

并行计算

(Parallel Computing)

分布式内存编程 - MPI

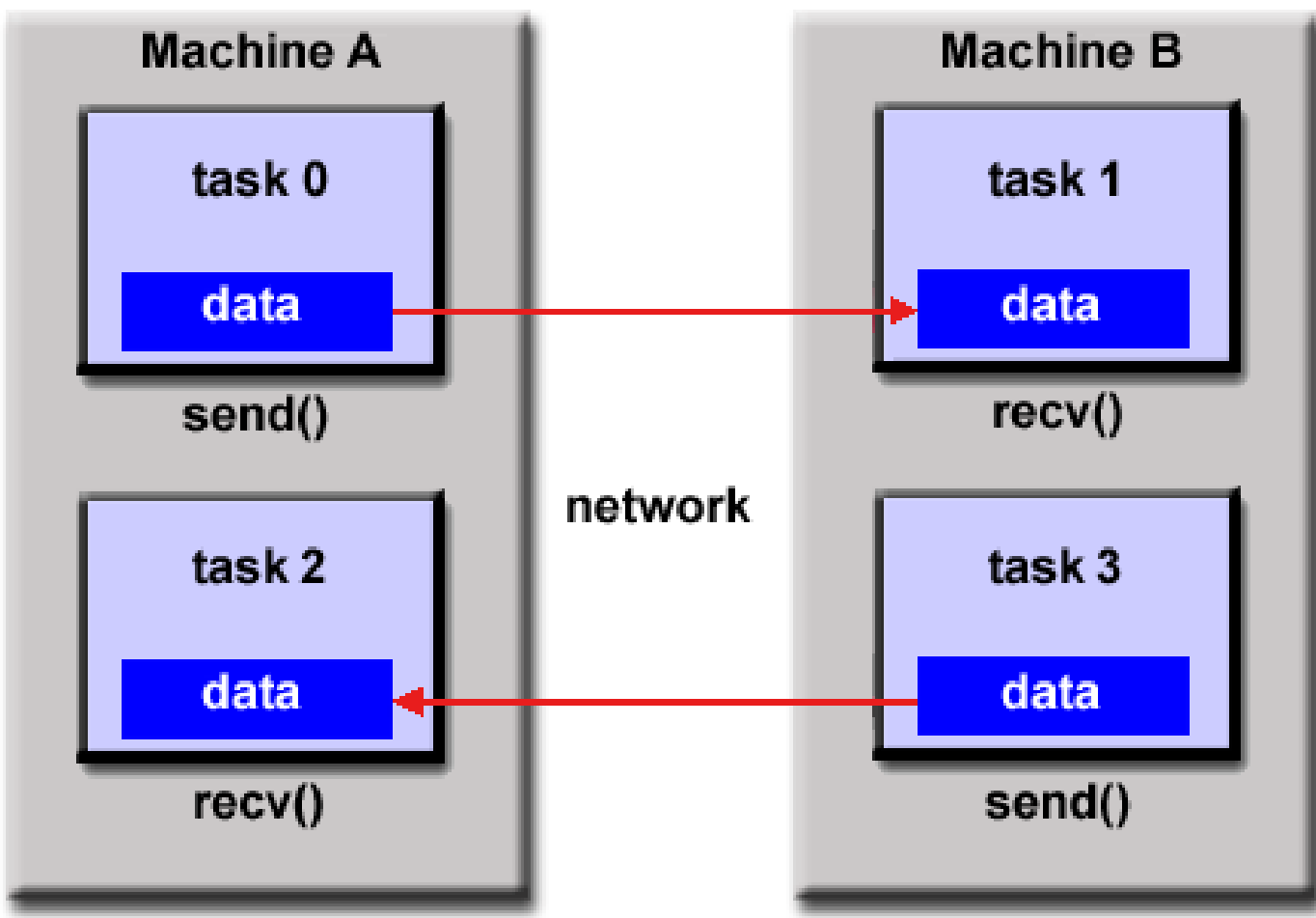
学习内容：

- 第一个 MPI 程序
- 梯形法则（Trapezoidal Rule）求面积
- 集体通信（Collective communication）

分布式内存编程 - MPI

学习内容：

- MPI 中的派生数据类型（derived datatypes）
- MPI 程序的性能评价
- 并行排序
- MPI 程序的安全性



MPI include file

Declarations, prototypes, etc.

Program Begins

·
·
·

Serial code

Initialize MPI environment

Parallel code begins

·
·
·

Do work & make message passing calls

·
·
·

Terminate MPI environment

Parallel code ends

·
·
·

Serial code

Program Ends

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

int main (int argc, char *argv[])
{
    int numtasks, rank, dest, source, rc, count, tag=1;
    char inmsg, outmsg='x';
    MPI_Status Stat;

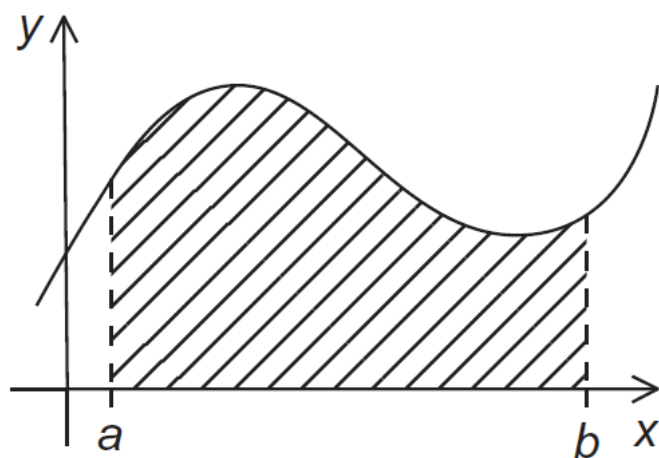
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numtasks);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    if (rank == 0) {
        dest = 1;
        source = 1;
        rc = MPI_Send(&outmsg, 1, MPI_CHAR, dest, tag, MPI_COMM_WORLD);
        rc = MPI_Recv(&inmsg, 1, MPI_CHAR, source, tag, MPI_COMM_WORLD, &Stat);
    }

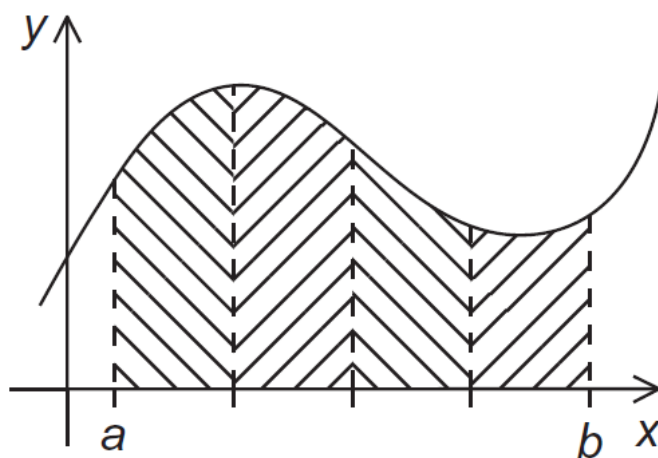
    else if (rank == 1) {
        dest = 0;
        source = 0;
        rc = MPI_Recv(&inmsg, 1, MPI_CHAR, source, tag, MPI_COMM_WORLD, &Stat);
        rc = MPI_Send(&outmsg, 1, MPI_CHAR, dest, tag, MPI_COMM_WORLD);
    }

    MPI_Finalize();
}
```

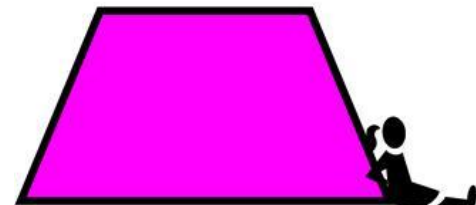
2. 梯形法则 (Trapezoidal Rule) 数值积分



(a)



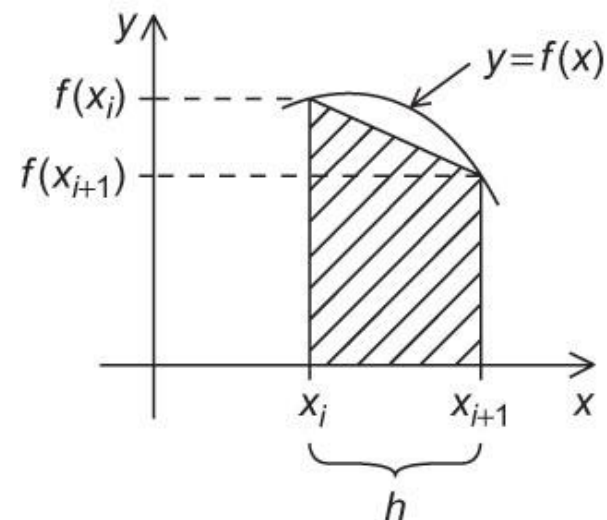
(b)



2. 梯形法则 (Trapezoidal Rule) 数值积分

$$\text{Area of one trapezoid} = \frac{h}{2}[f(x_i) + f(x_{i+1})]$$

$$h = \frac{b-a}{n}$$



$$x_0 = a, x_1 = a + h, x_2 = a + 2h, \dots, x_{n-1} = a + (n-1)h, x_n = b$$

$$\text{Sum of trapezoid areas} = h[f(x_0)/2 + f(x_1) + f(x_2) + \dots + f(x_{n-1}) + f(x_n)/2]$$

2. 梯形法则 (Trapezoidal Rule) 数值积分

● 串程序的伪代码

Sum of trapezoid areas = $h[f(x_0)/2 + f(x_1) + f(x_2) + \cdots + f(x_{n-1}) + f(x_n)/2]$

```
/* Input:  a, b, n */  
h = (b-a)/n;  
approx = (f(a) + f(b))/2.0;  
for (i = 1; i <= n-1; i++) {  
    x_i = a + i*h;  
    approx += f(x_i);  
}  
approx = h*approx;
```

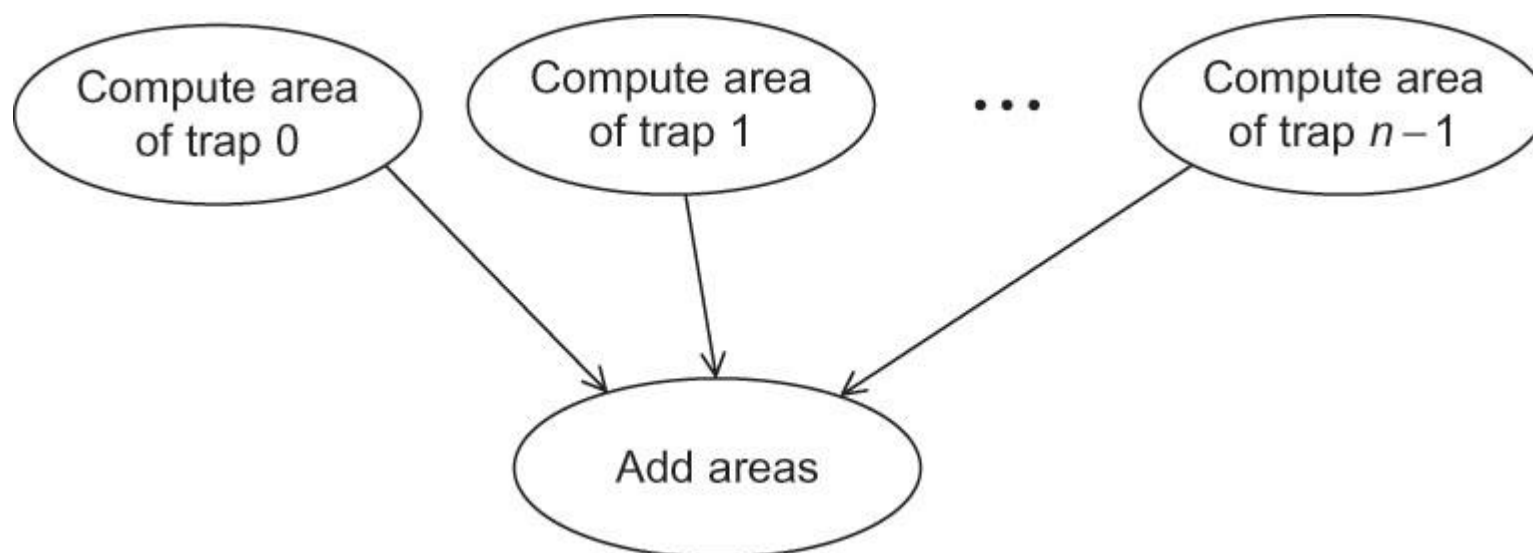

2.梯形法则 (Trapezoidal Rule) 数值积分

●并行化

- 将问题解决方案划分为任务
- 确认任务之间的通信
- 聚合任务为组合任务
- 将确定的任务映射到 core

2. 梯形法则 (Trapezoidal Rule) 数值积分

● 并行化



2. 梯形法则 (Trapezoidal Rule) 数值积分

● 并行化

```
1   Get a, b, n;
2   h = (b-a)/n;
3   local_n = n/comm_sz;
4   local_a = a + my_rank*local_n*h;
5   local_b = local_a + local_n*h;
6   local_integral = Trap(local_a, local_b, local_n, h);
7   if (my_rank != 0)
8       Send local_integral to process 0;
9   else /* my_rank == 0 */
10       total_integral = local_integral;
11       for (proc = 1; proc < comm_sz; proc++) {
12           Receive local_integral from proc;
13           total_integral += local_integral;
14       }
15   }
16   if (my_rank == 0)
17       print result;
```

2. 梯形法则 (Trapezoidal Rule) 数值积分

● 并行化

```
1  int main(void) {
2      int my_rank, comm_sz, n = 1024, local_n;
3      double a = 0.0, b = 3.0, h, local_a, local_b;
4      double local_int, total_int;
5      int source;
6
7      MPI_Init(NULL, NULL);
8      MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
9      MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);
10
11     h = (b-a)/n;          /* h is the same for all processes */
12     local_n = n/comm_sz; /* So is the number of trapezoids */
13
14     local_a = a + my_rank*local_n*h;
15     local_b = local_a + local_n*h;
16     local_int = Trap(local_a, local_b, local_n, h);
17
18     if (my_rank != 0) {
19         MPI_Send(&local_int, 1, MPI_DOUBLE, 0, 0,
20                 MPI_COMM_WORLD);
```

2. 梯形法则 (Trapezoidal Rule) 数值积分

● 并行化

```
21     } else {  
22         total_int = local_int;  
23         for (source = 1; source < comm_sz; source++) {  
24             MPI_Recv(&local_int, 1, MPI_DOUBLE, source, 0,  
25                     MPI_COMM_WORLD, MPI_STATUS_IGNORE);  
26             total_int += local_int;  
27         }  
28     }  
29  
30     if (my_rank == 0) {  
31         printf("With n = %d trapezoids, our estimate\n", n);  
32         printf("of the integral from %f to %f = %.15e\n",  
33             a, b, total_int);  
34     }  
35     MPI_Finalize();  
36     return 0;  
37 } /*  main  */
```

2. 梯形法则 (Trapezoidal Rule) 数值积分

● 并行化

```
1  double Trap(  
2      double left_endpt  /* in */,  
3      double right_endpt /* in */,  
4      int trap_count    /* in */,  
5      double base_len    /* in */) {  
6      double estimate, x;  
7      int i;  
8  
9      estimate = (f(left_endpt) + f(right_endpt))/2.0;  
10     for (i = 1; i <= trap_count-1; i++) {  
11         x = left_endpt + i*base_len;  
12         estimate += f(x);  
13     }  
14     estimate = estimate*base_len;  
15  
16     return estimate;  
17 } /* Trap */
```

3.处理 I/O

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

int main(void) {
    int my_rank, comm_sz;

    MPI_Init(NULL, NULL);
    MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);
    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);

    printf("Proc %d of %d > Does anyone have a toothpick?\n",
           my_rank, comm_sz);

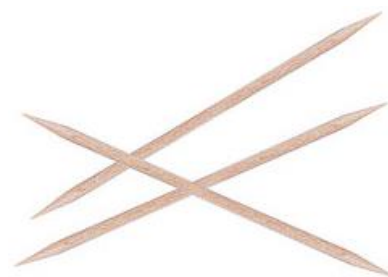
    MPI_Finalize();
    return 0;
} /* main */
```

*Each process just
prints a message.*

3.处理 I/O

```
Proc 0 of 6 > Does anyone have a toothpick?
Proc 1 of 6 > Does anyone have a toothpick?
Proc 2 of 6 > Does anyone have a toothpick?
Proc 4 of 6 > Does anyone have a toothpick?
Proc 3 of 6 > Does anyone have a toothpick?
Proc 5 of 6 > Does anyone have a toothpick?
```

unpredictable output



3.处理 I/O

●输入

- 大多数 MPI 的实现仅允许 MPI_COMM_WORLD 中的进程 0 访问 stdin
- 进程 0 必须读取数据 (scanf) 并发送给其他进程

```
...  
MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);  
MPI_Comm_size(MPI_COMM_WORLD, &comm_sz);  
Get_input(my_rank, comm_sz, &a, &b, &n);  
h = (b - a)/n;  
...
```

3.处理 I/O

●输入

```

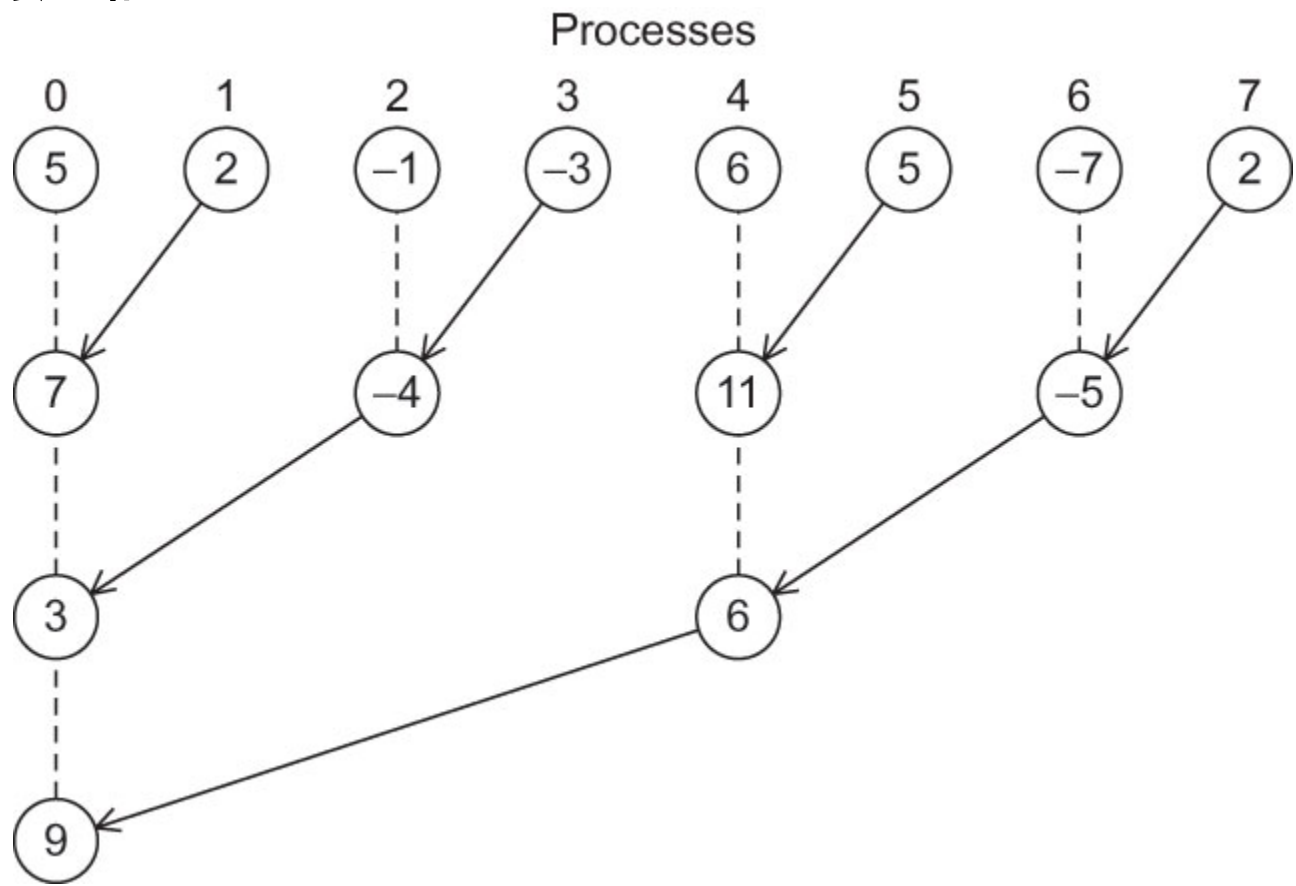
void Get_input(
    int      my_rank    /* in  */,
    int      comm_sz    /* in  */,
    double*  a_p        /* out */,
    double*  b_p        /* out */,
    int*      n_p        /* out */) {
    int dest;

    if (my_rank == 0) {
        printf("Enter a, b, and n\n");
        scanf("%lf %lf %d", a_p, b_p, n_p);
        for (dest = 1; dest < comm_sz; dest++) {
            MPI_Send(a_p, 1, MPI_DOUBLE, dest, 0, MPI_COMM_WORLD);
            MPI_Send(b_p, 1, MPI_DOUBLE, dest, 0, MPI_COMM_WORLD);
            MPI_Send(n_p, 1, MPI_INT, dest, 0, MPI_COMM_WORLD);
        }
    } else { /* my_rank != 0 */
        MPI_Recv(a_p, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD,
            MPI_STATUS_IGNORE);
        MPI_Recv(b_p, 1, MPI_DOUBLE, 0, 0, MPI_COMM_WORLD,
            MPI_STATUS_IGNORE);
        MPI_Recv(n_p, 1, MPI_INT, 0, 0, MPI_COMM_WORLD,
            MPI_STATUS_IGNORE);
    }
} /* Get_input */

```

4.集体通信（Collective communication）

●树结构的通信



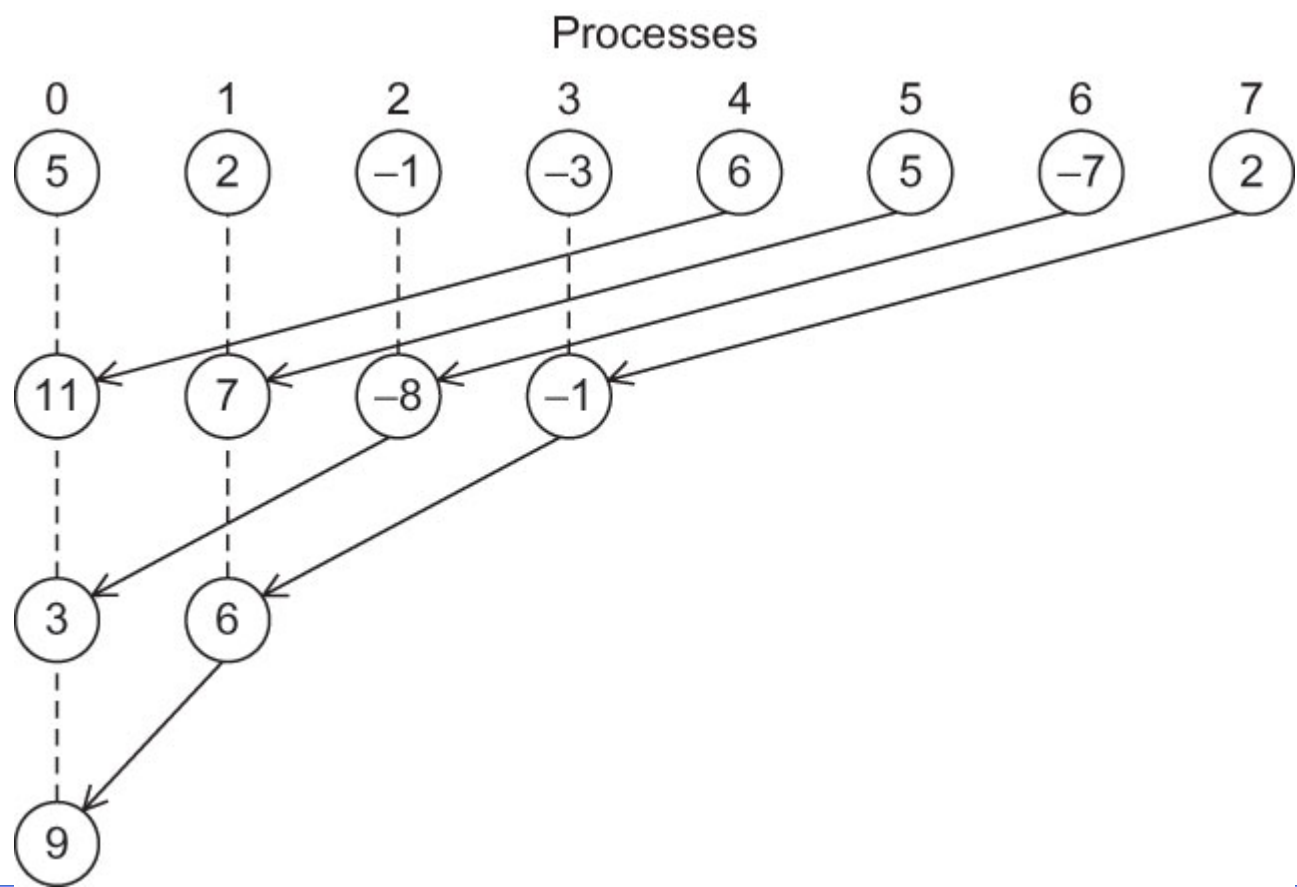
4.集体通信（Collective communication）

●树结构的通信

- a. 进程 1, 3, 5, 7 分别发送数据给进程 0, 2, 4, 6
- b. 进程 0, 2, 4, 6 累加接收到的值
- c. 进程 2 和 6 分别发送累加后的值给进程 0 和 4
- d. 进程 0 和 4 累加接收到的值
- e. 进程 4 发送新的累加值给进程 0
- f. 进程 0 累加接收到的值

4.集体通信（Collective communication）

●树结构的通信



4.集体通信（Collective communication）

●MPI_Reduce

- 实现 global_sum 的方法有很多，如何比较它们哪个最优？
- MPI 提供了 global sum 函数，由 MPI 的实现来完成优化工作
- 不同于 MPI_Send – MPI_Recv，global sum 函数需要多于两个进程之间进行通信，在梯形法则求面积的例子中，包含了 MPI_COMM_WORLD 中所有的进程
- 这种通信称为 collective communication
- MPI_Send – MPI_Recv 称为点对点通信

4.集体通信 (Collective communication)

●MPI_Reduce

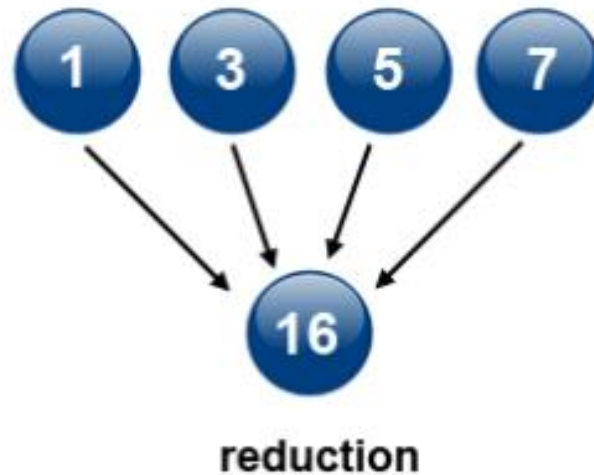
```
int MPI_Reduce(  
    void*      input_data_p    /* in */,  
    void*      output_data_p   /* out */,  
    int        count           /* in */,  
    MPI_Datatype datatype       /* in */,  
    MPI_Op      operator       /* in */,  
    int        dest_process    /* in */,  
    MPI_Comm     comm          /* in */);
```

```
MPI_Reduce(&local_int, &total_int, ①, MPI_DOUBLE, MPI_SUM, 0,  
          MPI_COMM_WORLD);
```

```
double local_x[N], sum[N];  
.  
.  
.  
MPI_Reduce(local_x, sum, N, MPI_DOUBLE, MPI_SUM, 0,  
          MPI_COMM_WORLD);
```

4.集体通信 (Collective communication)

●MPI_Reduce



4.集体通信（Collective communication）

●MPI_Reduce

Operation Value	Meaning
MPI_MAX	Maximum
MPI_MIN	Minimum
MPI_SUM	Sum
MPI_PROD	Product
MPI_LAND	Logical and
MPI_BAND	Bitwise and
MPI_LOR	Logical or
MPI_BOR	Bitwise or
MPI_LXOR	Logical exclusive or
MPI_BXOR	Bitwise exclusive or
MPI_MAXLOC	Maximum and location of maximum
MPI_MINLOC	Minimum and location of minimum

4.集体通信（Collective communication）

●集体通信 vs. 点对点通信

➤ 通信器中的所有进程必须调用同一 collective 通信函数

- 例如：一个进程中调用 MPI_Reduce，另一个进程中调用 MPI_Recv，将会产生错误

➤ 每个进程传递给集体通信的参数必须“兼容”

- 例如：如果一个进程传递 0 作为 dest_process 而另一个进程传递 1，则 MPI_Reduce 的结果是错误的

➤ output_data_p 参数只在 dest_process 中使用。其他进程也需传递实参给 output_data_p，即使只传递 NULL

4.集体通信（Collective communication）

●集体通信 vs. 点对点通信

- 点对点通信通过 tag 和通信器匹配
- 集体通信不使用 tag
- 它们仅通过通信器和调用的顺序来匹配

4.集体通信（Collective communication）

●集体通信 vs. 点对点通信

Time	Process 0	Process 1	Process 2
0	a = 1; c = 2	a = 1; c = 2	a = 1; c = 2
1	MPI_Reduce (&a, &b, ...)	MPI_Reduce (&c, &d, ...)	MPI_Reduce (&a, &b, ...)
2	MPI_Reduce (&c, &d, ...)	MPI_Reduce (&a, &b, ...)	MPI_Reduce (&c, &d, ...)

Multiple calls to MPI_Reduce

4.集体通信（Collective communication）

●集体通信 vs. 点对点通信

- 每个进程以运算符 MPI_SUM 调用 MPI_Reduce，目的进程为进程 0
- *b = 3, d = 6?*

Time	Process 0	Process 1	Process 2
0	a = 1; c = 2	a = 1; c = 2	a = 1; c = 2
1	MPI_Reduce (&a, &b, ...)	MPI_Reduce (&c, &d, ...)	MPI_Reduce (&a, &b, ...)
2	MPI_Reduce (&c, &d, ...)	MPI_Reduce (&a, &b, ...)	MPI_Reduce (&c, &d, ...)

4.集体通信（Collective communication）

●集体通信 vs. 点对点通信

- 内存位置的名称与匹配 MPI_Reduce 调用无关
- 调用的顺序决定了 MPI_Reduce 的匹配
- 因此， $b = 1+2+1 = 4$ ， $d = 2+1+2 = 5$

Time	Process 0	Process 1	Process 2
0	<code>a = 1; c = 2</code>	<code>a = 1; c = 2</code>	<code>a = 1; c = 2</code>
1	<code>MPI_Reduce (&a, &b, ...)</code>	<code>MPI_Reduce (&c, &d, ...)</code>	<code>MPI_Reduce (&a, &b, ...)</code>
2	<code>MPI_Reduce (&c, &d, ...)</code>	<code>MPI_Reduce (&a, &b, ...)</code>	<code>MPI_Reduce (&c, &d, ...)</code>

4.集体通信（Collective communication）

●集体通信 vs. 点对点通信

- 输入和输出使用同一缓冲区，为非法调用，结果不可预测

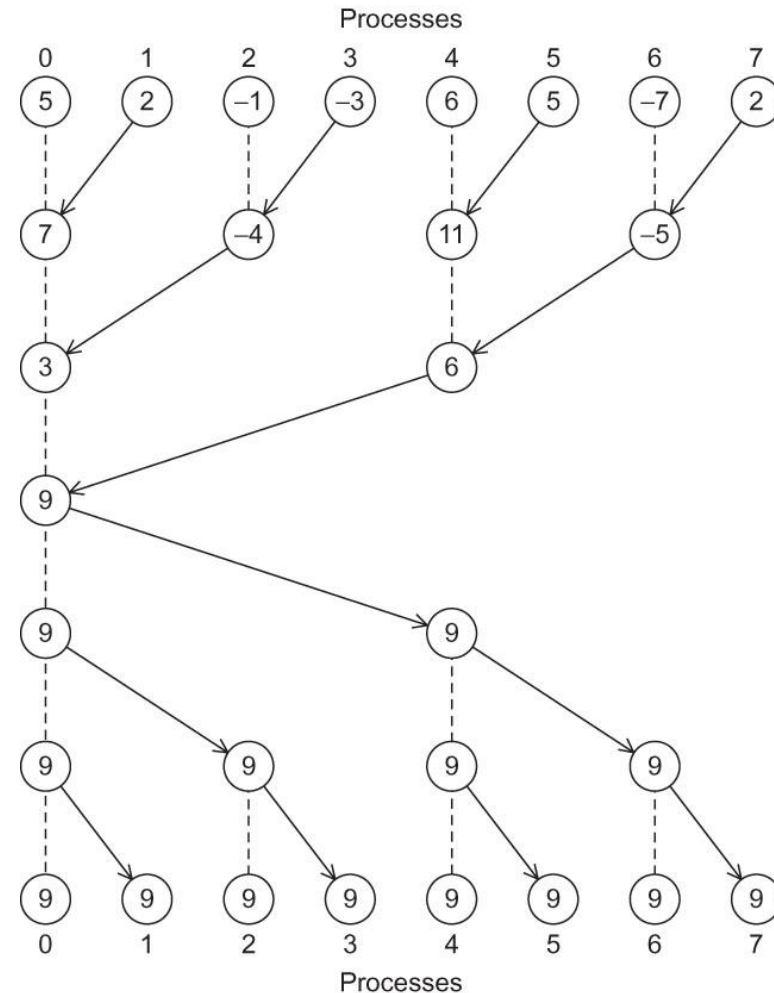
```
MPI_Reduce(&x, &x, 1, MPI_DOUBLE, MPI_SUM, 0, comm);
```

4.集体通信（Collective communication）

●MPI_Allreduce

- 当所有的进程都需要计算结果（如：sum）以进行后续的计算

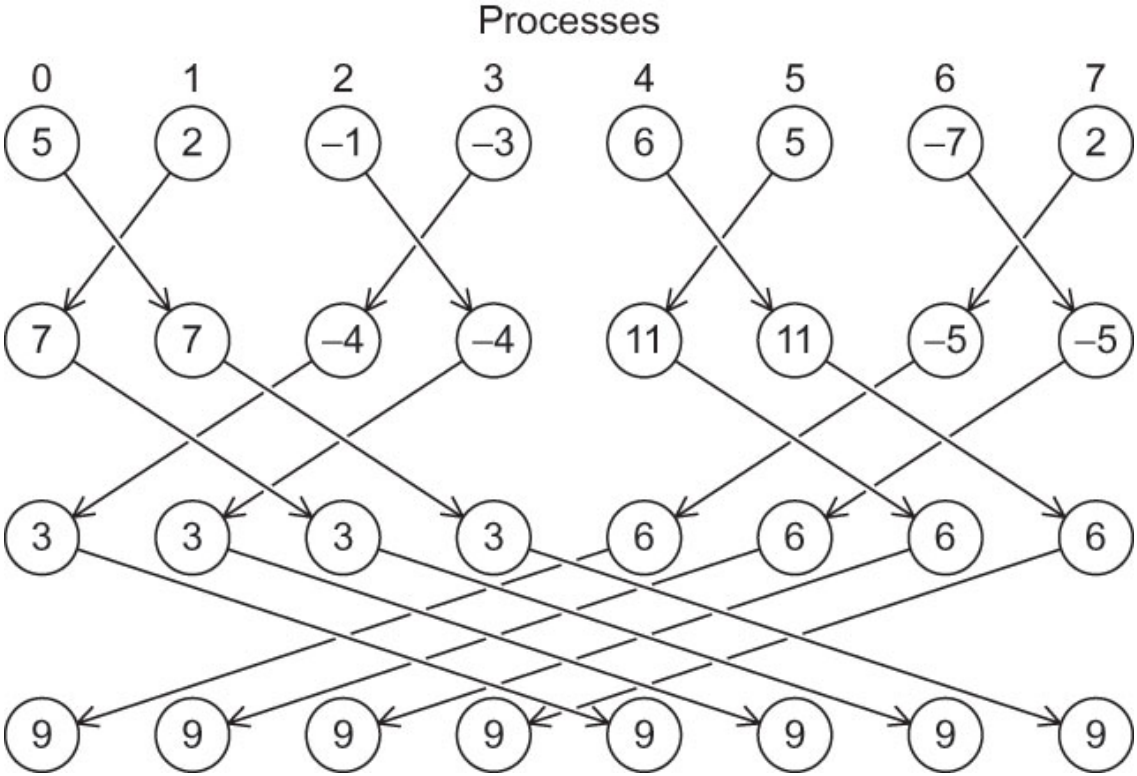
A global sum followed by distribution of the result.





4.集体通信（Collective communication）

●MPI_Allreduce



A butterfly-structured global sum.

4.集体通信（Collective communication）

●MPI_Allreduce

- 同样的，我们不希望去决定使用哪种通信结构，如何去优化代码
- MPI 提供了 MPI_Reduce 的变种，在所有进程中存放结果

```
int MPI_Allreduce(  
    void*          input_data_p    /* in    */,  
    void*          output_data_p   /* out   */,  
    int            count           /* in    */,  
    MPI_Datatype   datatype        /* in    */,  
    MPI_Op         operator        /* in    */,  
    MPI_Comm       comm            /* in    */);
```

4.集体通信（Collective communication）

●广播（Broadcast）

- 将属于一个进程的数据发送到通信器中的所有进程

```
int MPI_Bcast(  
    void*      data_p      /* in/out */,  
    int        count       /* in      */,  
    MPI_Datatype datatype   /* in      */,  
    int        source_proc  /* in      */,  
    MPI_Comm   comm        /* in      */);
```

- source_proc 进程将 data_p 指向的内容发送给 comm 通信器中的所有进程

4.集体通信 (Collective communication)

●广播 (Broadcast)

```
void Get_input(  
    int      my_rank    /* in */,  
    int      comm_sz    /* in */,  
    double*  a_p        /* out */,  
    double*  b_p        /* out */,  
    int*     n_p        /* out */) {  
  
    if (my_rank == 0) {  
        printf("Enter a, b, and n\n");  
        scanf("%lf %lf %d", a_p, b_p, n_p);  
    }  
    MPI_Bcast(a_p, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);  
    MPI_Bcast(b_p, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);  
    MPI_Bcast(n_p, 1, MPI_INT, 0, MPI_COMM_WORLD);  
} /* Get_input */
```

4.集体通信 (Collective communication)

●数据分布

$$\begin{aligned}\mathbf{x} + \mathbf{y} &= (x_0, x_1, \dots, x_{n-1}) + (y_0, y_1, \dots, y_{n-1}) \\ &= (x_0 + y_0, x_1 + y_1, \dots, x_{n-1} + y_{n-1}) \\ &= (z_0, z_1, \dots, z_{n-1}) \\ &= \mathbf{z}\end{aligned}$$

Compute a vector sum.

4.集体通信 (Collective communication)

●数据分布

```
1 void Vector_sum(double x[], double y[], double z[], int n) {  
2     int i;  
3  
4     for (i = 0; i < n; i++)  
5         z[i] = x[i] + y[i];  
6 } /* Vector_sum */
```

4.集体通信（Collective communication）

●数据分布

➤并行化

- 对应元素相加
- 任务之间不需要通信
- 映射到核心（core）：如何划分？
 - 块划分（block partition）
 - 循环划分（cyclic partition）
 - 块-循环划分（block-cyclic partition）

4.集体通信（Collective communication）

●数据分布

➤并行化

- 将12维向量映射到3个进程的不同划分

Process	Components											
	Block				Cyclic				Block-cyclic Blocksize = 2			
0	0	1	2	3	0	3	6	9	0	1	6	7
1	4	5	6	7	1	4	7	10	2	3	8	9
2	8	9	10	11	2	5	8	11	4	5	10	11

4.集体通信 (Collective communication)

●数据分布

➤并行化

```
1 void Parallel_vector_sum(  
2     double local_x[] /* in */,  
3     double local_y[] /* in */,  
4     double local_z[] /* out */,  
5     int local_n /* in */) {  
6     int local_i;  
7  
8     for (local_i = 0; local_i < local_n; local_i++)  
9         local_z[local_i] = local_x[local_i] + local_y[local_i];  
10 }
```

4.集体通信 (Collective communication)

●分发 (Scatter)

➤ *进程0 读入向量 x 和 y , 广播给其他进程?*

- 如果10个进程, 向量为10000维
- 每个进程为向量分配10000维空间, 仅操作其中的1000维
- 如果采用块划分, 我们希望进程0仅发送1000-1999维给进程1, 2000-2999给进程2,

4.集体通信（Collective communication）

●分发（Scatter）

➤MPI_Scatter 可以用来在进程 0 读入整个向量，向其他进程发送它们所需的数据

```
int MPI_Scatter(  
    void*          send_buf_p    /* in */,  
    int            send_count    /* in */,  
    MPI_Datatype   send_type     /* in */,  
    void*          recv_buf_p    /* out */,  
    int            recv_count    /* in */,  
    MPI_Datatype   recv_type     /* in */,  
    int            src_proc      /* in */,  
    MPI_Comm       comm         /* in */);
```

4.集体通信（Collective communication）

●分发（Scatter）

- 如果通信器 comm 包含 comm_sz 个进程，则 MPI_Scatter 将 send_buf_p 指向的数据分成 comm_sz 份，第一份发给进程 0，第二份发给进程 1，以此类推
- 其他进程将本地向量作为 rece_buf_p 参数，将 local_n 作为 recv_count 参数

```
int MPI_Scatter(  
    void*          send_buf_p    /* in */,  
    int            send_count    /* in */,  
    MPI_Datatype   send_type     /* in */,  
    void*          recv_buf_p    /* out */,  
    int            recv_count    /* in */,  
    MPI_Datatype   recv_type     /* in */,  
    int            src_proc      /* in */,  
    MPI_Comm       comm         /* in */);
```

4.集体通信（Collective communication）

●分发（Scatter）

- 假设使用块划分，进程0读入向量的所有 n 维数据到 `send_buf_p` 中
- 进程0获得第一块 $local_n = n / comm_sz$ 维数据，进程1获得下一块 $local_n$
- 每个进程传递本地向量作为 `recv_buf_p` 参数，`recv_count` 应为 $local_n$
- `send_count` 应为发送给每个进程的数据量，即： $local_n$

```
int MPI_Scatter(  
    void*      send_buf_p  /* in  */,  
    int        send_count  /* in  */,  
    MPI_Datatype send_type  /* in  */,  
    void*      recv_buf_p  /* out */,  
    int        recv_count  /* in  */,  
    MPI_Datatype recv_type  /* in  */,  
    int        src_proc    /* in  */,  
    MPI_Comm   comm        /* in  */);
```

4.集体通信（Collective communication）

```
1 void Read_vector(  
2     double    local_a[]    /* out */,  
3     int        local_n     /* in  */,  
4     int        n           /* in  */,  
5     char       vec_name[]  /* in  */,  
6     int        my_rank     /* in  */,  
7     MPI_Comm   comm        /* in  */) {  
8  
9     double* a = NULL;  
10    int i;  
11  
12    if (my_rank == 0) {  
13        a = malloc(n*sizeof(double));  
14        printf("Enter the vector %s\n", vec_name);  
15        for (i = 0; i < n; i++)  
16            scanf("%lf", &a[i]);  
17        MPI_Scatter(a, local_n, MPI_DOUBLE, local_a, local_n,  
18                    MPI_DOUBLE, 0, comm);  
19        free(a);  
20    } else {  
21        MPI_Scatter(a, local_n, MPI_DOUBLE, local_a, local_n,  
22                    MPI_DOUBLE, 0, comm);  
23    }  
24 } /* Read_vector */
```

4.集体通信（Collective communication）

●收集（Gather）

➤进程 0 收集向量的所有分量进行后续处理

```
int MPI_Gather(  
    void*          send_buf_p    /* in */,  
    int            send_count    /* in */,  
    MPI_Datatype    send_type    /* in */,  
    void*          recv_buf_p    /* out */,  
    int            recv_count    /* in */,  
    MPI_Datatype    recv_type    /* in */,  
    int            dest_proc     /* in */,  
    MPI_Comm        comm         /* in */);
```

4.集体通信（Collective communication）

●收集（Gather）

- 进程 0 中由 send_buf_p 引用的数据存放在 recv_buf_p 中第一块位置，进程 1 中由 send_buf_p 引用的数据存放在 recv_buf_p 中第二块位置，以此类推
- recv_count 为从每个进程接收到的数据量

```
int MPI_Gather(  
    void*      send_buf_p  /* in  */,  
    int        send_count  /* in  */,  
    MPI_Datatype send_type  /* in  */,  
    void*      recv_buf_p  /* out */,  
    int        recv_count  /* in  */,  
    MPI_Datatype recv_type  /* in  */,  
    int        dest_proc   /* in  */,  
    MPI_Comm    comm       /* in  */);
```


4.集体通信 (Collective communication)

●收集 (Gather)

```
1 void Print_vector(  
2     double    local_b[] /* in */,  
3     int        local_n  /* in */,  
4     int        n        /* in */,  
5     char       title[]  /* in */,  
6     int        my_rank  /* in */,  
7     MPI_Comm   comm     /* in */) {  
8  
9     double* b = NULL;  
10    int i;  
11  
12    if (my_rank == 0) {  
13        b = malloc(n*sizeof(double));  
14        MPI_Gather(local_b, local_n, MPI_DOUBLE, b, local_n,  
15                  MPI_DOUBLE, 0, comm);  
16        printf("%s\n", title);  
17        for (i = 0; i < n; i++)  
18            printf("%f ", b[i]);  
19        printf("\n");  
20        free(b);  
21    } else {  
22        MPI_Gather(local_b, local_n, MPI_DOUBLE, b, local_n,  
23                  MPI_DOUBLE, 0, comm);  
24    }  
25 } /* Print_vector */
```

4.集体通信 (Collective communication)

- 全收集 (Allgather) : 矩阵与向量相乘

$A = (a_{ij})$ is an $m \times n$ matrix

\mathbf{x} is a vector with n components

$\mathbf{y} = A\mathbf{x}$ is a vector with m components

$$y_i = a_{i0}x_0 + a_{i1}x_1 + a_{i2}x_2 + \cdots a_{i,n-1}x_{n-1}$$

i -th component of y

Dot product of the i -th row of A with x .

4.集体通信（Collective communication）

●全收集（Allgather）：矩阵与向量相乘

a_{00}	a_{01}	\cdots	$a_{0,n-1}$
a_{10}	a_{11}	\cdots	$a_{1,n-1}$
\vdots	\vdots		\vdots
a_{i0}	a_{i1}	\cdots	$a_{i,n-1}$
\vdots	\vdots		\vdots
$a_{m-1,0}$	$a_{m-1,1}$	\cdots	$a_{m-1,n-1}$

x_0
x_1
\vdots
x_{n-1}

$=$

y_0
y_1
\vdots
$y_i = a_{i0}x_0 + a_{i1}x_1 + \cdots a_{i,n-1}x_{n-1}$
\vdots
y_{m-1}

4.集体通信 (Collective communication)

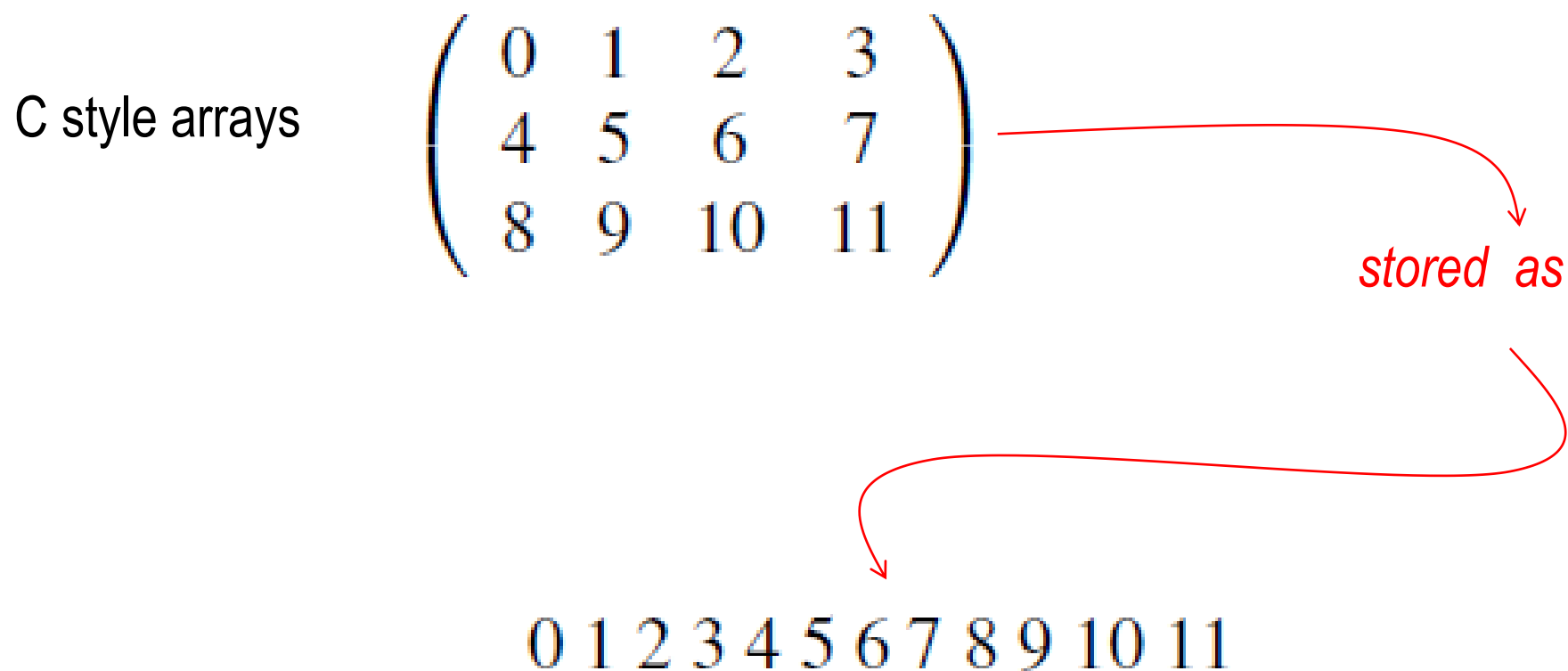
- 全收集 (Allgather) : 矩阵与向量相乘

```
/* For each row of A */  
for (i = 0; i < m; i++) {  
    /* Form dot product of ith row with x */  
    y[i] = 0.0;  
  
    for (j = 0; j < n; j++)  
        y[i] += A[i][j]*x[j];  
}
```

Serial pseudo-code

4.集体通信 (Collective communication)

- 全收集 (Allgather) : 矩阵与向量相乘



4.集体通信 (Collective communication)

●全收集 (Allgather) : 矩阵与向量相乘

```
1 void Mat_vect_mult(  
2     double A[] /* in */,  
3     double x[] /* in */,  
4     double y[] /* out */,  
5     int m /* in */,  
6     int n /* in */) {  
7     int i, j;  
8  
9     for (i = 0; i < m; i++) {  
10        y[i] = 0.0;  
11        for (j = 0; j < n; j++)  
12            y[i] += A[i*n+j]*x[j];  
13    }  
14 } /* Mat_vect_mult */
```

4.集体通信（Collective communication）

●全收集（Allgather）：矩阵与向量相乘

➤并行化

- $y[i] += A[i * n + j] * x[j];$
- 将 A 按行划分
- $y[i]$ 的计算包含了 A 中第 i 行的所有元素和 x 中的所有元素

4.集体通信（Collective communication）

●全收集（Allgather）：矩阵与向量相乘

➤并行化

- 为减少通信，将 x 中的所有元素分配给每个进程
- 实际应用中，矩阵与向量相乘的结果 y 可能是下一次相乘的输入
- 如何使得每个进程获得相乘的结果？
- `MPI_Gather + MPI_Bcast` ？

4.集体通信（Collective communication）

●全收集（Allgather）：矩阵与向量相乘

➤并行化

- 将每个进程的 `send_buf_p` 的内容拼接起来存放到每个进程的 `recv_buf_p`
- `recv_count` 为从每个进程接收到的数据量

```
int MPI_Allgather(  
    void*          send_buf_p /* in */,  
    int            send_count /* in */,  
    MPI_Datatype    send_type  /* in */,  
    void*          recv_buf_p /* out */,  
    int            recv_count  /* in */,  
    MPI_Datatype    recv_type  /* in */,  
    MPI_Comm        comm       /* in */);
```

4.集体通信（Collective communication）

●全收集（Allgather）：矩阵与向量相乘

➤并行化

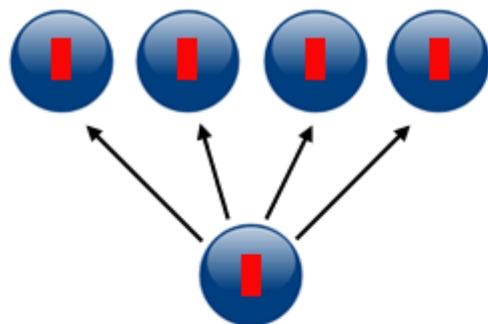
```
1 void Mat_vect_mult(  
2     double    local_A[] /* in */,  
3     double    local_x[] /* in */,  
4     double    local_y[] /* out */,  
5     int        local_m /* in */,  
6     int        n        /* in */,  
7     int        local_n /* in */,  
8     MPI_Comm   comm     /* in */) {  
9     double* x;  
10    int local_i, j;  
11    int local_ok = 1;  
12  
13    x = malloc(n*sizeof(double));  
14    MPI_Allgather(local_x, local_n, MPI_DOUBLE,  
15                x, local_n, MPI_DOUBLE, comm);  
16  
17    for (local_i = 0; local_i < local_m; local_i++) {  
18        local_y[local_i] = 0.0;  
19        for (j = 0; j < n; j++)  
20            local_y[local_i] += local_A[local_i*n+j]*x[j];  
21    }  
22    free(x);  
23 } /* Mat_vect_mult */
```

4.集体通信（Collective communication）

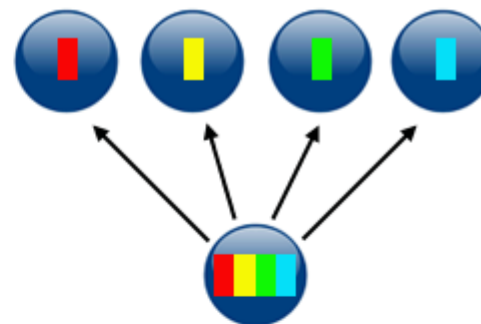
●集体操作的类型

- 同步（Synchronization）：进程等待组内其他进程到达同步点
- 数据移动（Data Movement）：broadcast, scatter/gather
- 集体计算（Collective Computation）：一个进程收集其他进程的数据，并对该数据进行操作（min, max, add, multiply 等）

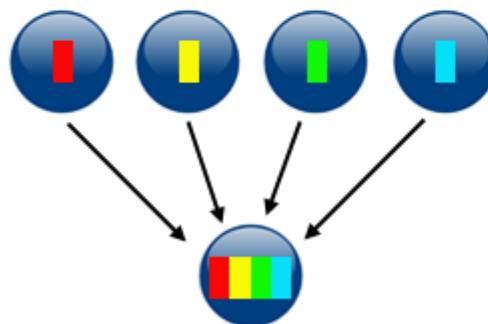
4.集体通信 (Collective communication)



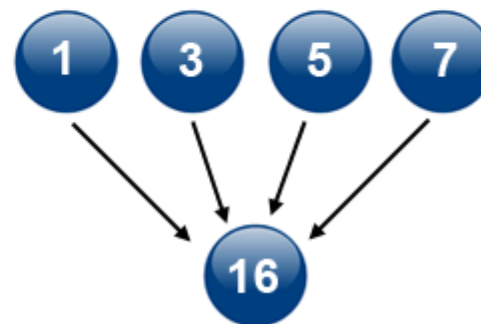
broadcast



scatter



gather



reduction