

# 正则采样排序PSRS的MPI算法

## 算法流程

假设有  $p$  个进程，有  $N$  条数据需要排序。

### 1. 均匀划分

将  $N$  条数据均匀划分为  $p$  段，每个进程处理一段数据。其中  $i$  ( $i = 0, 1, \dots, p - 1$ ) 号进程处理  $\lfloor \frac{i \times N}{p} \rfloor$  到  $\lfloor \frac{(i + 1) \times N}{p} \rfloor - 1$  行（行号从0开始计算）；

3个进程划分18个数据的结果如下，一个进程处理一段。

Processor 1						Processor 2						Processor 3					
48	39	6	72	91	14	69	40	89	61	12	21	84	58	32	33	72	20

在实现上，由于数据较大（32G），为了充分并行节省时间，由各进程读取所需要的数据，同样达到均匀划分效果。实现代码如下：

```
1  #define BLOCK_LOW(my_rank, comm_sz, T) ((my_rank)*(T)/(comm_sz))
2  #define BLOCK_SIZE(my_rank, comm_sz, T) (BLOCK_HIGH(my_rank,comm_sz,T) -
   BLOCK_LOW(my_rank,comm_sz,T) + 1)
3
4  // 打开文件
5  ifstream fin(fileName, ios::binary);
6  // 计算各进程读取文件的起始行号 and 大小
7  unsigned long myDataStart = BLOCK_LOW(my_rank, comm_sz, dataLength);
8  unsigned long myDataLength = BLOCK_SIZE(my_rank, comm_sz, dataLength);
9  // 将文件指针移动到起始行号
10 fin.seekg((myDataStart)*sizeof(unsigned long), ios::beg);
11 unsigned long *myData = new unsigned long[myDataLength];
12 // 读取数据
13 for(unsigned long i=0; i<myDataLength; i++)
14     fin.read((char*)&myData[i], sizeof(unsigned long));
15 fin.close();
```

### 2. 局部排序

各进程对各自的数据进行排序。

例子如下，各段数据各自有序：

6	14	39	48	72	91
---	----	----	----	----	----

12	21	40	61	69	89
----	----	----	----	----	----

20	32	33	58	72	84
----	----	----	----	----	----

代码如下，可以调用算法库：

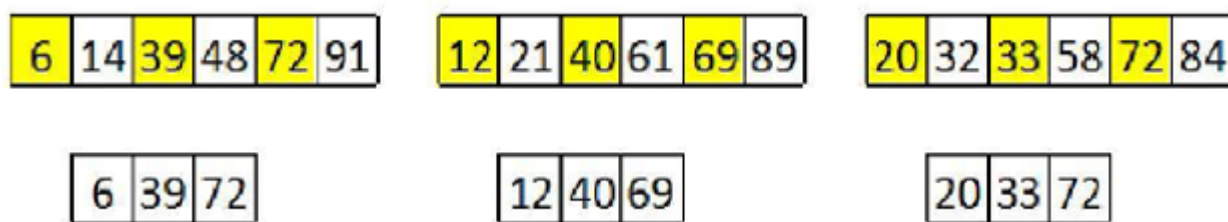
```
1 | sort(myData, myData+myDataLength);
```

### 3. 选取样本

$p$  个进程中，每个进程需要选取出  $p$  个样本（regular samples），选取规则为

$\lfloor \frac{i \times dataLength}{p} \rfloor, i = 0, 1, \dots, p-1$ 。 `dataLength` 是进程各自的数据长度。注意此时选取的  $p$  个样本也是有序的。

如下图三个进程中，每个进程选取出三个样本：



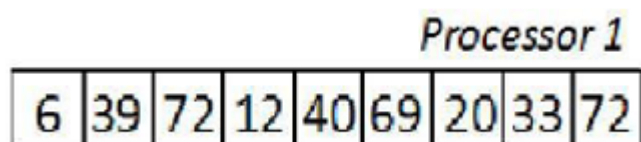
代码如下：

```
1 | unsigned long *regularSamples = new unsigned long[comm_sz];  
2 | for(int index=0; index<comm_sz; index++)  
3 |     regularSamples[index] = myData[(index*myDataLength)/comm_sz];
```

### 4. 样本排序

用一个进程对  $p$  个进程的共  $p \times p$  个样本进行排序，此时样本都是局部有序的，使用归并能减少时间复杂度。

如下图，processor 1 将 9 个样本排序：



在实现中，由 0 号进程负责样本排序，首先需要 `Gather` 操作，将样本收集起来，然后要进行归并。

`Gather` 操作：

```
1 | unsigned long *gatherRegSam;  
2 | if(my_rank == 0)  
3 |     gatherRegSam = new unsigned long[comm_sz*comm_sz];  
4 | // sendbuf, sendcount, sendtype, recvbuf, recvcount, recvtype, root, comm  
5 | MPI_Gather(regularSamples, comm_sz, MPI_UNSIGNED_LONG, gatherRegSam, comm_sz,  
    MPI_UNSIGNED_LONG, 0, MPI_COMM_WORLD);
```

归并函数：

```
1 | struct mmdata {  
2 |     // 待归并数组序号
```

```

3     int stindex;
4     // 在数组中的序号
5     int index;
6     unsigned long stvalue;
7     mmdata(int st=0, int id=0, unsigned long stv =
0):stindex(st),index(id),stvalue(stv){}
8 };
9
10 // mmdata比较运算符重载, 在优先队列的表现是mmdata中stvalue小的排在队列前面
11 bool operator<( const mmdata & One, const mmdata & Two) {
12     return One.stvalue > Two.stvalue;
13 }
14
15 // 各进程regularSamples二维数组, 各进程regularSamples长度, 待归并数组数量, 结果数组, 待归并总
数据量
16 void multiple_merge(unsigned long* starts[], const int lengths[], const int Number,
unsigned long newArray[], const int newArrayLength) {
17     priority_queue< mmdata> priorities;
18
19     // 将每个待归并数组的第一个数加入优先队列, 同时保存它所在待归并数组序号和数字在数组中的序号
20     for(int i=0; i<Number;i++) {
21         if (lengths[i]>0) {
22             priorities.push(mmdata(i,0,starts[i][0]));
23         }
24     }
25
26     int newArrayindex = 0;
27     while (!priorities.empty() && (newArrayindex<newArrayLength)) {
28         // 拿下最小的数据
29         mmdata xxx = priorities.top();
30         priorities.pop();
31
32         // 将拿下的数据加入到结果数组中
33         newArray[newArrayindex++] = starts[xxx.stindex][xxx.index];
34
35         // 如果被拿下数据不是所在的待归并数组的最后一个元素, 将下一个元素push进优先队列
36         if ( lengths[xxx.stindex]>(xxx.index+1)) {
37             priorities.push(mmdata(xxx.stindex, xxx.index+1, starts[xxx.stindex]
[xxx.index+1]));
38         }
39     }
40 }

```

归并样本:

```

1 if(my_rank == 0) {
2     // start用于存储gatherRegSam中各进程RegularSamples开始下标, 相当于二维数组
3     unsigned long **starts = new unsigned long* [comm_sz];
4     // gatherRegSam中各进程RegularSamples长度, 都一样是comm_sz
5     int *lengths = new int[comm_sz];
6     for(int i=0; i<comm_sz; i++) {
7         starts[i] = &gatherRegSam[i*comm_sz];
8         lengths[i] = comm_sz;

```

```

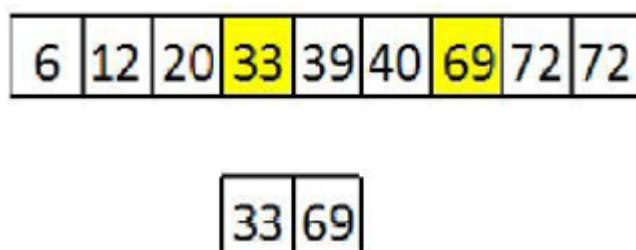
9      }
10
11     // 因为各进程的的ragularSamples就是有序的，因此只需要将gatherRegSam中的各进程数据归并即可
12     unsigned long *sortedGatRegSam = new unsigned long[comm_sz*comm_sz];
13     multiple_merge(starts, lengths, comm_sz, sortedGatRegSam, comm_sz*comm_sz);
14 }

```

## 5. 选取主元

一个进程从排好序的样本中选取  $p - 1$  个主元。选取方法是  $i \times p, i = 1, 2, \dots, p - 1$ 。

例子如下，从9个样本中选取2个主元：



选取主元代码如下：

```

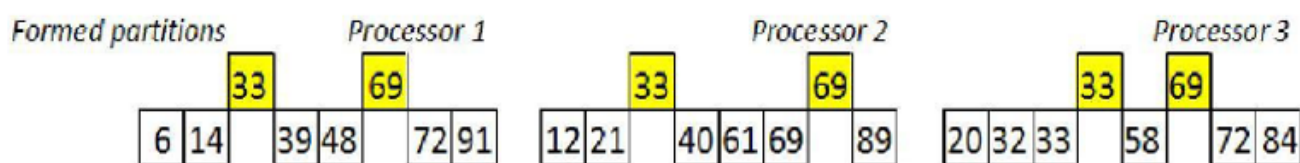
1  for(int i=0; i<comm_sz-1; i++)
2      privots[i] = sortedGatRegSam[(i+1)*comm_sz];

```

## 6. 主元划分

$p$  个进程的数据按照  $p - 1$  个主元划分为  $p$  段。

例子如下，3个进程的数据被都被2个主元划分为3段：



在具体实现中，0号进程要将  $p - 1$  个主元广播到其它所有进程。然后所有进程将按照主元划分。

```

1  // 将主元广播到其它进程
2  MPI_Bcast(privots, comm_sz-1, MPI_UNSIGNED_LONG, 0, MPI_COMM_WORLD);
3  // partStartIndex保存每段开始下标
4  int *partStartIndex = new int[comm_sz];
5  // partLength保存每段的长度
6  int *partLength = new int[comm_sz];
7  unsigned long dataIndex = 0;
8  for(int partIndex = 0; partIndex<comm_sz-1; partIndex++) {
9      partStartIndex[partIndex] = dataIndex;
10     partLength[partIndex] = 0;
11
12     while((dataIndex<myDataLength) && (myData[dataIndex]<=privots[partIndex])) {
13         dataIndex++;
14         partLength[partIndex]++;

```

```

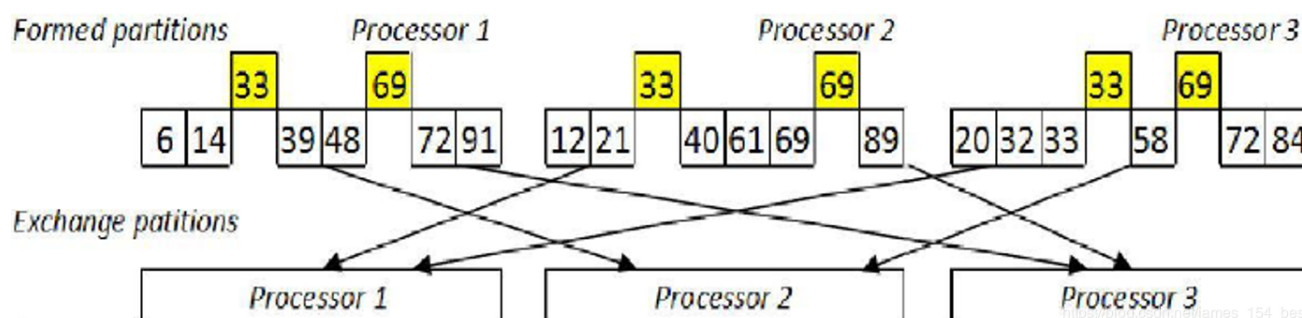
15     }
16 }
17 // 最后一段数据补齐 (防止主元溢出)
18 partStartIndex[comm_sz-1] = dataIndex;
19 partLength[comm_sz-1] = myDataLength - dataIndex;

```

## 7. 全局交换

进程  $i$  ( $i = 0, 1, \dots, p-1$ ) 将第  $j$  ( $j = 0, 1, \dots, p-1$ ) 段发送给进程  $j$ 。也就是每个进程都要给其它所有进程发送数据段，并且还要从其它所有进程中接收数据段，所以称为全局交换。

例子如下：



在具体实现中，可以使用 `MPI_Alltoall` 和 `MPI_Alltoallv` 简化操作，进程  $i$  先使用 `MPI_Alltoall` 将第  $j$  段数据的长度发送给进程  $j$ （或者说进程  $j$  使用 `MPI_Alltoall` 接收进程  $i$  的第  $j$  段数据的长度），进程  $j$  将这个长度保存在 `recvRankPartLen[i]` 中，进程  $j$  知道长度后便可以使用 `MPI_Alltoallv` 接收进程  $i$  的第  $j$  段数据了。

```

1 // 接收各进程数据段长度到recvRankPartLen中
2 int *recvRankPartLen = new int[comm_sz];
3 MPI_Alltoall(partLength, 1, MPI_INT, recvRankPartLen, 1, MPI_INT, MPI_COMM_WORLD);
4
5 // 计算接收到的数据段在缓存中的摆放位置 (距离起始地址的偏移量)
6 int rankPartLenSum = 0;
7 int *rankPartStart = new int[comm_sz];
8 for(int i=0; i<comm_sz; i++) {
9     rankPartStart[i] = rankPartLenSum;
10    rankPartLenSum += recvRankPartLen[i];
11 }
12 unsigned long *recvPartData = new unsigned long[rankPartLenSum];
13 // 发送数据段和接收数据段
14 MPI_Alltoallv(myData, partLength, partStartIndex, MPI_UNSIGNED_LONG, recvPartData,
15    recvRankPartLen, rankPartStart, MPI_UNSIGNED_LONG, MPI_COMM_WORLD);

```

需要注意的是

```

1 int MPI_Alltoallv(const void *sendbuf, const int *sendcounts, const int
  *sdispls, MPI_Datatype sendtype, void *recvbuf, const int *recvcounts, const int
  *rdispls, MPI_Datatype recvtype, MPI_Comm comm)

```



