# Design and Realization of Platform of Mass Data Processing Based on High-performance Computer

## Jing Nie

*Institute of Information Engineering, Nanning College of Vocational Technolog, Nanning, 530008, Guangxi Zhuang Autonomous Region, China*

*Corresponding author's e-mail: 353355590@qq.com*

**Abstract**

Today when network information technology develops rapidly, people propose higher requirements on the speed and quality of information processing. For the purpose of satisfying such requirements, we can only rely on the support from high-performance computers. At present, high-performance computers are mainly applied in the field of science and seldom applied in people's daily life relatively speaking. The platform of mass data processing can effectively improve the ability of parallel processing of network information and facilitate the storage, management, processing and utilization of information data to become more standardized and proceduralized. The paper proposes the platform of mass data processing based on high-performance computers, analyzes and inquires into the problems in system application, and puts forward corresponding solutions finally.

*Keywords:* high-performance computer; mass data processing; programming model; design

## 1 Introduction

With the development of information technology, the current society is featured by "information explosion". The development of network information processing is accelerating and the scope of data is increasing. In general, the present network information demonstrates the following characteristics: First, the quantity of data is quite huge and keeps increasing with the passage of time. Second, data do not remain unchanged. They are changing all the time. Third, there is no specific pattern and order. Fourth, it manifests certain insecurity. Fifth, it shows strong ability of guiding public opinions. In all, the requirements on application proposed by people's production and life exceed the current computer platform's performances[1-3]. Therefore, information technology developers are searching for a greater platform so as to realize the processing of mass data and satisfy people's requirements. Under the urge of the time, cloud computer technology begins to enter people's sight and shows its superiority in the field of computer information.

At present, the research on and application of cloud computer technology is mainly dominated by Google, Yahoo, Facebook and other large companies that provide outstanding products and services for customers by using the technology to process mass information data, realize data processing or satisfy the application demands, optimize the application services and so on. These large companies all choose to use the framework of MapReduce to set up a platform for processing customers' mass information[4-7]. The paper mainly explores into the simpler and more effective method of processing mass data on the high-performance computer platform, carries out relevant test and evaluation through applying MapReduce system, and analyzes the major factors in processing mass data effectively based on high-performance computers.

## 2 Platform of processing mass network information data

The platform of mass data processing based on high-performance computer mainly aims to conduct more in-depth analysis on and processing of mass information data through high-performance computer's technology, and improve the ability of finding, processing and analyzing information. Through excavation of parallel data, construction of a better-improved information processing platform, agglomerated high-speed inter-linkage, middleware of unified views and other technology, the goals can be well realized.

The platform of mass data processing is composed of the level of obtaining, organizing, storing and integrating data and customers' interface level. The core part is the module of data excavation. The platform provides statistics, analysis and excavation for mass information data processing.

In terms of obtaining data, the information that passes by the network is connected to the processing platform[8-10]. After screening and integration, the information enters corresponding organizational level. In terms of organizing data, the obtained mass information data are processed online, including extracting texts and features of information, identifying information, rapid scanning and classifying, etc. In terms of storing data, it refers to implementing parallel processing of the data with three-level granularity by making use of distributed parallel DB structure and then storing the organized data through the middleware in the unified views. In terms of integrating data, it refers to extracting the original data in DB, organizing mass data through using multi-dimensional data models, and providing checking service for development personnel. The following diagram shows the overall systematic structure of high-performance computer's platform of mass data processing.
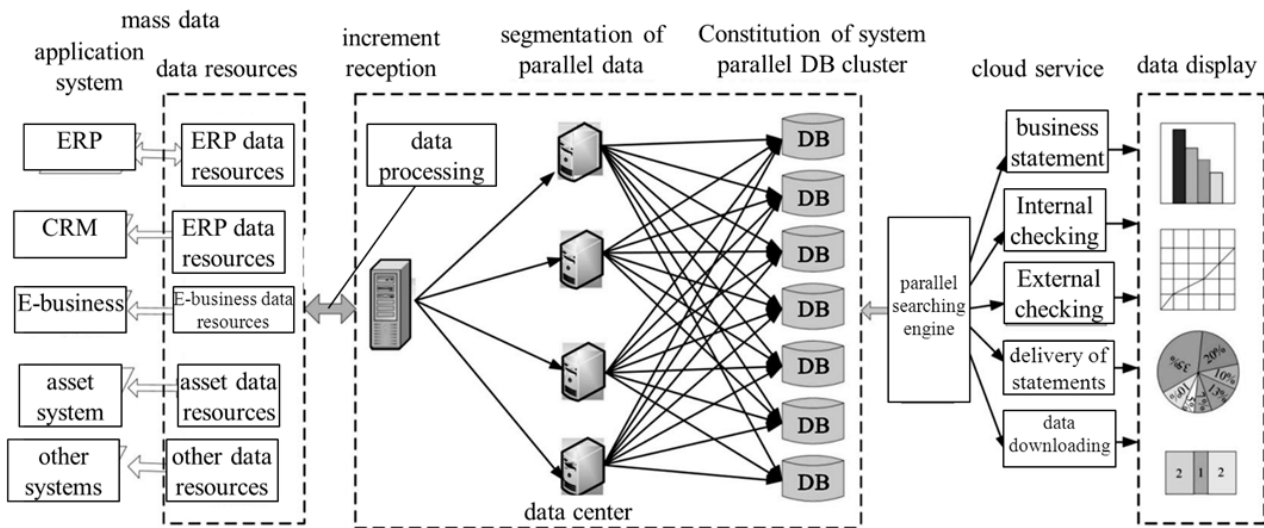
Figure 1 Overall Framework of High-performance Computer

## 3 MapReduce programming model

The paper mainly analyzes the structure and programming model of Apache Hadoop – the source version of MapReduce. Apache Hadoop is composed of two major parts, including MapReduce programming model and HDFS. Each node in HDFS has local storage which integrates the nodes into a very large distributed file system that can provide overall views for different nodes. Generally speaking, HDFS requires relatively large data block storage space for storing files. Additionally, maintenance measures should be taken to ensure the system's security and reliability.

In MapReduce model, Map can process data, extract the information data needed by the system and form the middleware value. In Hadoop model, Reduce provides convenience for the system through the tasks at three stages. First of all, it is shuffle, the stage of extracting the intermediate results from the nodes. Second, it is merge, the stage of combining and processing the values. Finally, it is Reduce, the stage of further integrating the processed results and getting the final data after processing.

## 4 Construction of MapReduce model in high-performance computer system

### 4.1 FEATURES OF MAPREDUCE MODEL

At present, although high-performance computers possess stronger calculation and processing ability, MapReduce model applied in large international enterprises are still being operated in low-end servers. However, the usage of centralized storage sub-systems distinguishes the utilization of MapReduce model in high-performance computers and the utilization in ordinary cluster.

From the comparative analysis, we can know that the systematic cluster file system used in high-performance computers is Lustre, which can provide a complete overall view for computers and overcome the deficiencies of HDFS. However, at present, in the application of the sys-

tem in enterprises, HDFS is still maintained because Lustre cannot realize the function of getting access to data like HDFS interface.
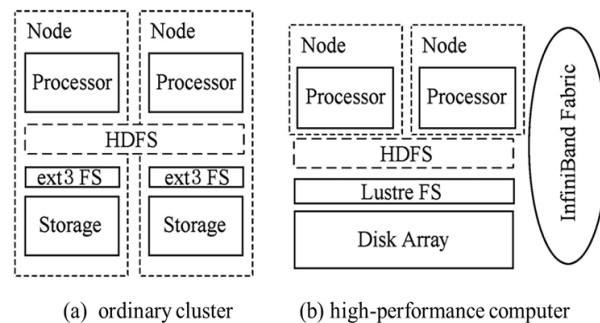


FIGURE 2 Comparison between Ordinary Cluster and High-performance Computer

### 4.2 OPTIMAL DESIGN OF MAPREDUCE MODEL

Pertinent to the application of MapReduce in high-performance computers, the information data are stored in sharing cluster file system instead of carrying out relevant optimal design of the problems in the local disks. Take the ordinary MapReduce as example, first, it has to undergo the first stage, shuffle, i.e., the intermediate result generated in the map requires inter-node network transmission in the process of transmitting to different nodes. Therefore, a lot of network broad band is occupied. However, when the data are stored in cluster file system, what can be seen at different nodes are the same names and space. The time of data transmission in the network can be ignored. It is very easy and convenient to read the data, which saves the quantity of data transmission to a large degree. In order to realize the advantage, the applied HDFS should be revised to realize cluster file system. In the research in the paper, the application of cluster file system is realized and relevant evaluation on the performances is made.

112

## 5 Experiment of high-performance computer's mass data processing

In order to realize the processing of mass data on high-performance computer, the performance of MapReduce system should be detected through experiments. The link and factors that affect MapReduce should be found out under the background of high-performance computers. First of all, the paper predicts the factors that may cause the problems in MapReduce system. First, the calculation speed of the computer's processor. Second, the writing and reading speed of the storage system in high-performance computer and the speed of transmission of network data. Pertinent to the problems caused by the three factors, the paper proposes corresponding solutions.

### 5.1 EXPERIMENTAL BACKGROUND

1. Configuration of high-performance computer: 100 nodes are configured with two-way and six-core processor with the frequency of 2.93GHz. The capacity of main storage is 50GB. Different nodes, nodes and storage systems share the broadband and realize network interconnection. The storage sub-systems are composed of 170 600GB pipeline disk arranges. The file management system has 180 disks in total which are divided into 30 DAID6 logical volumes that are under management respectively. The evaluation method is IOzone. The total speed of reading and writing of the systematic process is respectively 36.97GB/s and 1.12GB/s.

### 5.2 EXPERIMENTAL RESULT

After the complete experiment, the paper evaluates the performance of processing mass data in high-performance computers by applying MapReduce system. After high-performance computers apply MapReduce, the data are put into cluster file system through distributed file system. When Terasort programme is operated, the 100GB data are put in order and the test result is obtained. Then, by making use of the HDFS and cluster file system that are optimized, optimize the data transmission at the shuffle stage. Then, sequence the 100GB data again to get the corresponding test results. Finally, on the basis of not using the cluster file system, put the data into HDFS in the local disks directly, sequence the 100GB data, and get the final testing result. The result shows that the quantity of map tasks and quantity of reduce tasks are changing with the time. The completion time of the whole task under different environment also varies. The quantity of communications at the calculation nodes and utilization ratio of high-performance computers' CPU are also different.

### 5.3 EXPERIMENTAL CONCLUSIONS

According to the experiment above, the precision of the prediction before the experiment cannot be concluded. In general, the main factor affecting mass data processing in high-performance computers is the writing and reading speed of the storage system. The calculation speed of the computer's processor and transmission speed of network data are not major factors affecting the effectiveness. Through further comparison and analysis, when there are a lot of map tasks and reduce tasks, the occurrence of con-

current competition caused by a number of programme operation reduces the performances of the cluster file system – Lustre, which is far from the evaluation on specific storage sub-system. Moreover, it is even found that when the cluster file system, Luster is operated, errors will occur when the number of operating programmes increases. Such a result has also appeared in other research reports. At current, it is generally agreed that the reduction of the cluster file system – Lustre's performances is caused by the too fierce competition of tasks for resources. In the cluster file system, the file data are stored in the logical volume in the form of band type. The ability of parallel reading and writing also effectively improves in the process of getting access to the data.

In the experiment, about 1MB cluster file system is striped, while 64MB distributed file system is striped, which indicates that when the block of one distributed file system is put into the cluster file system, it will occupy 20 storage targets. In the experiment, MapReduce system assigns tasks for each distributed file system, which implies that all tasks will be exchanged with OST in the system. This is also the main reason for the competition. Finally, the performances of the cluster file system reduce. It fails to give its functions into the maximal play.

As a result, in terms of processing mass data in high-performance computers, the effectiveness and reliability of data processing should be ensured. The competition for resources in the cluster file system must be solved so as to improve the system's data processing speed. The paper proposes three different solutions. First of all, when intensive application processing is carried out, the scope of storage sub-system should be expanded so as to provide adequate services for the cluster file system. Second, pertinent to the problem that the striped block of the cluster file system has a small capacity, it can be increased to the size of the distributed file system, i.e., 64MB. Third, abandon the existence of the distributed file system. Use the MapReduce system as the cluster file system to provide services, which can save the block's reflection process and better solve problems compared to the other two methods.

## 6 Conclusions

Above all, the paper analyzes and explores into the application of the platform of mass data processing in high-performance computers, proposes the optimal design of cluster file system, analyzes the possibility of the realization of MapReduce in high-performance computers, and affects the factors affecting the performances of MapReduce through relevant experiment. With the three plans, the problem is solved. The cluster file system is also facilitated to give its functions into the maximal play.

## References

[1] Guo Zhi-Liang, Gao Chun-Hai1, Ma Lian-Chuan, Lü Ji-Dong 2011 Formal verification of safety computer platform based on timed automata model *Tiedao Xuebao/Journal of the China Railway Society* **33**(6) 68-73

[2] Tan Huai-Liang, He, Zai-Hong 2006 A TCP implementation on an embedded bare computer platform *Hunan Daxue Xuebao/Journal of Hunan University Natural Sciences* **33**(3) 119-23

[3] An Peng, Shao Beibei, Zhang Jian 2009 High reliability computer platform using quadruple modular redundancy *Qinghua Daxue Xuebao/Journal of Tsinghua University* **49**(11) 1737-40

[4] Wu Hongling, Zeng Xiaofei 2011 Application of computer technology in efficiency analysis of China Life Insurance Company *Journal of Computers* **6**(9) 1832-41

[5] Deng Yonggang, Tang Chenghao, Sun Gangqiang, Huang Sheng 2005 The application of the computer simulation technology in derrick detection *Drilling and Production Technology* **28**(3) 82-84+5

[6] Shi Jianyong, Li Yinqing, Chen Huchuan 2008 Application of Computer Integration Technology for Fire Safety Analysis *Tsinghua Science and Technology* **13**(SUPPL.1) 387-92

[7] Li Ni, Chen Zheng, Gong Guang-Hong, Peng Xiao-Yuan 2010 Application of multi-core parallel computing technology in scene matching simulation *Xi Tong Gong Cheng Yu Dian Zi Ji Shu/Systems Engineering and Electronics* **32**(2) 428-32

[8] Kinjo Keita, Aizawa Akiko, Ozaki, Tomonobu 2010 Extraction of social structure change from survey data: Application of inductive logic programming for cohort network analysis T*ransactions of the Japanese Society for Artificial Intelligence* **25**(3) 452-63

[9] Om Romny, Kunleang Heng 2005 A survey on data model for obtaining power by using geographic information system (GIS) *International Energy Journal* **6**(1) 2139-46

[10] Kao Chiang 2014 Network data envelopment analysis: A review European Journal of Operational Research **239**(1) 1-16

**Authors**

**Jing Nie, born in 1984, Nanning, Guangxi Zhuang Autonomous Region, P.R. China**

**Current position, grades:** Master degree  the lecturer of Nanning College of  Vocational Technology, China
**Scientific interest:** mass data, computer platform design
**Publications:** more than 10 papers
**Experience:** teaching experience of 8 years, has completed three scientific research projects

Information and Computer Technologies