# Parallel Computing and Programming Report

Wanghua Shi
s-whua@sjtu.edu.com
Shanghai Jiao Tong University
Shanghai, China

## Abstract

This report is the projects of course "Parallel Computing and Programming". There are three projects including seven questions. The first project is MPI programming. The second project is OMP programming. The third project is a simple hadoop application.

For each question, I design experiments to verify the correctness and performance, and then analysis the result. The number of code line is about 2000(C/C++) and 1000(python). The python code is for drawing pictures. The project code can be obtained at https://github.com/shiwanghua/SharedFiles/tree/main/%E5%B9%B6%E8%A1%8C%E8%AE%A1%E7%AE%97%E4%B8%8E%E5%B9%B6%E8%A1%8C%E7%AE%97%E6%B3%95.

***CCS Concepts:*** • **MPI → OpenMP**.

***Keywords:*** MPI,OpenMP,hadoop,parallelism

**ACM Reference Format:**

## 1 Project One

"p" represents the number of processes ot threads. Without loss of generality, use randomly generated data of double type to do experiments.

### 1.1 *MPI_Allgather*

**1.1.1 Core Idea.** We can use the implementation of question 6.10 in howework 3 to implement customed *MPI_Allgather*. In 6.11, I realize *MPI_Bcast*. Calling the self-defined *MPI_Bcast* for p times based on each process can realize the same function of *MPI_Allgather*. In self-define *MPI_Bcast* function, the data from root should be stored in receive buffer, and then send data to other p-1 processes based on receive buffer.

### 1.1.2 Experiments.

**1. Verification.** Set p=8, each process generate 10 integer, the i-th integer is i times of the process id. After the *Swhua_Allgather* function is finished, output the receiving data in process 0:

**Figure 1.** Validate *Swhua_Bcast* function.

The "size" means the size of data transferred by each process. The running time is about 70.6 us. Line i is from process i. The j-th in i-th line is j times of number i. Hence, the function is correct.

**2. Process Number Exp.** Let each process send data of 256MB. Run *MPI_Allgather* and *Swhua_Allgather* with 1 to 29 processes and measure the running time. The result is shown in picture 2.
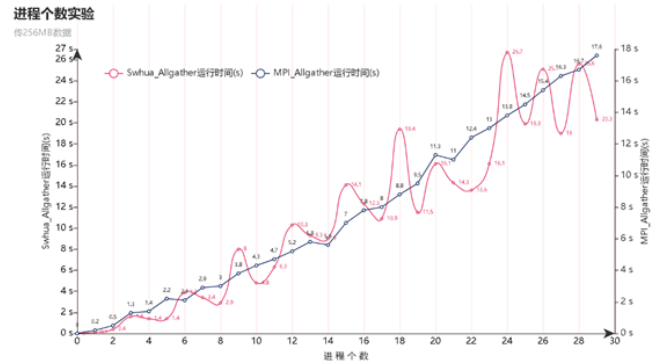


**Figure 2.** Running time of different process number.

The more the processes, the more the data needed to be broadcast. When p = 24, the running time of *Swhua_Bcast* reaches its maximum value 26.7 s. Overall, my version is lower than the MPI version. Besides, my version is not so stable as the data gets bigger.

**3. Process Number Exp2.** To ensure the data needed to transfered is a constant, I design another Process Number Exp. All processes sends 256MB data in total. So each process needs to broadcast $\frac{256MB}{p}$ data. Do experiments with p = 1 to 30. The result is as follows:
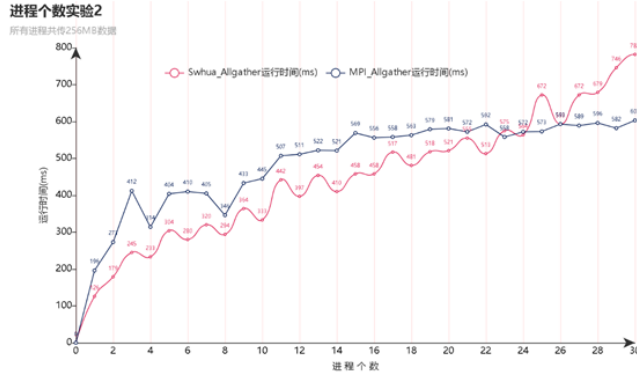
**Figure 3.** Running time of different process number(Fixed data size in total).

When p < 23, my Bcast function is faster. Since the data size is constant, my function becomes a little bit more stable than before.

*4. DataTransmission Exp.* Set p = 16. Let each process send 8B, 1KB, 1MB,64MB,256MB,512MB, the running time in millisecond is represented in figure 4.
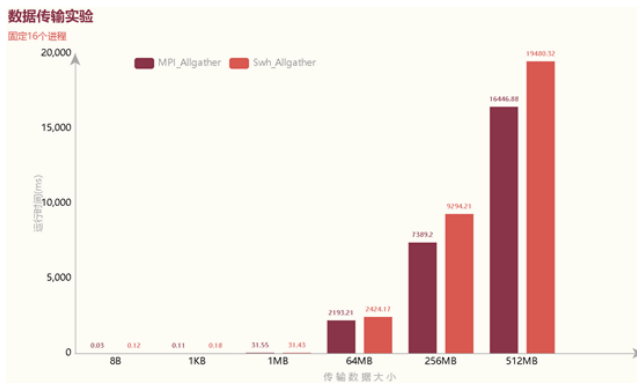


**Figure 4.** Running time in different data size.

The first three experiments have low running time although the data size is magnified many times, which means the growth of running time is not as fast as the growth of data volume. In the last three cases, the circumstances are the opposite. My implementation is all slower than *MPI_Allgather* in the six situations.

Moreover, nonblocking function *MPI_Isend* can be used in Bcast function. But I got slower elapsed time.

## 1.2 Gemm

### 1.2.1 Core Idea.
For block matrix multiplication, use *MPI_Dims_create* [1] to get the row number and column number of grids. Use *MPI_Cart_create* to create communication domain. Use *MPI_Cart_coords* to get the grid coordinate of current process. Use *MPI_Comm_split* to create single row and single column internal communication domain.

Use *MPI_Cart_rank* to get the grid ID of current process. Let process 0 be root process. Root process read or generate data and send data to relative process by *MPI_Send* or *MPI_Scatterv*. After every process get their data and finish calculating, for each row, use *MPI_Gatherv* and row communicators to collect the partial results of each column to the first process of each row. At last, use the first column communicator to collect the row results of the first column to root process and save the final result.

Convolutional operation can also be handled by this way. Set the step size to 1, the padding parameter to 0. The kernel size is marked "k" and original matrix has size (n × n). The outcome matrix has size ($n-k+1$, $n-k+1$). But it is necessary to send (k-1) more rows and columns to calculate the result in the last row and the last column. For max pooling operation, the procedure is basically the same as convolutional operation but the kernel is not needed.

### 1.2.2 Experiments.

*1. Verification.* Generate three matrices a1, a2 and a3 in natural number sequence, each of which has 4 rows and 8 columns. Generate a matrix b in natural number sequence which has 4 rows and 8 columns. Generate a matrix k in natural number sequence which has 2 rows and 2 columns. Calculate $a1 \times b^T$, convolution of k on a2, (2×2) max pooling of a3. The results appear as shown in fig5.

It is easy to see that every number is correct. The running time of matrix multiplication is 9.93ms, convolution is 6.087ms, and pooling is 3.884ms, which is faster and faster and logical.

*2. Process Number Exp.* Run experiments with 1 to 50 processes. In each time, set the shape of a and b to 1024 × 1024, set the kernel shape to 4 × 4. The running time curve is shown in fig6.

The row and column number of the grid under each process number is printed out on screen. For matrix multiplication, the fastest running time is 3.57 s when there are 10 processes. At this time, there are 5 grids in each column and 2 grids in each row; For convolution operation, situation of 17 processes is the fastest, followed by 7 processes, both less than 1.9 s; For pooling operation, the fastest time is reached at 19 processes, and the time is further reduced to 1.69 s. In general, the changing trend of the running time of the three operations is roughly the same.

## 1.3 Wordcount

### 1.3.1 Core Idea.
Manager/Work mode is adopted in this question. Let process 0 be root again. Root process read text from files and send segmental text to other p-1 processes. When all the text is sent out, root process send finish message to inform other processes to do reduce operation. Once receiving a finish message from one process, root process will handle the partial result of that process. After reducing
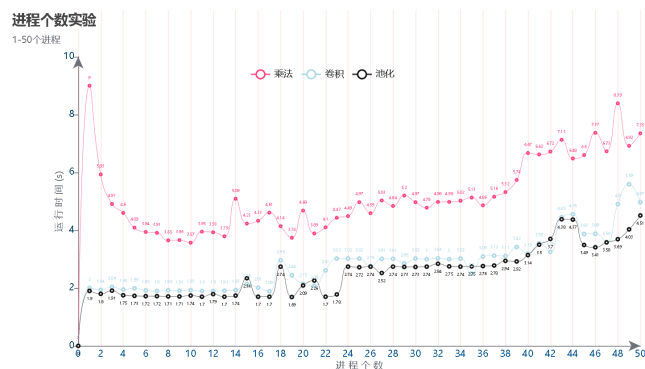
**Figure 5.** Verification.



**Figure 6.** Runing time in different process number.

all segmental results, root process sorts and saves the word count result to local file named "*final_result.txt*".

Other p-1 processes will apply for text and find words from text and count the words continuously. Once receiving a finish message, they will stop map operation and save their segmental results to temporary files. At last, they will send a message to root process to show that they have finished their tasks.

For the big file, I try two ways to distribute the text information: send a fixed length of text each time (Special treatment at the end); send one average length of text to each process, each process only needs to receive one time. Besides, there is a **space** at the end of each line in big file, and a word at the end of one line may be divided into two parts by a special **hyphens**, whose size is 2 bytes.

For small files, I also try two ways to distribute text: for each request from processes, root process send one small file based on original file sequence; arrange files in descending order according to file length, ensuring that the number of characters to be processed by each process is close. In addition, there are two **empty** files. For sorted method, once a empty file is encountered, the program can end immediately.

### 1.3.2 Experiments.

*1. Verification.* The file folder "*big_file*" and "*small_file*" are put in the same path as "1.3.c" file. The sorted word count results will be saved in "*big_file/final_result.txt*" and "*small_file/final_result.txt*".

According to my word recognition algorithm, I get the following results. In the big file, there are 147792 words, 10290 of which are different words. The top five words are shown in fig 7.



**Figure 7.** Top 5 words count in the big file.

In addition, I find there are 8 words "echo", 88 words "against" and 596 words "from" in the big file by using the search function (ctrl + F) of VSCode, which are the same as mine.

For the small files, there are 346680 words, 35386 of which are different words. The top five words are shown in fig 8.

*2. Big file Process Number Exp.* Set p = 1 to 50. Root process is not included in p in this experiment. For the fixed length version, the buffer size is set to 1024B. Record the runing time of two word count algorithms for the big file. The result is fig 9. From the picture we can easily see that the

**Figure 8.** Top 5 words count in small files.



**Figure 10.** Runing time of two small files word count algorithms in different process number.

average sending version is better in most cases. This is because the number of communications has decreased. When p = 5, the running time reaches its minimum, 52 ms.
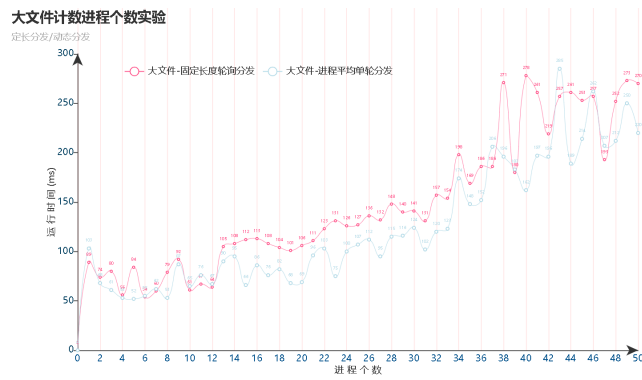


**Figure 9.** Runing time of two big file word count algorithms in different process number.

**3. Small files Process Number Exp.** Set p = 1 to 50 again. Root process is not included in p in this experiment. Record the runing time of two word count algorithms for small files. The result is fig 10. In most cases, the sort version has lower running time and better stability, which guarantees the fairness for every process. When p = 2 and 3, the running time reaches its minimum, 181 ms.

**4. Stand-alone/Cluster Exp.** Run fixed length big file word count algorithm in stand-alone mode and cluster mode. Set p = 1 to 50. Root process is not included in p. The experimental result is in fig 11. Cluster mode is slightly better than stand-alone mode. When p = 3, cluster mode finishes word count task in 43 ms. Nine milliseconds less than before.

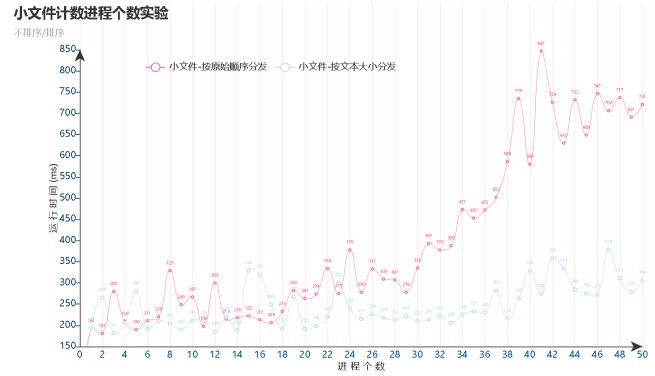**5. Stand-alone/Cluster Exp2.** Run four word count algorithms in stand-alone mode and cluster mode. Set p = 5.
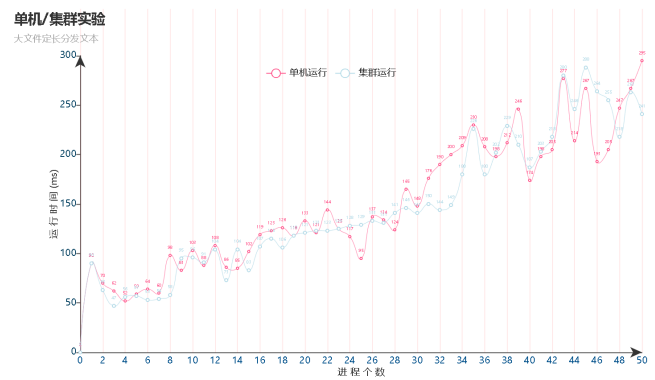


**Figure 11.** Runing time of big file word count task in stand-alone and cluster mode.

Four processes ask for data and do the word count recognition algorithms. The output is in fig 12. Besides the average version of big file, cluster mode shows better performance in the other three cases. Maybe the delay for sending between hosts is too long. When the root process send data to one process, those who are behind the process have to wait.
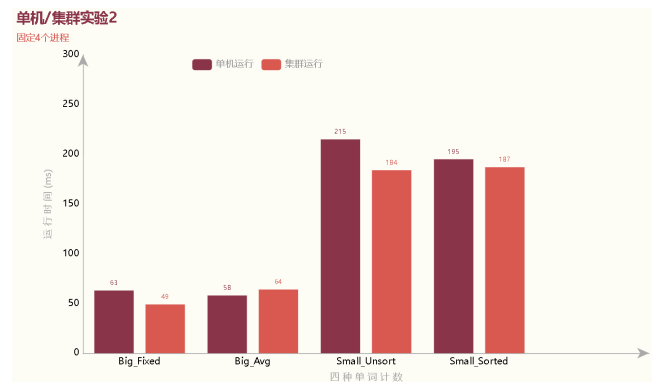


**Figure 12.** Runing time of four algorithms in stand-alone and cluster mode.

## 2 Project Two

### 2.1 Monte Carlo

Input the number of threads and points, count the number of points(n) in unit circle. Use "#pragma omp parallel for" to realize parallelization. Record the running time, the estimated $\pi$ and the accuracy of estimated $\pi$.

Use one single thread to run the simulation program with different n: from 10, 100, 1000, ... to 10000000000. The output is shown in fig 13. When n = 1000, it is accurate enough to estimate $\pi$.
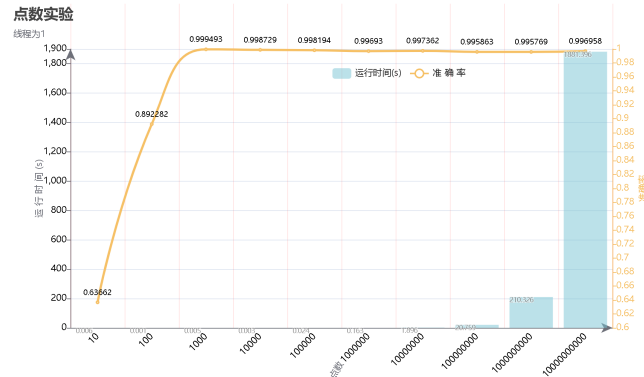


**Figure 13.** Runing time and accuracy in different number of points.

### 2.2 Quick Sort

Once determining a final position of one element, use two threads to sort the left and right subarray. When the array length is 1000000, the running time of omp version is bigger (see fig 14). This is because the required threads are too many. From the figure we can see that the generated number are different.



**Figure 14.** Output when n = 1000000.

When the array length is 10000, less threads are required, so the running time decreases (see fig 15).



**Figure 15.** Output when n = 10000.

If when the length of subarray is smaller than 50, directly use insertion sort instead, then less threads are required, the running time can also be reduced (see fig 16).



**Figure 16.** Runing time while using insertion sort optimization.

### 2.3 PageRank

Set the damping coefficient to 0.85. Use "guided" schedule mode. Use "unordered_set" to store the link number so that there are no duplicate link numbers. Run program in 1 to 50 threads. The running time is in fig 17. The minimum running time is 3735 ms when p = 35.
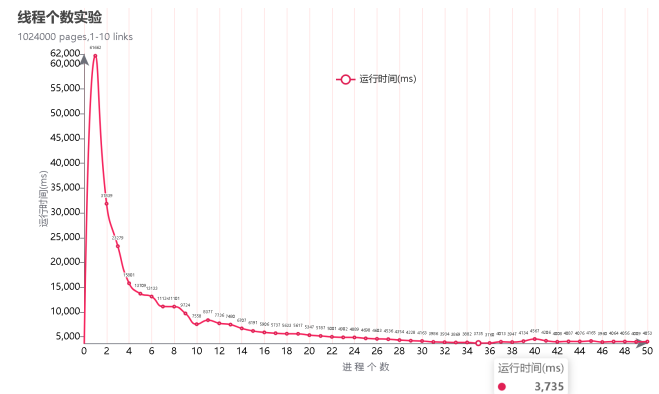


**Figure 17.** Runing time in different number of threads.

## 3   Project Three

Use "load" and "write.table" function of R language to convert RDA file to TXT file. In "sbin" directory of hadoop package, run command "./start-all.sh" to start the cluster. Run command "hadoop fs -mkdir /txtData" to make a directory in hdfs. Run command "hadoop fs -put Ram-N-weatherData-4326bdc/txtData /" to upload TXT files to hdfs. Run command "hadoop fs -ls /txtData" to show the files uploaded. Run command "mvn clean install package" to get the jar package. If it is not the first time to run the program, run command "hadoop fs -rm -r /Result /tmp" to remove the result of the last run of the program. Run command "hadoop jar target/Project3-final.jar GetMaxMinTemperature" to upload and run the jar package. Wait for the program to finish running. Run command "hadoop fs -get /Result ./" to download output files from hdfs. Run command "cat Result/part-r-00000" to show the final maximum and minimum temperature. All the outputs are shwon in fig 18 and fig 19.



**Figure 18.** Output1



**Figure 19.** Output2

The maximum temperature is 91.4. The minimum temperature is 23.0. The running time is 8550 ms.

## Acknowledgments

To teachers and assistants, for providing cluster experimental environment and answering the questions and doubts.

## References

[1] Michael J. Quinn. 2003.   *Parallel Programming in C with MPI and OpenMP* (1st. ed.).  McGraw-Hill Education, New York, NY.