- MPI通信模式
- 非阻塞通信



点到点通信

- 对于某一消息
 - 唯一发送进程
 - 唯一接收进程





MPI_Send

MPI_Send(buffer, count, datatype, destination, tag, communicator)

- `MPI_Send(&N, 1, MPI_INT, i, i, MPI_COMM_WORLD);`
- 第一个参数指明消息缓存的起始地址，即存放要发送的数据信息。
- 第二个参数指明消息中给定的数据类型有多少项，数据类型由第三个参数给定。
- 数据类型要么是基本数据类型，要么是导出数据类型，后者由用户生成指定一个可能是由混合数据类型组成的非连续数据项。
- 第四个参数是目的进程的标识符(进程编号)。
- 第五个是消息标签。
- 第六个参数标识进程组和上下文，即通信域。通常，消息只在同组的进程间传送。但是MPI允许通过intercommunicators在组间通信。



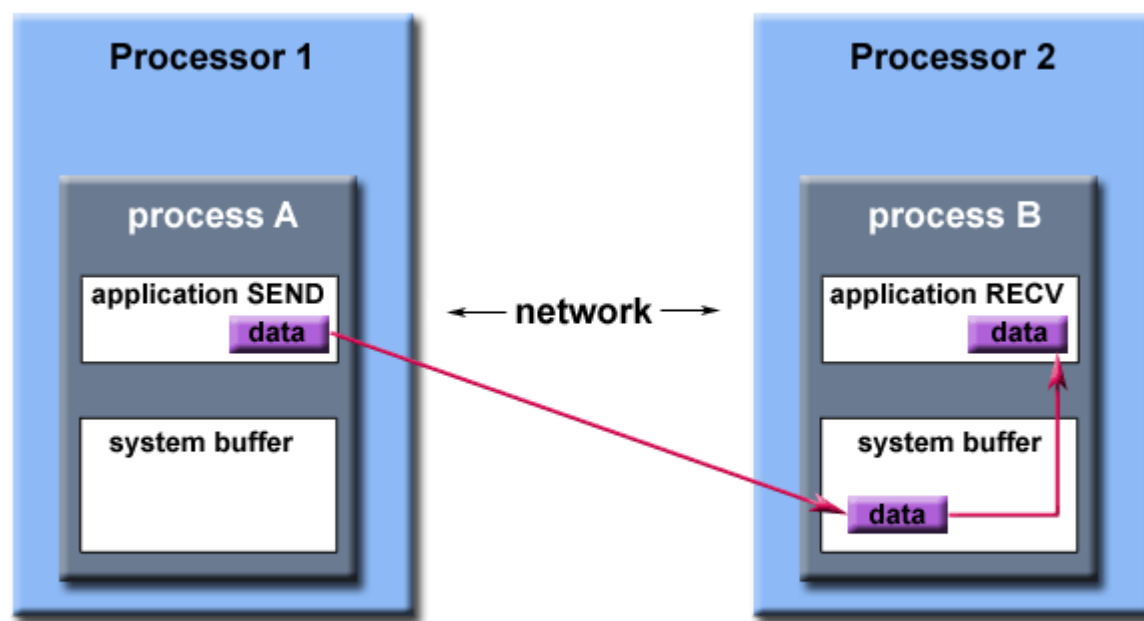
MPI_Receive

MPI_Recv(address, count, datatype, source, tag, communicator, status)

- `MPI_Recv(&tmp, 1, MPI_INT, i, i, MPI_COMM_WORLD, &Status)`
- 第一个参数指明接收消息缓冲的起始地址，即存放接收消息的内存地址。
- 第二个参数指明给定数据类型可以被接收的最大项数，接收到的实际项数可能少一些。
- 第三个参数指明接收的数据类型。
- 第四个参数是源进程标识符 (编号)。
- 第五个是消息标签。
- 第六个参数标识一个通信域。
- 第七个参数是一个指针，指向一个结构： `MPI_Status Status`
 - 存放有关接收消息的各种信息。 (`Status.MPI_SOURCE`, `Status.MPI_TAG`)
 - `MPI_Get_count(&Status, MPI_INT, &C)` 读出实际接收到的数据项数。



消息的接收（系统缓存）



Path of a message buffered at the receiving process



标签的使用

为什么要使用消息标签(Tag)?

这段代码需要传送A的前32个字节进入X，传送B的前16个字节进入Y。但是，如果消息B尽管后发送但先到达进程Q，就会被第一个recv()接收在X中。

使用标签可以避免这个错误。

未使用标签

Process P:

```
send(A,32,Q)  
send(B,16,Q)
```

Process Q:

```
recv(X, 32, P)  
recv(Y, 16, P)
```

使用了标签

Process P:

```
send(A,32,Q,tag1)  
send(B,16,Q,tag2)
```

Process Q:

```
recv (X, 32, P, tag1)  
recv (Y, 16, P, tag2)
```



标签的使用

Process P:

```
send (request1,32, Q)
```

Process R:

```
send (request2, 32, Q)
```

Process Q:

```
while (true) {  
    recv (received_request, 32, Any_Process);  
    process received_request;  
}
```

使用标签的另一个原因是可以简化对下列情形的处理。

假定有两个客户进程P和R，每个发送一个服务请求消息给服务进程Q。

Process P:

```
send(request1, 32, Q, tag1)
```

Process R:

```
send(request2, 32, Q, tag2)
```

Process Q:

```
while (true){  
    recv(received_request, 32, Any_Process, Any_Tag, Status);  
    if (Status.Tag==tag1) process received_request in one way;  
    if (Status.Tag==tag2) process received_request in another way;  
}
```



组通信

- 一到多 (Broadcast, Scatter)
- 多到一 (Reduce, Gather)
- 多到多 (Allreduce, Allgather)

MPI_Bcast() - 从root进程广播到所有其他进程

MPI_Gather() - 从本组内进程收集数据

MPI_Scatter() - 将缓冲区内数据散播到组内进程

MPI_Alltoall() - 所有进程的数据发送到所有进程

MPI_Reduce() - 将所有进程的数据归约为单值

MPI_Reduce_scatter() - 归约数据并进行散播

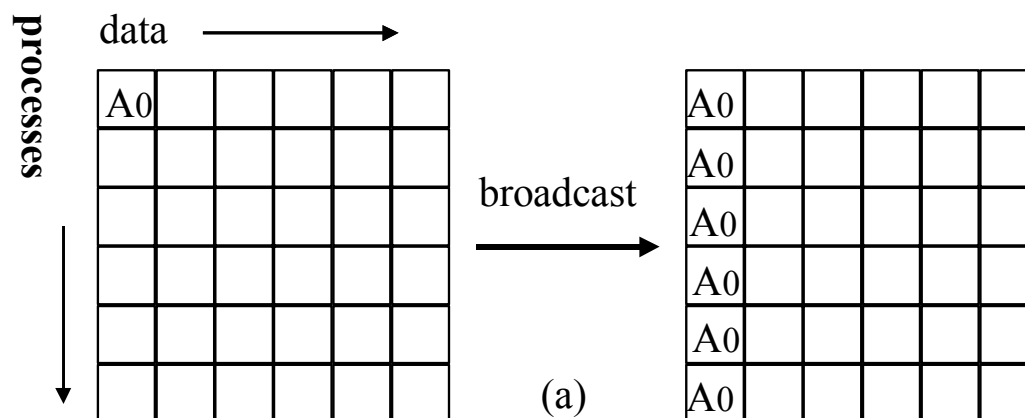


广播 (Broadcast)

MPI_Bcast(Address, Count, Datatype, *Root*, *Comm*)

标号为Root的进程发送相同的消息给标记为Comm的通信子中的所有进程。

消息的内容如同点对点通信一样由三元组(Address, Count, Datatype)标识。对Root进程来说，这个三元组既定义了发送缓冲也定义了接收缓冲。对其它进程来说，这个三元组只定义了接收缓冲。





MPI_Bcast() 语法

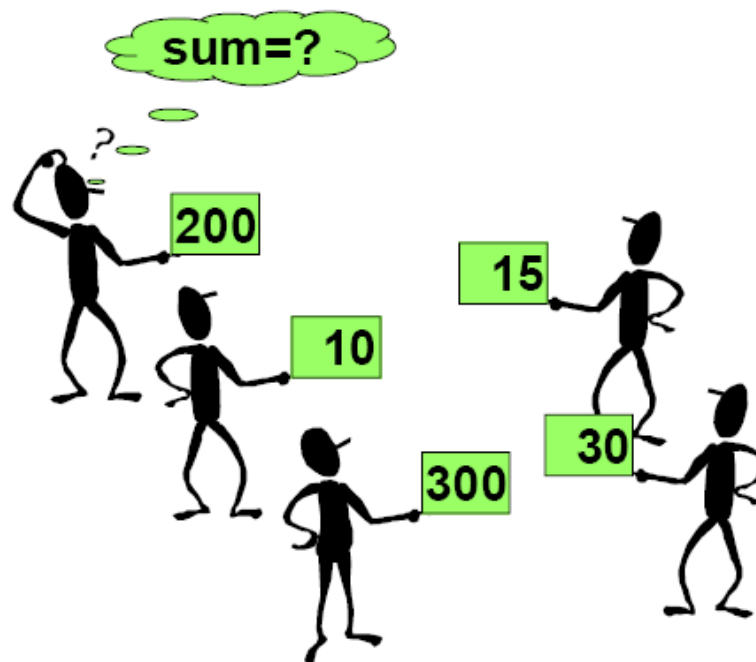
```
MPI_Bcast(mess, count, MPI_INT,  
          root, MPI_COMM_WORLD);
```

mess	消息缓冲区指针
count	发送数据的数量
MPI_INT	数据类型
root	发送数据的进程号
MPI_COMM_WORLD	通信域



归约 (Reduce)

- 所有进程向同一进程发送消息，与broadcast的消息发送方向相反。
- 接收进程对所有收到的消息进行归约处理。
- 归约操作：
 - MAX, MIN, SUM, PROD, LAND, BAND, LOR, BOR, LXOR, BXOR, MAXLOC, MINLOC





MPI_Reduce() 语法

```
MPI_Reduce(&dataIn, &result, count,  
           MPI_DOUBLE, MPI_SUM, root,  
           MPI_COMM_WORLD);
```

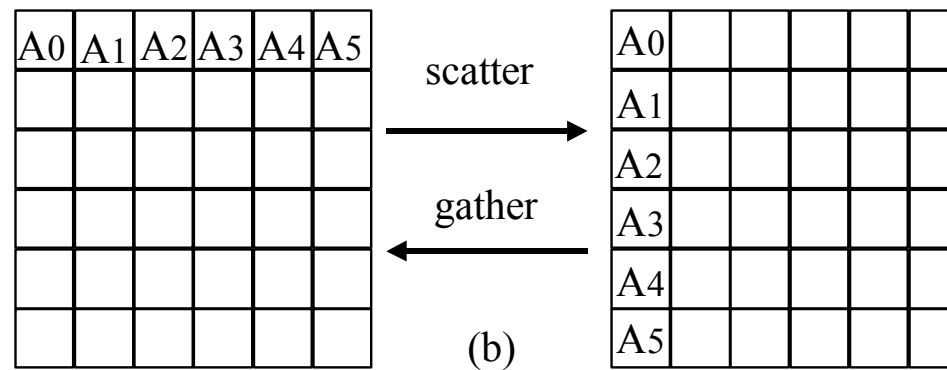
dataIn	每个进程发送的数据
result	归约操作的结果
count	dataIn和result中数据的数目
MPI_DOUBLE	dataIn和result的数据类型
MPI_SUM	归约操作
root	接收数据的进程号
MPI_COMM_WORLD	通信域



Scatter and Gather

`MPI_Scatter` (`SendAddress`, `SendCount`, `SendDatatype`,
`RecvAddress`, `RecvCount`, `RecvDatatype`, `Root`, `Comm`)

`MPI_Gather` (`SendAddress`, `SendCount`, `SendDatatype`,
`RecvAddress`, `RecvCount`, `RecvDatatype`, `Root`, `Comm`)



`MPI_Scatter`

Root进程发送给所有n个进程各发送一个不同的消息，包括自己。这n个消息在Root进程的发送缓冲区中按标号的顺序有序地存放。每个接收缓冲由三元组(`RecvAddress`, `RecvCount`, `RecvDatatype`)标识。非Root进程忽略发送缓冲。对Root进程，发送缓冲由三元组(`SendAddress`, `SendCount`, `SendDatatype`)标识。

`MPI_Gather`

Root进程接收各个进程(包括它自己)的消息。这n个消息的连接按序号rank进行，存放在Root进程的接收缓冲中。每个发送缓冲由三元组(`SendAddress`, `SendCount`, `SendDatatype`)标识。非Root进程忽略接收缓冲。对Root进程，发送缓冲由三元组(`RecvAddress`, `RecvCount`, `RecvDatatype`)标识。



多对多通信

- **MPI_Allreduce**

- 语法与reduce类似，但无root参数
- 所有进程都将获得结果

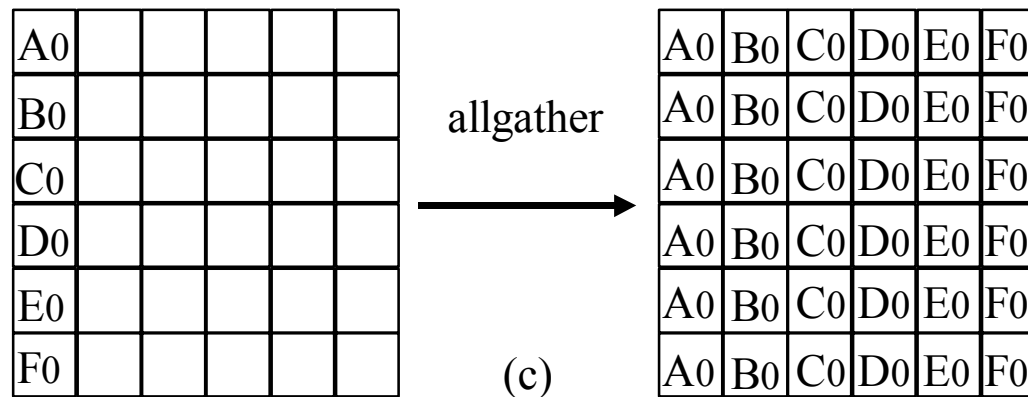
- **MPI_Allgather**

- 语法与gather类似，但无root参数
- 所有进程都将获得结果数组



Allgather

MPI_Allgather (SendAddress, SendCount, SendDatatype,
***RecvAddress, RecvCount, RecvDatatype, Comm*)**



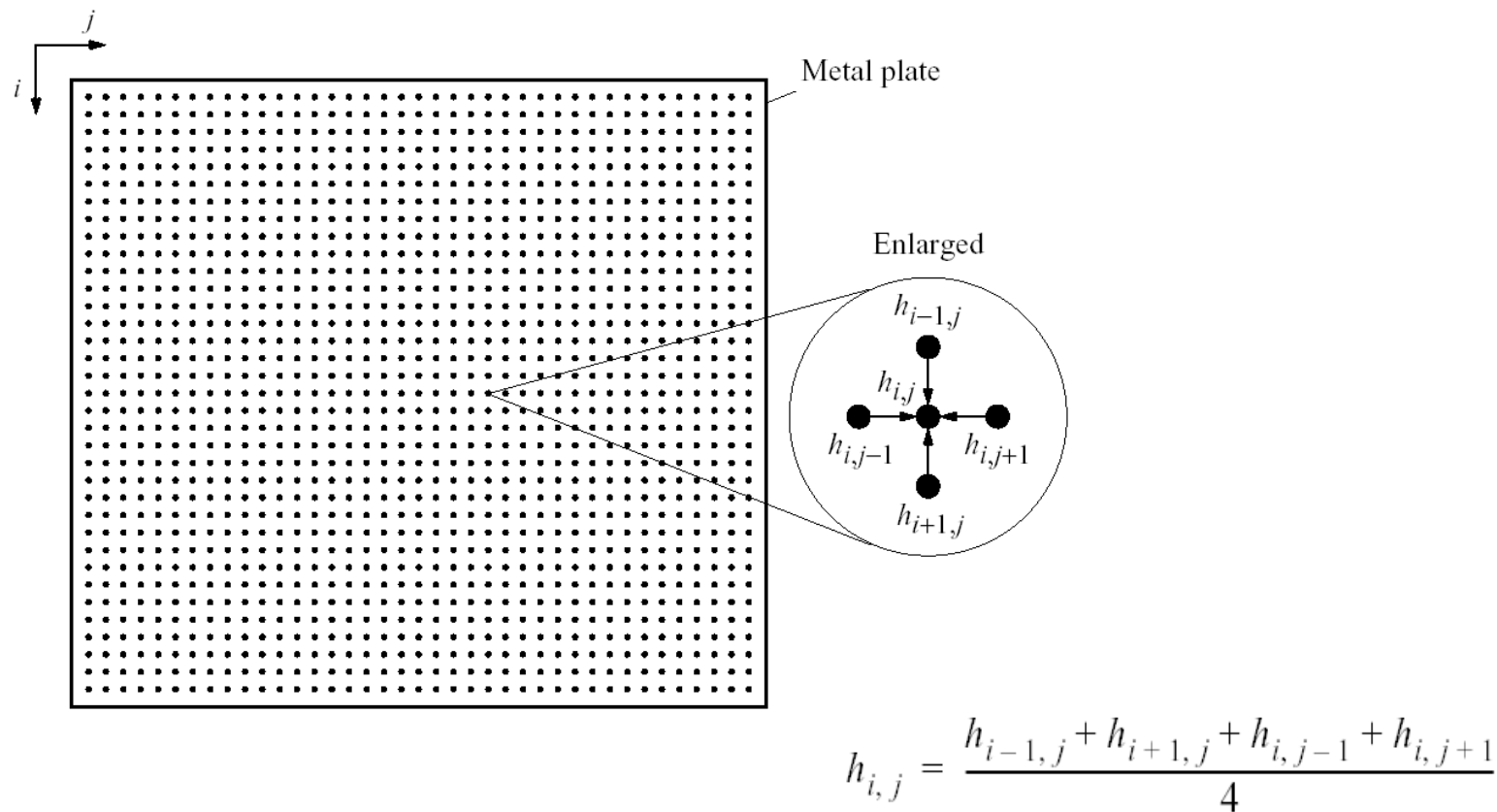


Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- **MPI并行程序的基本模式**
 - **对等模式**
 - 主从模式
- MPI数据及进程
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信



对等模式示例：MPI实现Jacobi迭代





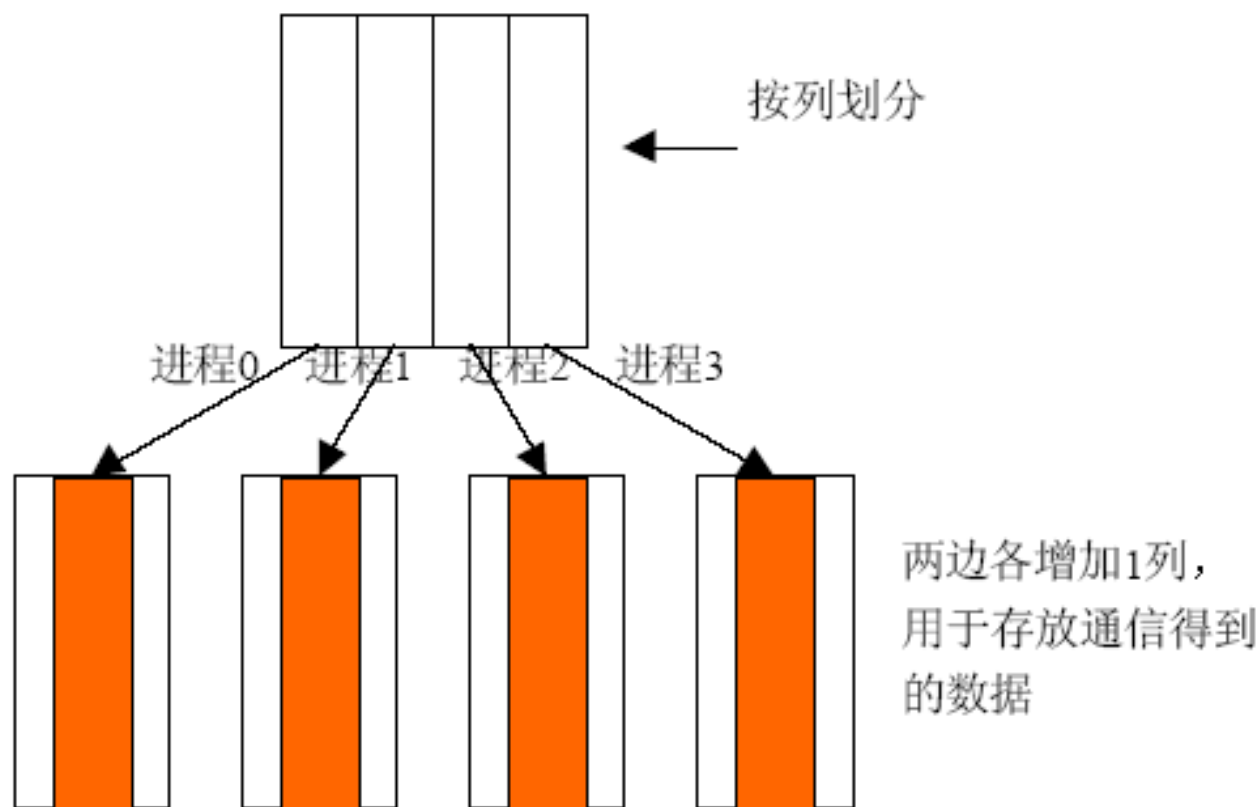
Jacobi迭代

■ 伪代码描述:

```
...  
REAL A(N+1,N+1), B(N+1,N+1)  
...  
DO K=1,STEP  
  DO J=1,N  
    DO I=1,N  
      B(I,J)=0.25*(A(I-1,J)+A(I+1,J)+A(I,J+1)+A(I,J-1))  
    END DO  
  END DO  
  DO J=1,N  
    DO I=1,N  
      A(I,J)=B(I,J)  
    END DO  
  END DO  
END DO
```

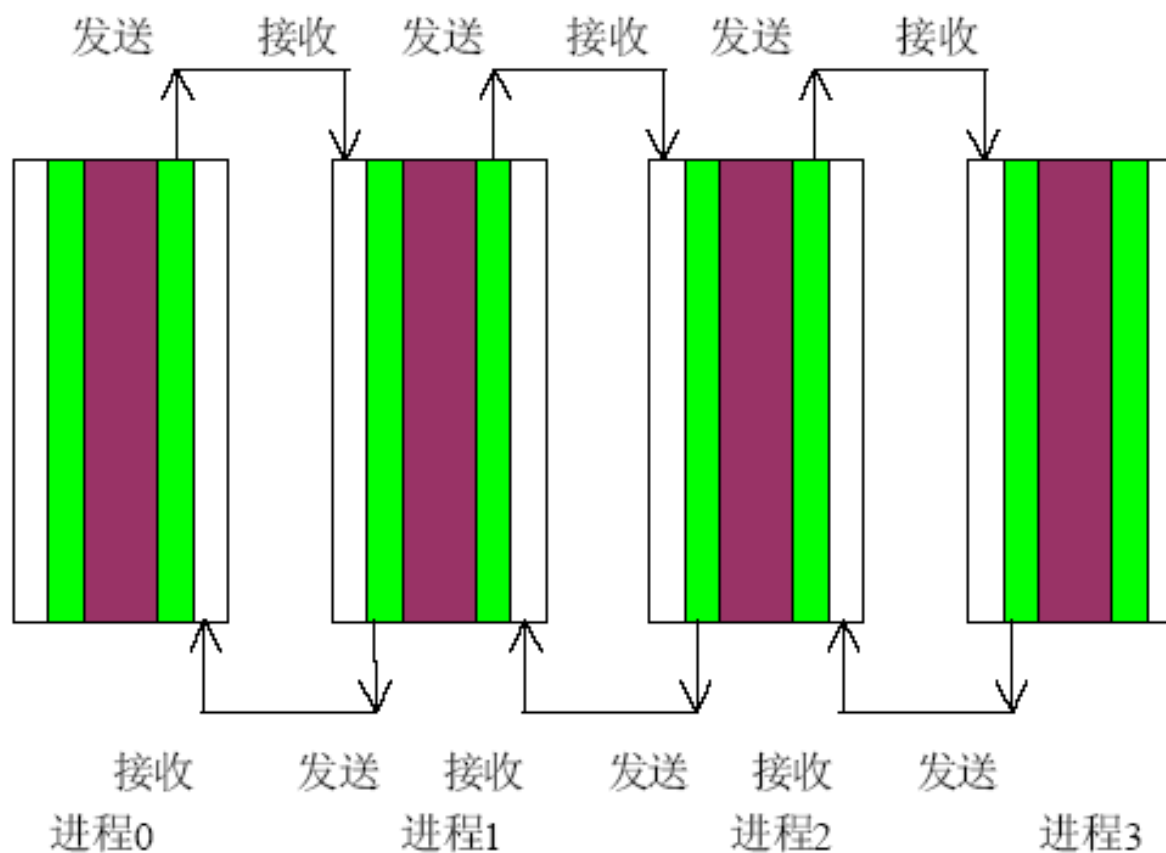


Jacobi迭代：数据划分



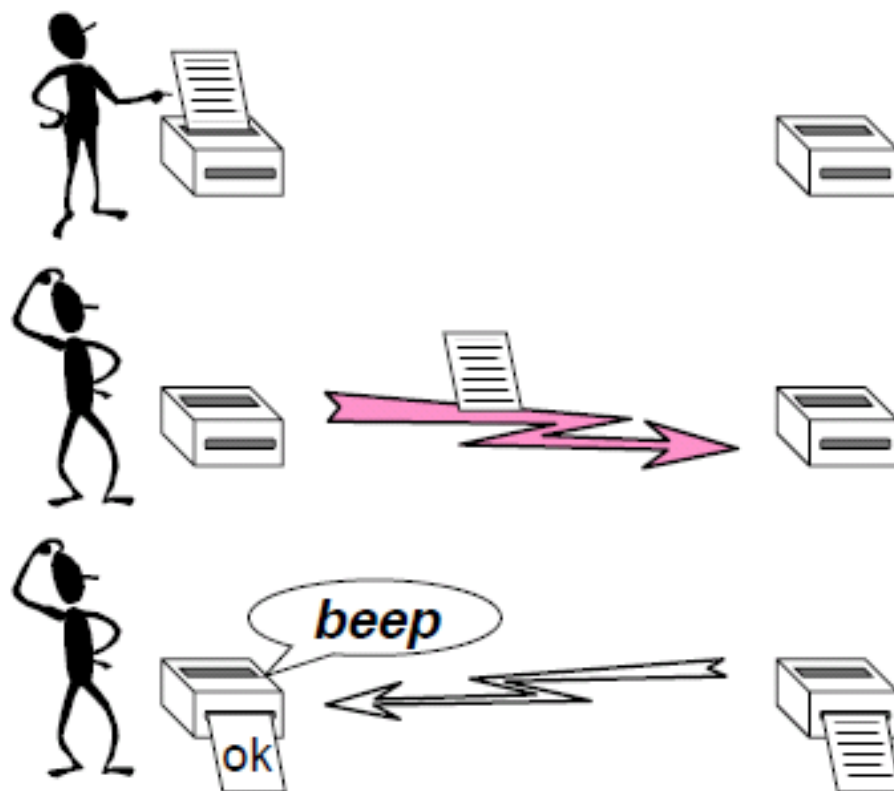


Jacobi迭代：通信



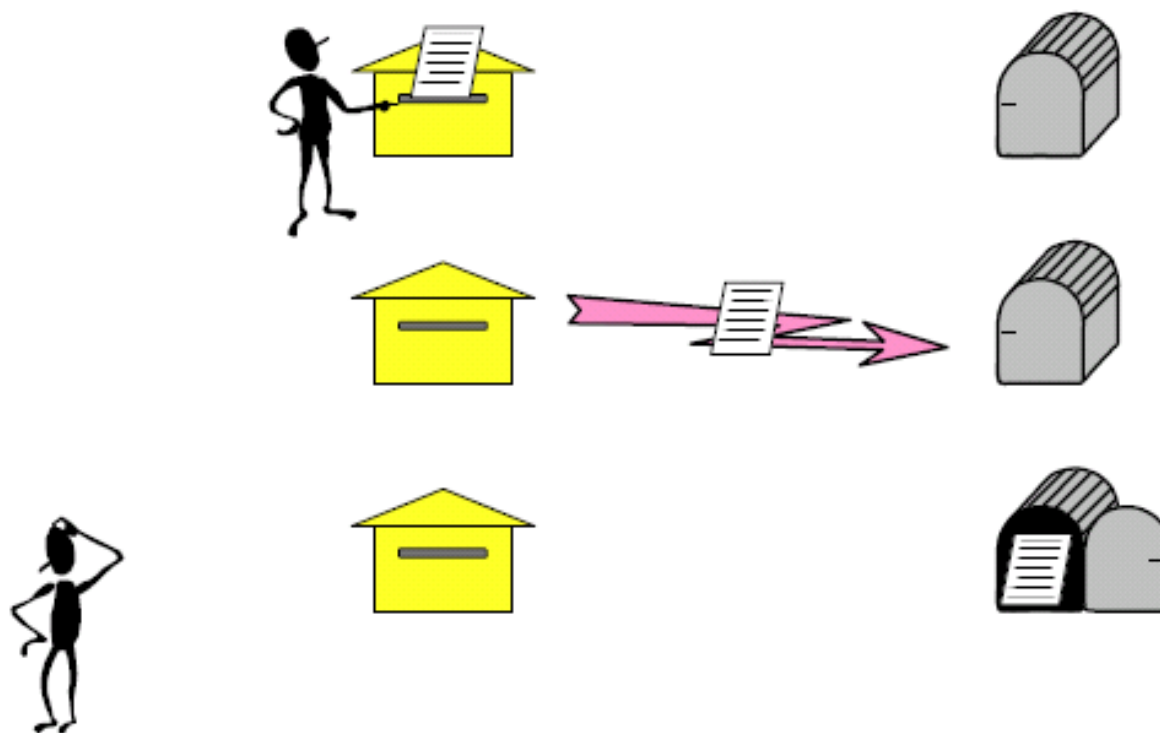


同步发送（阻塞）





异步发送（缓存，非阻塞）





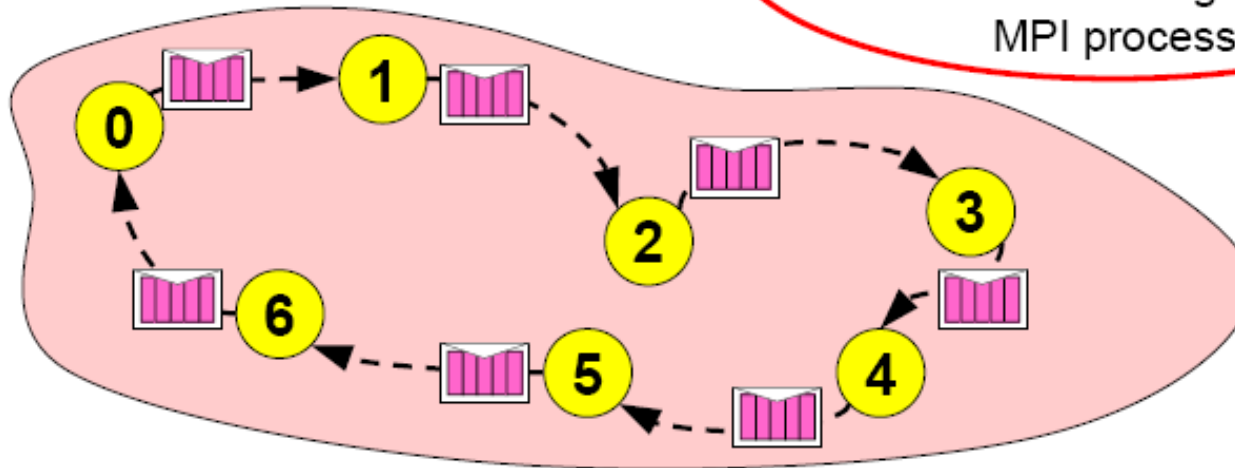
死锁

- Code in each MPI process:

```
MPI_Ssend(..., right_rank, ...)
```

```
MPI_Recv( ..., left_rank, ...)
```

Will block and never return,
because MPI_Recv cannot
be called in the right-hand
MPI process



- Same problem with standard send mode (MPI_Send),
if MPI implementation chooses synchronous protocol ■



捆绑发送接收

- 在上面的**Jacobi**迭代例子中，每一个进程都要向相邻的进程发送数据，同时从相邻的进程接收数据。
- **MPI**提供了捆绑发送和接收操作，可以在一条**MPI**语句中同时实现向其它进程的数据发送和从其它进程接收数据操作。



捆绑发送接收

- 捆绑发送和接收操作把发送一个消息到一个目的地和从另一个进程接收一个消息合并到一个调用中，源和目的可以是相同的。
- 捆绑发送接收操作虽然在语义上等同于一个发送操作和一个接收操作的结合，但是它可以有效地避免由于单独书写发送或接收操作时，由于次序的错误而造成的死锁。这是因为该操作由通信系统来实现，系统会优化通信次序从而有效地避免不合理的通信次序，最大限度避免死锁的产生。



MPI_Sendrecv

- `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)`

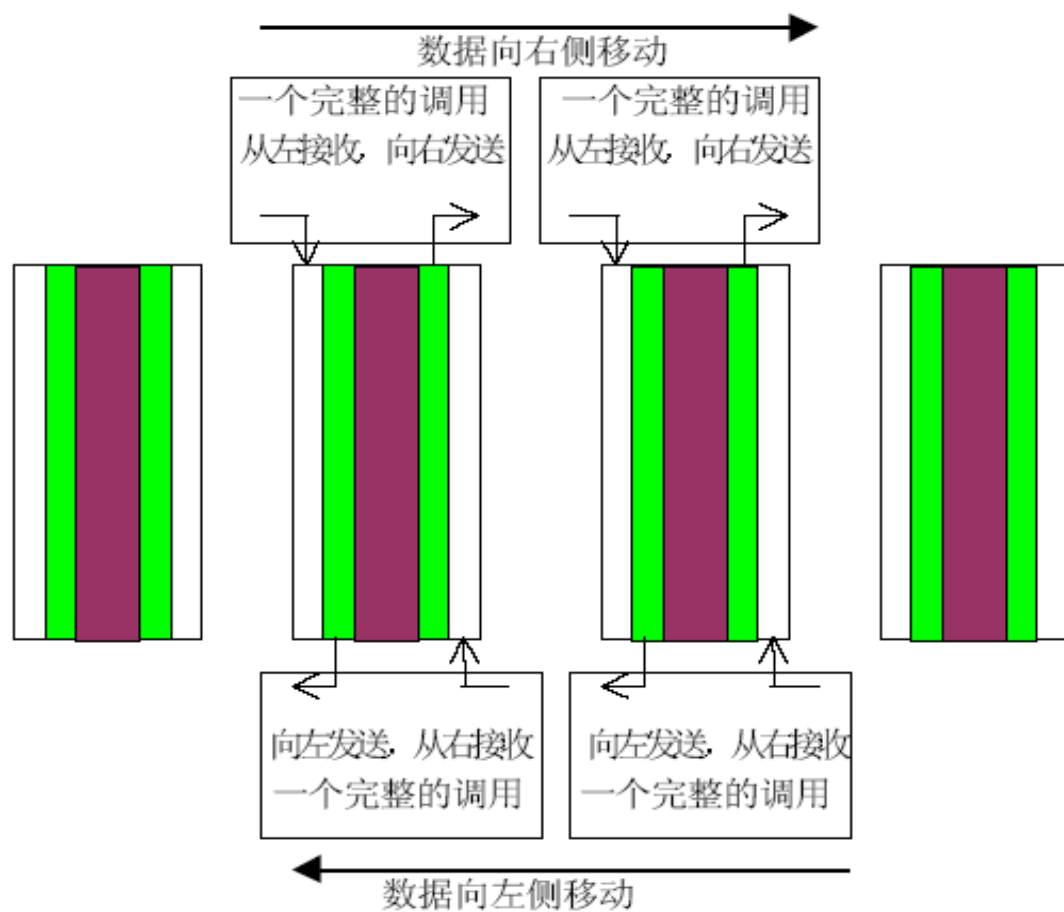


捆绑发送接收

- 捆绑发送接收操作是不对称的，即一个由捆绑发送接收调用发出的消息可以被一个普通接收操作接收，一个捆绑发送接收调用可以接收一个普通发送操作发送的消息。
- 该操作执行一个阻塞的发送和接收，接收和发送使用同一个通信域，但是可能使用不同的标识。
- 发送缓冲区和接收缓冲区必须分开，他们可以是不同的数据长度和不同的数据类型。



用MPI_Sendrecv实现Jacobi迭代





虚拟进程

- 虚拟进程（**MPI_PROC_NULL**）是不存在的假想进程，在**MPI**中的主要作用是充当真实进程通信的目或源。
- 引入虚拟进程的目的是为了在某些情况下编写通信语句的方便。当一个真实进程向一个虚拟进程发送数据或从一个虚拟进程接收数据时，该真实进程会立即正确返回，如同执行了一个空操作。



虚拟进程

- 一个真实进程向虚拟进程MPI_PROC_NULL发送消息时，会立即成功返回。
- 一个真实进程从虚拟进程MPI_PROC_NULL的接收消息时，也会立即成功返回，并且对接收缓冲区没有任何改变。

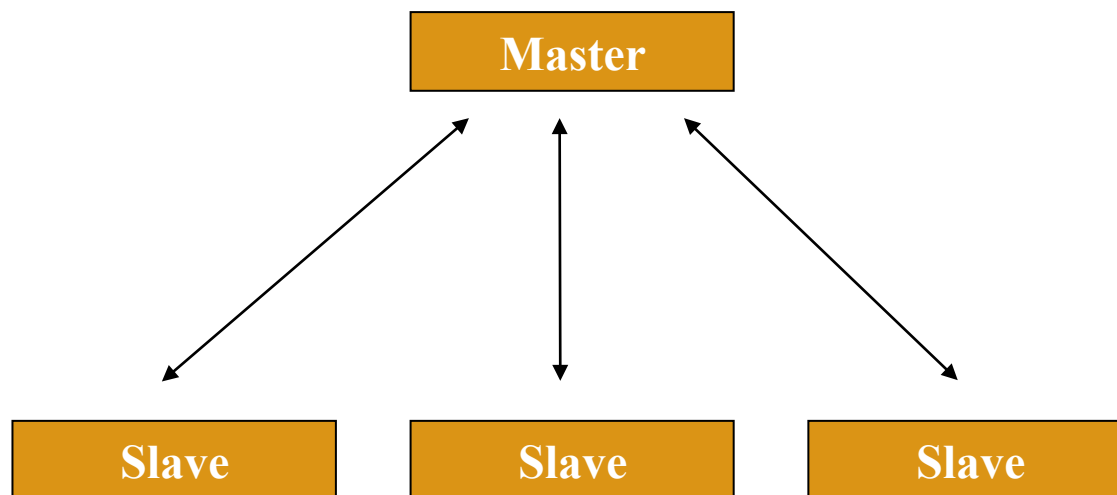


Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- **MPI并行程序的基本模式**
 - 对等模式
 - **主从模式**
- MPI数据及进程
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信

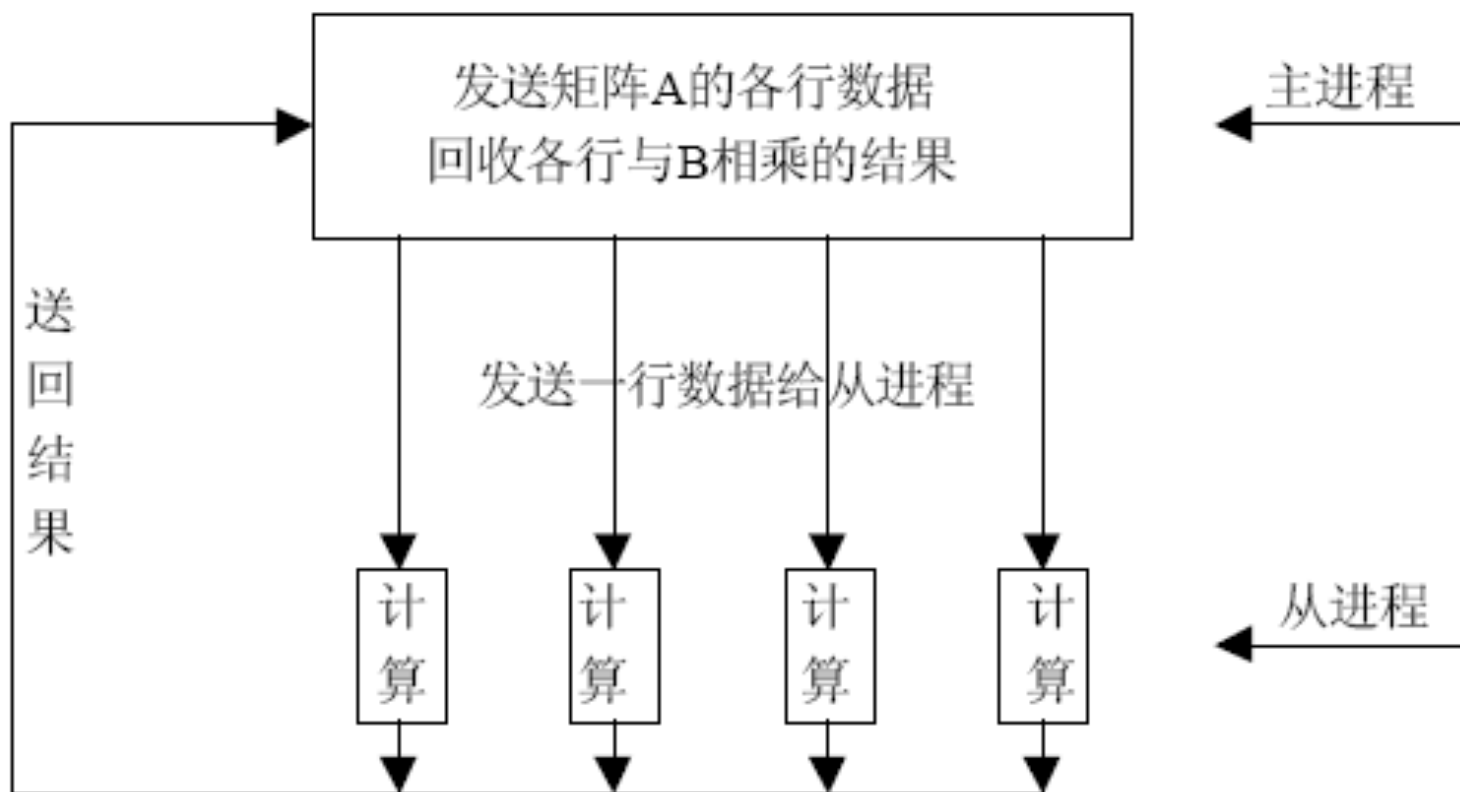


主从模式的MPI程序设计





矩阵向量乘





Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- MPI并行程序的基本模式
 - 对等模式
 - 主从模式
- **MPI数据及进程**
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信

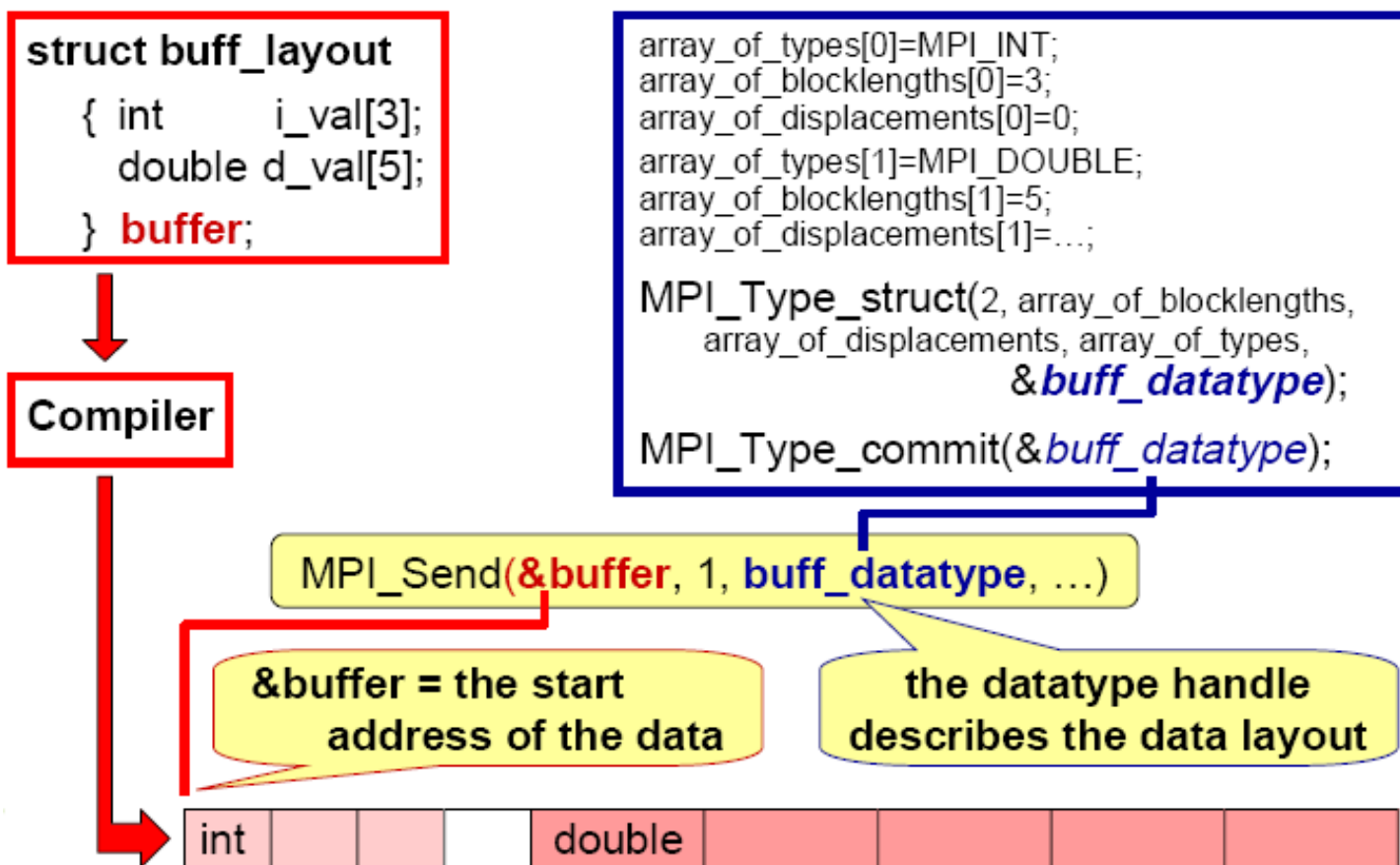


MPI基本数据类型

MPI Datatype	C datatype
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_SHORT	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_BYTE	
MPI_PACKED	

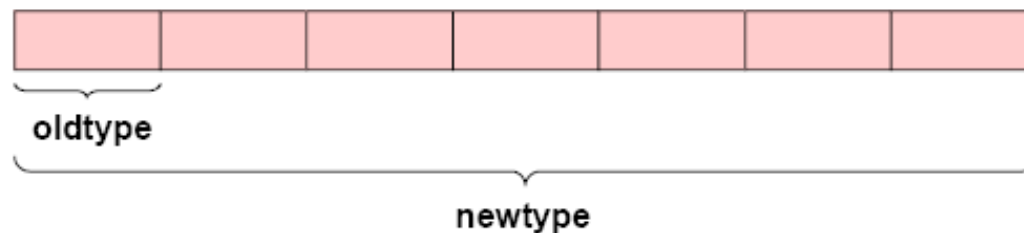


自定义数据类型





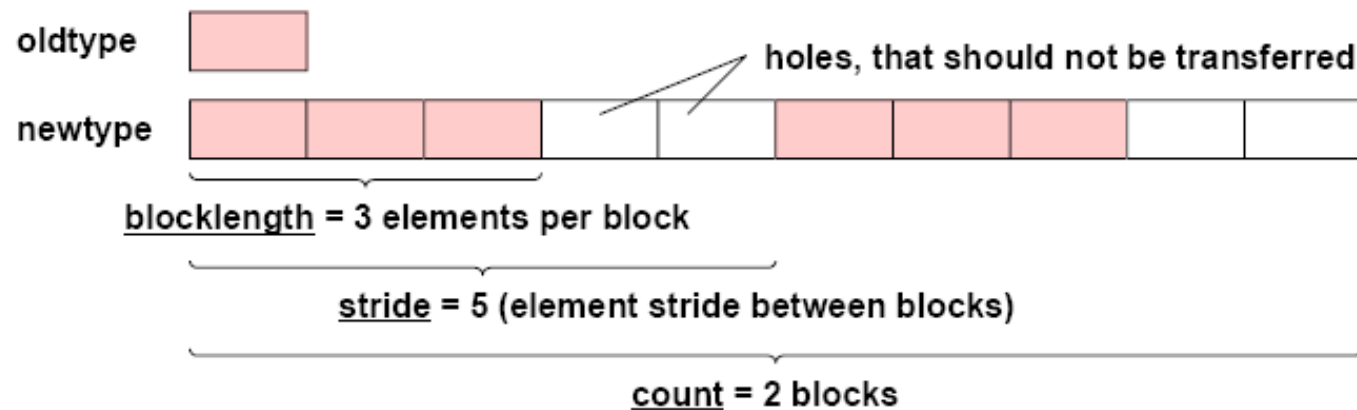
自定义数据类型：连续数据



- C: `int MPI_Type_contiguous(int count, MPI_Datatype oldtype, MPI_Datatype *newtype)`
- Fortran: `MPI_TYPE_CONTIGUOUS(COUNT, OLDTYPE, NEWTYPE, IERROR)`
`INTEGER COUNT, OLDTYPE`
`INTEGER NEWTYPE, IERROR`



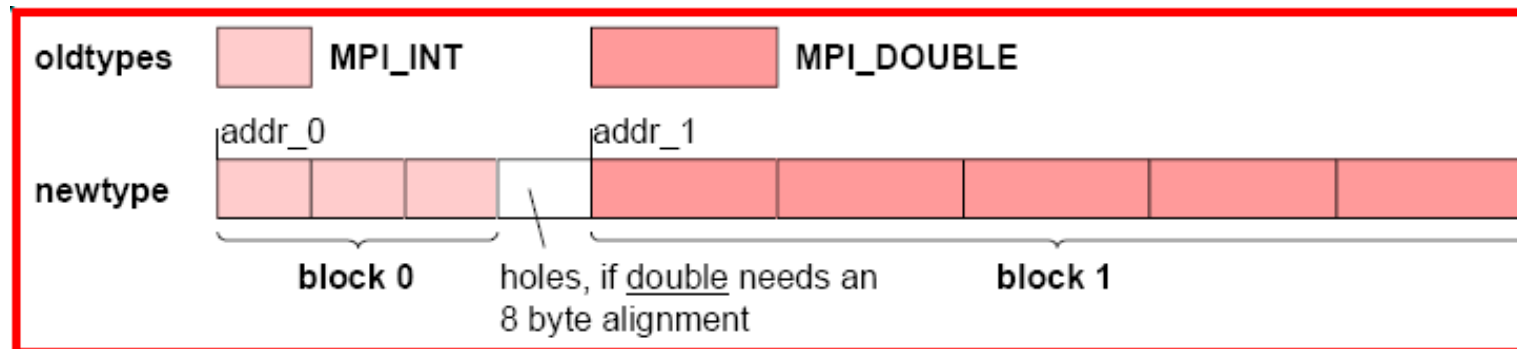
自定义数据类型： 向量



- C: `int MPI_Type_vector(int count, int blocklength, int stride, MPI_Datatype oldtype, MPI_Datatype *newtype)`
- Fortran: `MPI_TYPE_VECTOR(COUNT, BLOCKLENGTH, STRIDE, OLDTYPE, NEWTYPE, IERROR)`
`INTEGER COUNT, BLOCKLENGTH, STRIDE`
`INTEGER OLDTYPE, NEWTYPE, IERROR`



自定义数据类型：结构体



- C: `int MPI_Type_struct(int count, int *array_of_blocklengths, MPI_Aint *array_of_displacements, MPI_Datatype *array_of_types, MPI_Datatype *newtype)`
- Fortran: `MPI_TYPE_STRUCT(COUNT, ARRAY_OF_BLOCKLENGTHS, ARRAY_OF_DISPLACEMENTS, ARRAY_OF_TYPES, NEWTYPE, IERROR)`



Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- MPI并行程序的基本模式
 - 对等模式
 - 主从模式
- **MPI数据及进程**
 - 自定义数据类型
 - **虚拟进程拓扑**
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信



虚拟进程拓扑

- 在许多并行应用程序中，进程的线性排列不能充分地反映进程间在逻辑上的通信模型(通常由问题几何和所用的算法决定)。进程经常被排列成二维或三维网格形式的拓扑模型，而且通常用一个图来描述逻辑进程排列。这种逻辑进程排列称为**虚拟拓扑**。



虚拟进程拓扑

- 拓扑是组内通信域上的额外、可选的属性，它不能附加在组间通信域(**inter-communicator**)上。
- 拓扑能够提供一种方便的命名机制，对于有特定拓扑要求的算法使用起来直接、自然而方便。
- 拓扑还可以辅助运行时系统将进程映射到实际的硬件结构之上。

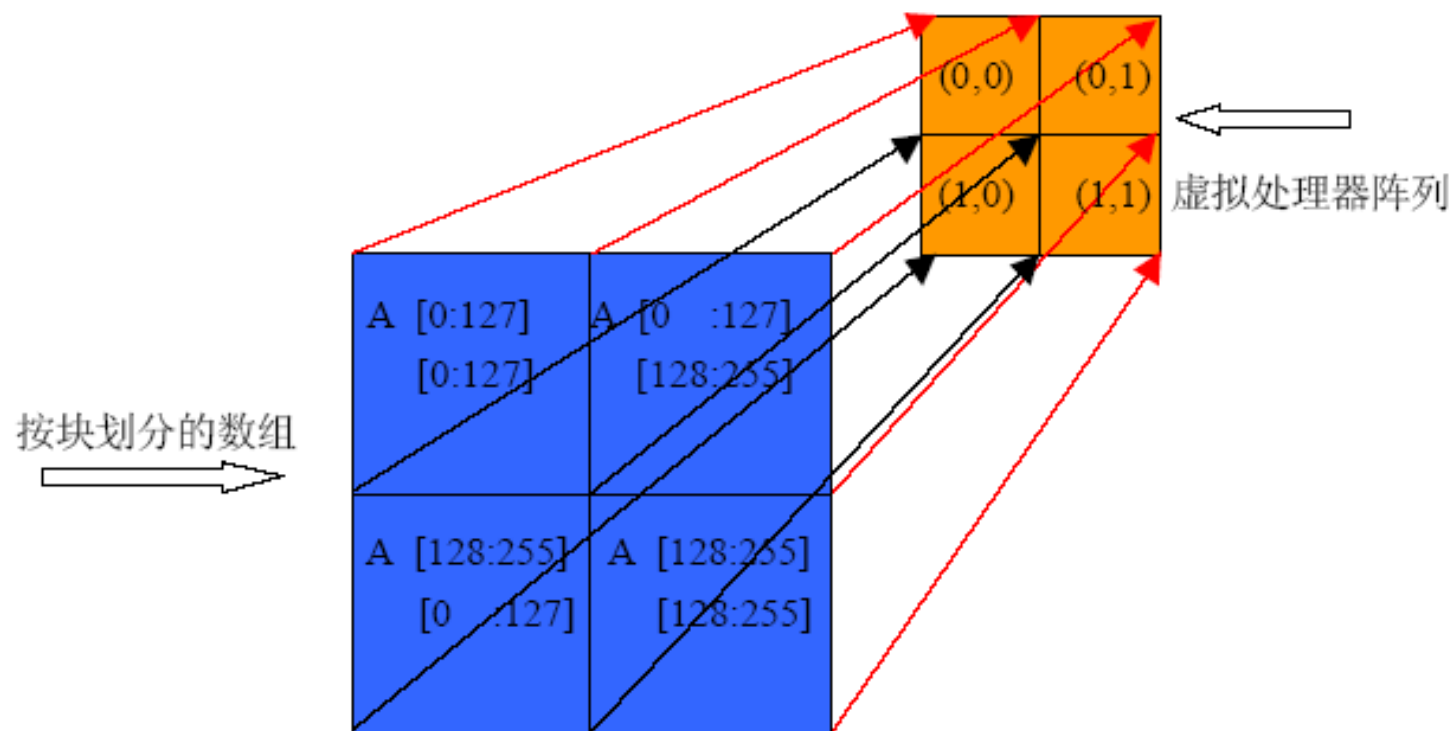


虚拟进程拓扑

- MPI 提供两种拓扑
 - 笛卡儿拓扑
 - 图拓扑
- 分别用来表示简单规则的拓扑和更通用的拓扑。



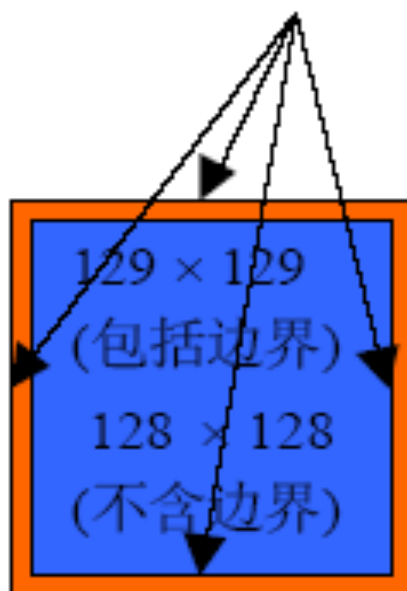
Jacobi迭代





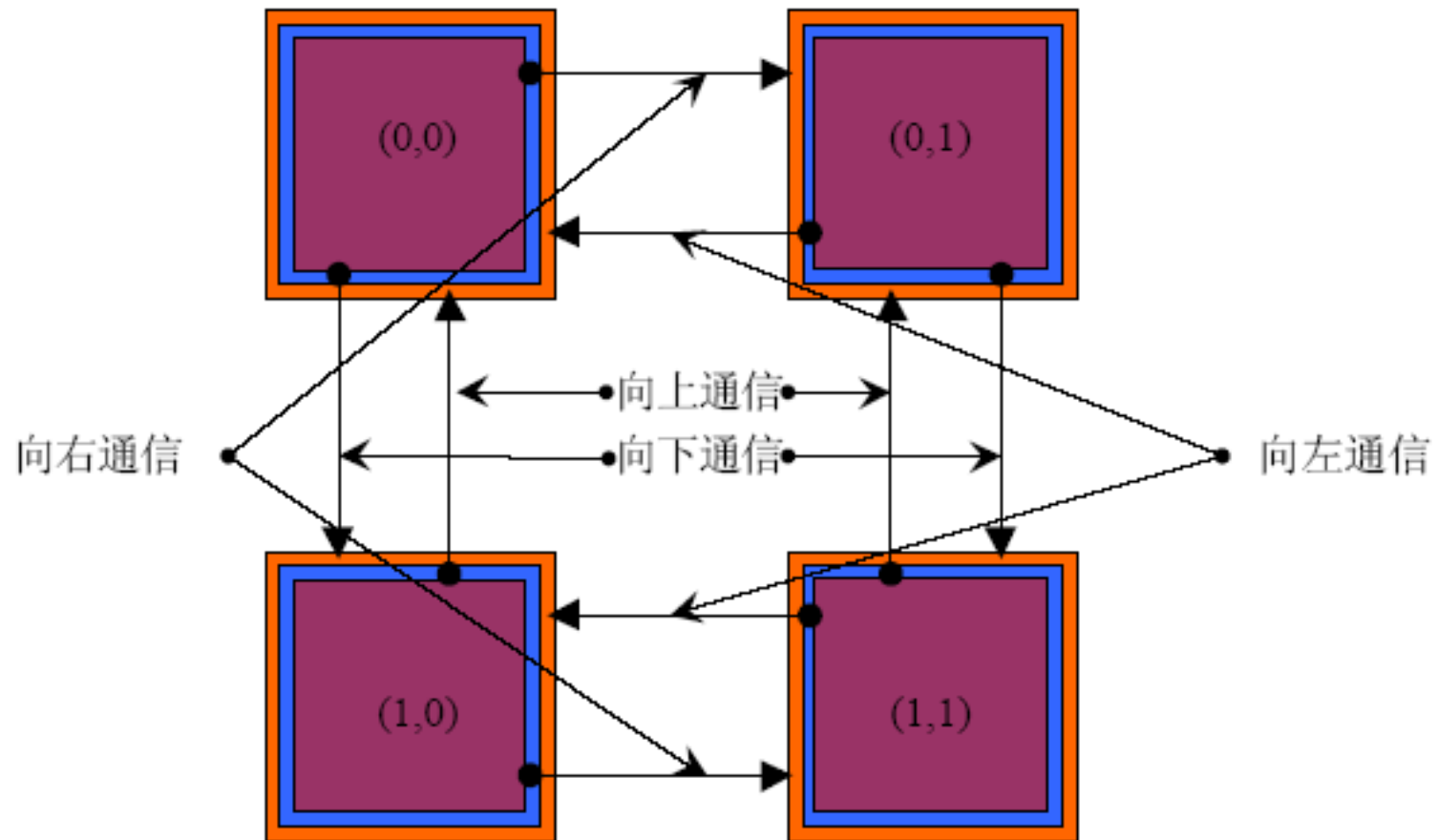
Jacobi迭代

预留出来用来存放通信得到数据的边界





Jacobi迭代





Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- MPI并行程序的基本模式
 - 对等模式
 - 主从模式
- MPI数据及进程
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - **MPI通信模式**
 - 非阻塞通信



MPI通信模式

- 标准通信模式
- 缓存通信模式
- 同步通信模式
- 就绪通信模式



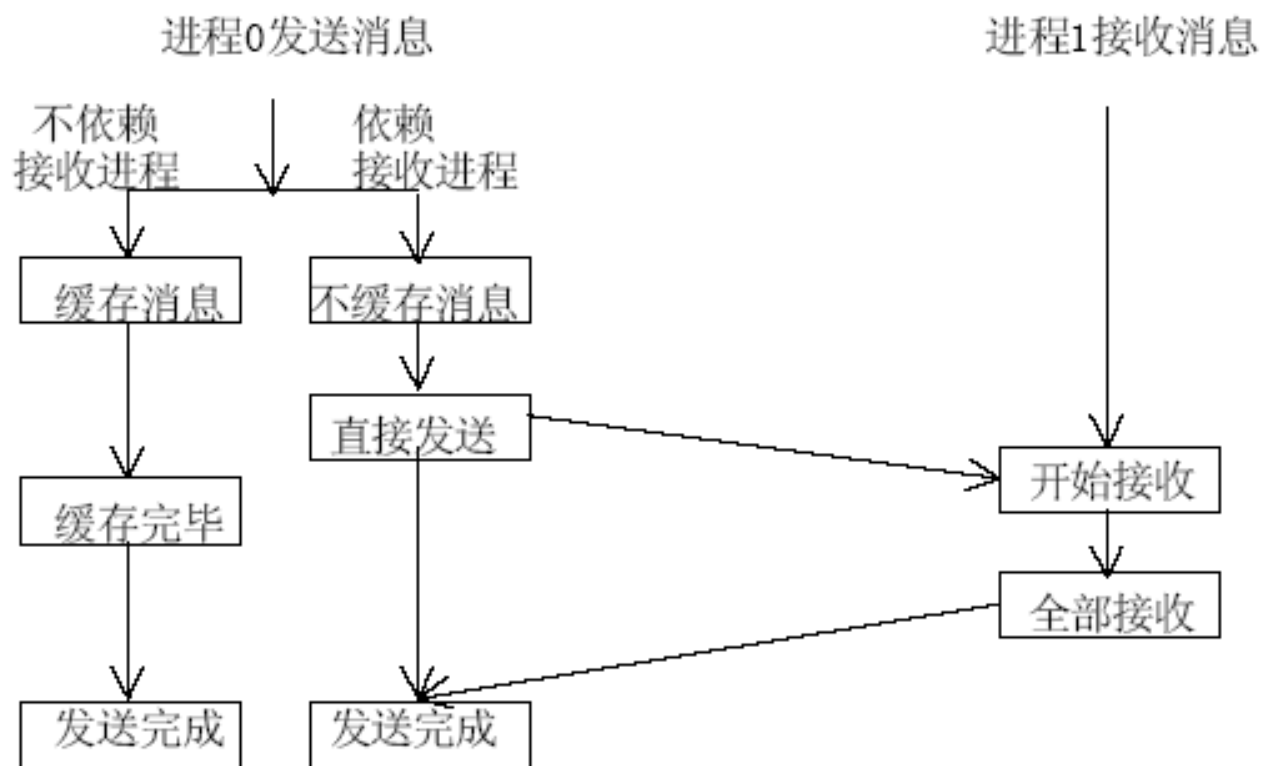
标准通信模式

■ MPI_Send

- 在MPI采用标准通信模式时，是否对发送的数据进行缓存是由MPI自身决定的，而不是由并行程序员来控制。
- 如果MPI决定缓存将要发出的数据，发送操作不管接收操作是否执行，都可以进行，而且发送操作可以正确返回，而不要求接收操作收到发送的数据。



标准通信模式



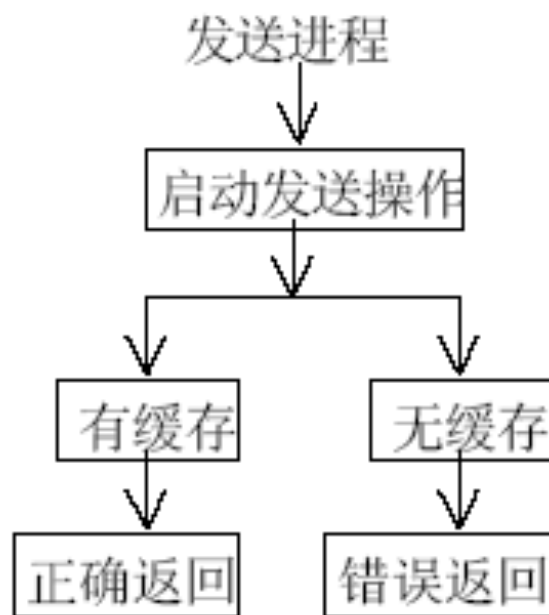


缓存通信模式

- **MPI_Bsend**
 - 由用户直接对通信缓冲区进行申请、使用和释放。
 - 因此缓存模式下对通信缓冲区的合理与正确使用是由程序设计人员自己保证的。
- **MPI_BSEND**的各个参数的含义和**MPI_SEND**的完全相同，不同之处仅表现在通信时是使用标准的系统提供的缓冲区还是用户自己提供的缓冲区。



缓存通信模式





缓存通信模式

- **MPI_BUFFER_ATTACH**将大小为**size**的缓冲区递交给MPI， 这样该缓冲区就可以作为缓存发送时的缓存来使用。
- **MPI_BUFFER_DETACH**将提交的大小为**size**的缓冲区**buffer**收回。该调用是阻塞调用它一直等到使用该缓存的消息发送完成后才返回，这一调用返回后用户可以重新使用该缓冲区，或者将这一缓冲区释放。



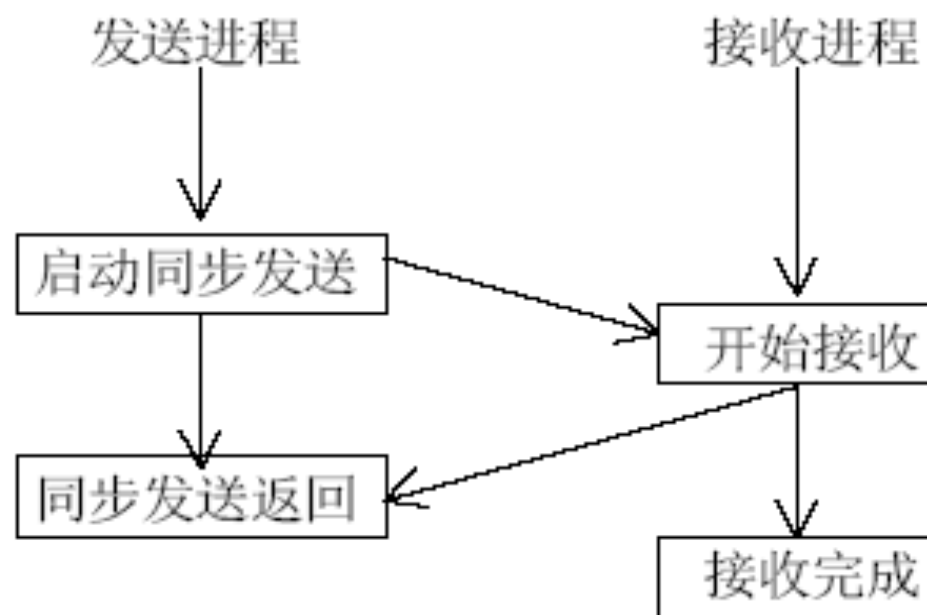
同步通信模式

■ MPI_Ssend

- 同步通信模式的开始不依赖于接收进程相应的接收操作是否已经启动，但是**同步发送**必须等到相应的接收进程开始后可以正确返回。
- 因此，同步发送返回后，意味着发送缓冲区中的数据已经全部被系统缓冲区缓存并且已经开始发送。
- 这样，当同步发送返回后发送缓冲区可以被释放或重新使用。



同步通信模式





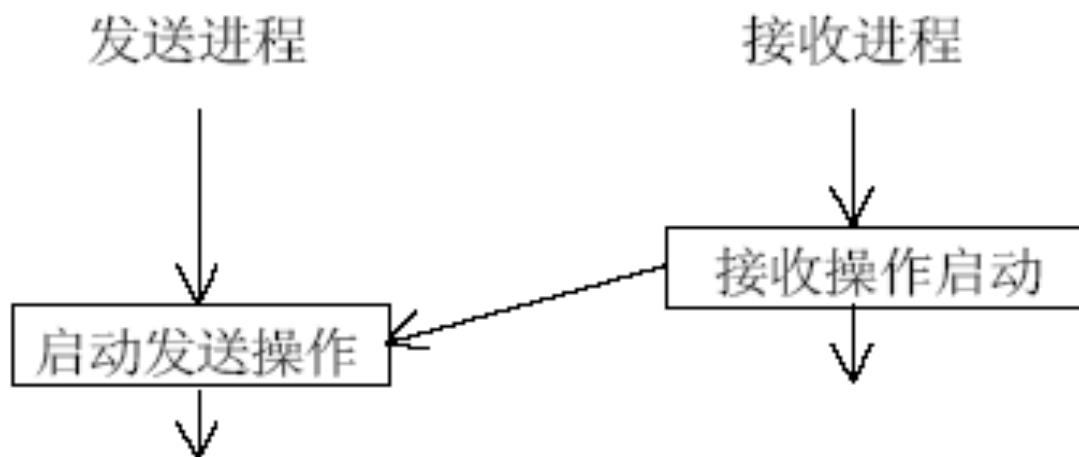
就绪通信模式

■ MPI_Rsend

- 在就绪通信模式中，只有当接收进程的接收操作已经启动时，才可以在发送进程启动发送操作，否则，当发送操作启动而相应的接收还没有启动时，发送操作将出错。
- 对于非阻塞发送操作的正确返回，并不意味着发送已完成；但对于阻塞发送的正确返回，则发送缓冲区可以重复使用。



就绪通信模式





通信模式小结

- 标准通信模式之外的其它通信模式，是**MPI**提供给程序员附加的并程序序通信手段。
- 它是当程序员认为标准通信模式不能满足要求，或者不能很好地满足给定的要求时，所采取的措施。
- 它要求程序员对程序和通信过程有更准确更深刻地理解，在此基础上，根据不同的需要，使用不同的通信模式，往往可以达到优化和提高效率的目的。

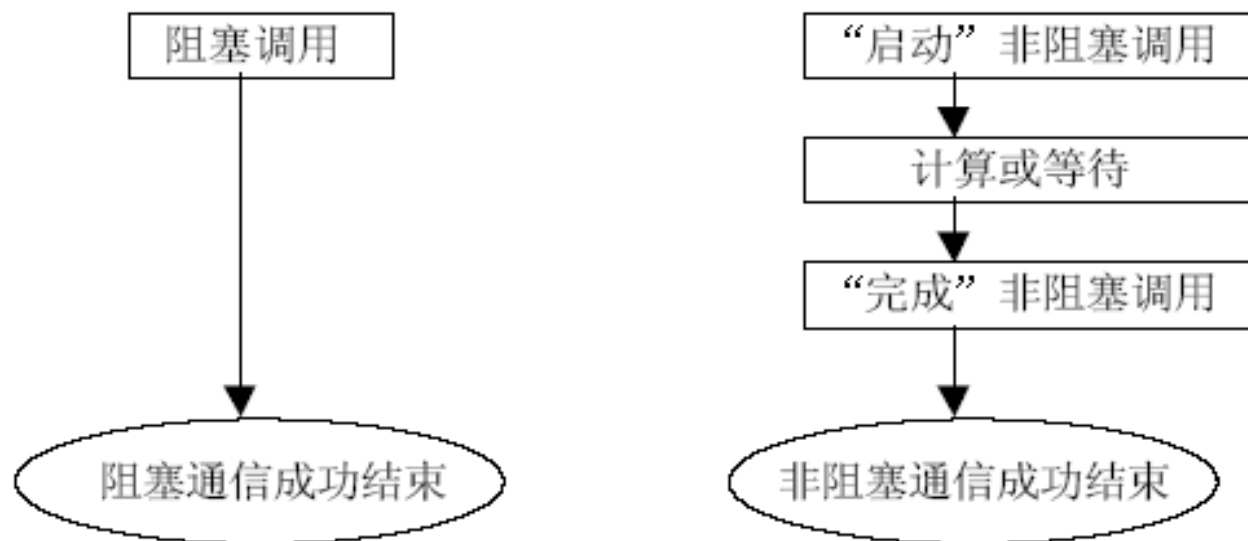


Outline

- MPI概述
 - MPI基本调用
 - 点到点通信与组通信
- MPI并行程序的基本模式
 - 对等模式
 - 主从模式
- MPI数据及进程
 - 自定义数据类型
 - 虚拟进程拓扑
- 消息传递方式
 - MPI通信模式
 - 非阻塞通信



非阻塞通信





非阻塞MPI通信模式

通信模式		发送	接收
标准通信模式		MPI_ISEND	MPI_Irecv
缓存通信模式		MPI_IBSEND	
同步通信模式		MPI_ISSEND	
就绪通信模式		MPI_IRSEND	
重复 非阻塞通信	标准通信模式	MPI_SEND_INIT	MPI_RECV_INIT
	缓存通信模式	MPI_BSEND_INIT	
	同步通信模式	MPI_SSEND_INIT	
	就绪通信模式	MPI_RSEND_INIT	

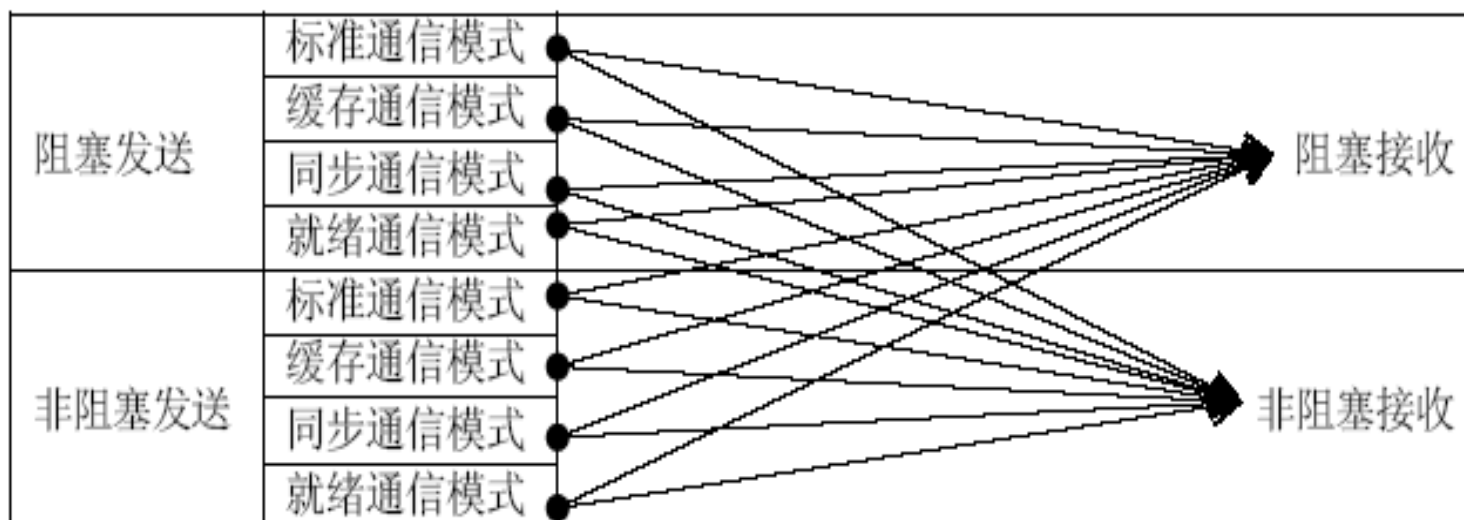


非阻塞通信的完成与检测

非阻塞通信的数量	检测	完成
一个非阻塞通信	MPI_TEST	MPI_WAIT
任意一个非阻塞通信	MPI_TESTANY	MPI_WAITANY
一到多个非阻塞通信	MPI_TESTSOME	MPI_WAITSOME
所有非阻塞通信	MPI_TESTALL	MPI_WAITALL

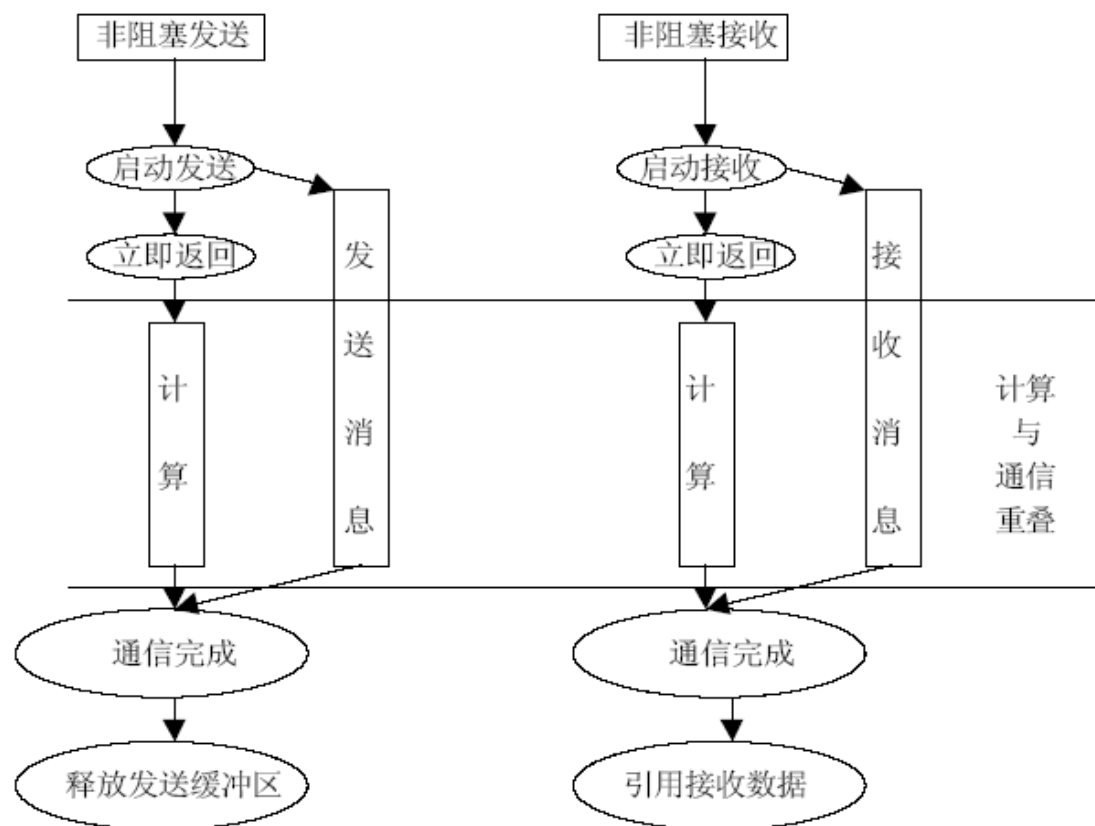


不同类型的发送与接收的匹配





非阻塞标准发送和接收





MPI_Isend

```
■ int MPI_Isend(void* buf,  
                int count,  
                MPI_Datatype datatype,  
                int dest,  
                int tag,  
                MPI_Comm comm,  
                MPI_Request *request)
```



非阻塞通信与其它三种通信模式的组合

- 对于阻塞通信的四种消息通信模式：标准通信模式，缓存通信模式，同步通信模式和接收就绪通信模式，非阻塞通信也具有相应的四种不同模式。
- **MPI**使用与阻塞通信一样的命名约定，前缀**B**、**S**、**R**分别表示缓存通信模式、同步通信模式和就绪通信模式。
- 前缀**I(immediate)**表示这个调用是非阻塞的。



非阻塞通信的完成

- 对于非阻塞通信，通信调用的返回并不意味着通信的完成，因此需要专门的通信语句来完成或检查该非阻塞通信。
- 不管非阻塞通信是什么样的形式，对于完成调用是不加区分的。
- 当非阻塞完成调用结束后，就可以保证该非阻塞通信已经正确完成了。



单个非阻塞通信的完成

- `int MPI_Wait(MPI_Request *request, MPI_Status *status)`
- `int MPI_Test(MPI_Request*request, int *flag, MPI_Status *status)`



多个非阻塞通信的完成（1）

- `int MPI_Waitany(int count,
MPI_Request *array_of_requests,
int *index,
MPI_Status *status)`
- `MPI_WAITANY`返回后`index=i`，即`MPI_WAITANY`完成的是非阻塞通信对象表中的第`i`个对象对应的非阻塞通信，则其效果等价于调用了`MPI_WAIT(array_of_requests[i],status)`



多个非阻塞通信的完成（2）

- `int MPI_Waitall(int count,
MPI_Request *array_of_requests,
MPI_Status *array_of_statuses)`



多个非阻塞通信的完成（3）

- `int MPI_Waitsome(int incount,
MPI_Request *array_of_request,
int *outcount,
int *array_of_indices,
MPI_Status *array_of_statuses)`



测试非阻塞通信对象

- `int MPI_Test` (`MPI_Request *request`, `int *flag`,
`MPI_Status *status`)
- `int MPI_Testany` (`int count`, `MPI_Request *array_of_requests`, `int *index`,`int *flag`, `MPI_Status *status`)
- `int MPI_Testall` (`int count`, `MPI_Request *array_of_requests`, `int *flag`,`MPI_Status *array_of_statuses`)
- `int MPI_Testsome` (`int incount`,`MPI_Request *array_of_request`,
`int *outcount`, `int *array_of_indices`, `MPI_Status *array_of_statuses`)



非阻塞通信对象

- 非阻塞通信的取消
 - int MPI_Cancel(MPI_Request *request)
- 非阻塞通信对象的释放
 - int MPI_Request_free(MPI_Request * request)
- 消息到达的检查
 - int MPI_Probe(int source,int tag,MPI_Comm comm,MPI_Status *status)
 - int MPI_Iprobe(int source,int tag,MPI_Comm comm,int *flag, MPI_Status *status)