# The Optimization of Task Assignments on Hadoop Platform for Large-Number Image Processing

G.R. Zhang

Key Laboratory of OptoElectronic Science and Technology
for Medicine of Ministry of Education
College of Photonic and Electronic Engineering,
Fujian Normal University
Fuzhou, China

L.P. Huang

Key Laboratory of OptoElectronic Science and Technology
for Medicine of Ministry of Education
College of Photonic and Electronic Engineering
Fujian Normal University
Fuzhou, China

Q.X. Wu

Key Laboratory of OptoElectronic Science and Technology
for Medicine of Ministry of Education
College of Photonic and Electronic Engineering
Fujian Normal University
Fuzhou, China

B.S. Chen

Key Laboratory of OptoElectronic Science and Technology
for Medicine of Ministry of Education
College of Photonic and Electronic Engineering
Fujian Normal University
Fuzhou, China

*Abstract*——**This paper proposes a large-number images processing model based on Hadoop platform. According to this model, the task allocation strategies are investigated for large-number image processing on the isomorphic Hadoop clusters and heterogeneous Hadoop clusters. Firstly, a series of different experiments is performed on the isomorphic Hadoop clusters. The obtained results show that the equal allocation of all images to each node is proved to be an efficient approach. For the heterogeneous clusters, the genetic algorithm is used as a task allocation strategy to optimize the processing task allocation for a large-number of images. Experimental results show that the GA-based optimization can significantly speed up the image processing so that the proposed approach is promising to apply to process big data of images.**

*Keywords-Hadoop cluster; Image processing; Hadoop streaming; genetic algorithm*

## I. INTRODUCTION

As images from surveillance systems and medical image domain become a source of big data, more and more scientists pay an attention to cloud computing for image processing. Hadoop is an open source distributed computing platform [1]. Derived from an implementation of Google's MapReduce [2], Hadoop consists of two chief components: Hadoop MapReduce and Hadoop Distributed File System (HDFS). HDFS and MapReduce as the core of Hadoop provide users with transparent bottom details of the distributed system infrastructure. Hadoop clusters' distributed storage and distributed computing have high reliability, expandability and efficiency. The architecture of Hadoop file system is suitable for storing and managing large volume of medical images [3]. Hadoop is also used in massive image retrieval works. The massive image retrieval system based on Hadoop in dealing with large data retrieval, compared with the single-node retrieval system, can effectively reduce the search time, and improve the retrieval speed [4].

This paper is based on the proposed large-number image processing model on Hadoop platform. Through the model, we can process large-number images efficiently compared with single node. At first, we do some experiments on the isomorphic Hadoop cluster, and get some results. Besides, as the Hadoop cluster nodes can be heterogeneous, thus the computing speed of each node may different. If the huge number images to be processed are still equably assigned to the heterogeneous cluster, the slowest node in the computing will determine the entire job processing time. In order to save the task operating time, this paper presents task allocation strategy based on genetic algorithm. And using the strategy in the heterogeneous cluster can make the cost time of cluster nodes convergence, save the entire job time.

Genetic Algorithm (GA) is a simulation mechanism of Darwinian's natural selection and genetics computational models of biological evolution process of biological evolution, the process is a way to search the optimal solution by simulating natural evolution, which was originally proposed by J.Holland professor from Michigan University at 1975 [5]. GA is begin with the representation of a population containing potential solutions, and then a certain number of individuals via genes coding consist of the population, each individual is actually entity with chromosomes characteristics [6]. GA is a constantly iterative process [7]. After the initial population generated, the whole population in accordance with the rules of survival of the fittest, through continuous genetic and evolution, eventually produce a near optimal solution of the problem [8, 9].

## II. THE EXPERIMENTAL PLATFORM

Due to the limited experimental conditions, this paper is simulated on virtual machines.

### A. Simulation Using Virtual Machines

In this paper, the platform is composed of six machines, including one NameNode and five DataNode. The 6-nodes

cluster configuration information is shown in the following tables:

TABLE I.    HARDWARE CONFIGURATION

| Name | Amount(num) | Detail Configuration |
|------|-------------|----------------------|
| NameNode | 1 | 1CPU*2.0GHz 4GB RAM |
| DataNode | 5 | 1CPU*2.0GHz 1GB RAM |
| Disk space | | 20GB |

TABLE II.    SOFTWARE CONFIGURATION

| Software Name | Version |
|---------------|---------|
| Red Hat AS | CentOS 6.3 |
| JDK | Jdk-6u31-linux-i586 |
| Hadoop | 1.0.0 |

TABLE III.    NET WORK CONFIGURATION

| Machine name | Detail configuration |
|--------------|----------------------|
| Master | 192.168.11.155 |
| Slave1~Slave5 | 192.168.11.156~192.168.11.160 |

The Hadoop cluster is isomorphic. The configurations of the DataNode nodes are all the same. This cluster is used to do the first experiment get a best task assignment strategy under the isomorphic Hadoop cluster.

## III.    THE CORE TECHNOLOGY USED IN THE PROGRAM

This program uses the Hadoop Streaming technology. Hadoop Streaming technology can help users to create and run a special kind of map/reduce jobs. These special jobs can be performed through an executable file or script file which acts as a mapper or reducer. The program is to use a shell script file to act as a mapper implementation. The shell script can call executable files of image processing. The image files are used as input data for the executable file to process and then the processing results are uploaded to file system HDFS [10].

In the program, image files to be processed should be uploaded to the HDFS at first, and then a filelist is made in a specified directory of Hadoop Streaming's input. The directory contains a set of filelist, files in which the contents are the HDFS paths of image files to be processed, and each line in the file is an image path. Inputsplit is a text file, which is regarded as the input of mapper. The shell script mapper contains the file read-line operation, read the path of the image and get an image, and then call the image processing executable file to deal with the image. Finally, the image processing result is saved back to HDFS. In general, reducer is used for combining the output from mapper, and then exports the final results. In this scenario, mapper is the shell script. The output of mapper is no longer intermediate results. Shell script already contains the operations of getting image files from HDFS, image processing, and upload the results to HDFS and other operations. The map can directly manipulate the image processing results, so the task number of reduce can be set to 0.

## IV.    THE EXPERIMENTAL RESULTS AND ANALYSIS

### A.    The experimental process on isomorphic Hadoop cluster

The experimental data are one thousand BMP format image files; each image size is 256×256. The job is to perform binary processing for all images (this job can be replaced with cancer image identification or people recognition in videos). The binary image processing is compiled by the C language. Virtual machine cluster contains six nodes. A JobTracker node is responsible for scheduling the work, and other five TaskTracker nodes are responsible for performing the calculations. The images are assigned to nodes through one file in filelist directory. The file stored 1000 lines text, each line is an image file path. The inputformat is set NLlineInputFormat in the experiment (default value of N is 1), which represents input by line. The document file is the input of job; we can change the size of N by setting the value of linespermap. For instance, if we set linespermap=2, it indicates that every two image files' path severed as an input slice which will be allocated to the idle node. It also represents that each map has to process two image files. In this paper, we use the same dataset of the images to do the experiments. Each time we set a different value of linesmap. Analyzed the experimental results, we obtained a better assignment strategy. The experimental results are shown in the following table and figure:

TABLE IV.    THE EXPERIMENTAL RESULTS

| Linespermap | Amount of map(num) | Job time(second) |
|-------------|--------------------|------------------|
| 2 | 500 | 790 |
| 4 | 250 | 657 |
| 8 | 125 | 613 |
| 16 | 63 | 597 |
| 32 | 32 | 593 |
| 64 | 16 | 582 |
| 100 | 10 | 540 |
| 200 | 5 | 527 |

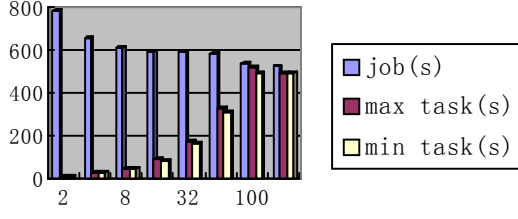| Linespermap | Max task time(second) | Min task time(second) |
|-------------|-----------------------|-----------------------|
| 2 | 18 | 15 |
| 4 | 33 | 31 |
| 8 | 54 | 50 |
| 16 | 93 | 84 |
| 32 | 177 | 171 |
| 64 | 336 | 315 |
| 100 | 522 | 498 |
| 200 | 498 | 496 |

FIGURE I. THE EXPERIMENTAL RESULTS

To make full use of the resources of five DataNode, the number of map tasks is not less than the amount of the calculated nodes. Experimental results show that, when the value of linespermap is small, there are a lot of map tasks. And each map task process small amounts of images so as to the image processing accounted for only a small proportion time of the whole task. It means that map task assignment and its start time take up most of the job time, which makes the job inefficient. While increasing the linespermap, the number of map is decreasing, and the job time is shrinking. When the value of linespermap is 200, the job's finished time is shortest. At this time, the job has five map tasks, each slave machine is allocated one, and each task is gonging to process 200 images. The result indicates that when dealing with large quantities of images on the isomorphic Hadoop cluster, equably assigned to the cluster is more efficient than the way divided the task too thin.

### B. Optimization of Task Allocation By Ga on Heterogeneous Hadoop Cluster

As the Hadoop architecture is succinct. Using the common computers can be built as Hadoop clusters. It can be that some nodes in the cluster have better calculate ability than the others. It means that if the task equably assigned to the cluster's DataNode, the job may not be efficient.

In order to verify that the GA can optimize task allocation, we build a heterogeneous Hadoop cluster. The cluster also includes 5 DataNode nodes, the hardware configuration of these nodes are not exactly the same. The cluster configuration information is shown in the following tables:

TABLE V. CLUSTER HARDWARE CONFIGURATION

| DataNode Name | Detail Configuration |
|---|---|
| Slave1~Slave4 | 2CPU*2.0GHz 1GB RAM |
| Slave5 | 1CPU*2.0GHz 1GB RAM |

Wherein, Slave1~Slave4 contain two processor, Slave5 has only one processor. The rest of the configuration is identical. The experimental data is 1000 BMP format image files which are the same as former experiment. The job is also to perform binary processing for all images. In this experiment, filelist directory has five files, each file contains 200 images' path respectively. The task is equably assigned to the heterogeneous cluster. The experimental result before optimization is showing in the following table. The average cost time of 5 DataNode process an image can be calculated from the experiment.

TABLE VI. THE EXPERIMENTAL RESULTS BEFORE OPTIMIZATION

| Name | The amount of images(num) | Cost time(second) | Average time per image(second) |
|---|---|---|---|
| Slave1 | 200 | 445 | 2.23 |
| Slave2 | 200 | 448 | 2.24 |
| Slave3 | 200 | 448 | 2.24 |
| Slave4 | 200 | 448 | 2.24 |
| Slave5 | 200 | 544 | 2.72 |

According to the Hadoop Map/Reduce Administration pages, to complete the job costs 559 seconds totally in this experiment.

This paper uses genetic algorithm to optimize the allocation. The encoding mode is using unsigned binary integer to represent the number of images distributed to the 5 DataNode. For example, the assignments in the table x = [200, 200, 200, 200, 200] will be coded as a chromosome X = [0011001000 0011001000 0011001000 0011001000 0011001000]. The phenotype x and genotype X can convert to each other by the encoding and decoding process. In the algorithm, the fitness function f(x) shows below:

$$f(x) = \frac{1}{\sum_{i=1}^{N} a_i x_i} \qquad (1)$$

In the formula, N represents the amount of DataNode, $a_i$ is the average cost time of the DataNode processing an image, $x_i$ means the amount of pictures assigned to node i for processing.

In this experiment, the initial population is generated by the programs, then after the operation of selection, crossover and mutation, we get the optimized chromosome:

X = [0011011001 0011001001 0011000110 0011001000 0010111000]

Convert to phenotype x = [217, 201, 198, 200, 184].

After the optimization of GA, the task allocation strategy assign the number of images to each slave node are: 217, 201, 198, 200, 184. The optimization result is shown in the following table:

TABLE VII. THE OPTIMIZATION RESULT

| Name | The amount of images(num) | Cost time(second) |
|---|---|---|
| Slave1 | 217 | 468 |
| Slave2 | 201 | 451 |
| Slave3 | 198 | 445 |
| Slave4 | 200 | 450 |
| Slave5 | 184 | 474 |

The job finished in 486 seconds in this experiment.

Experimental results show that in the heterogeneous cluster, using the optimization of GA can make the processing time of nodes more balanced, improve the problem of fast nodes waiting for the slow one a long time to complete the job. And the job time saves 13% compared with unoptimized.

## V. Summary and Future Work

This paper is based on the proposed large-number images processing model on Hadoop platform. We study the task allocation strategy for large-number image processing on the isomorphic Hadoop cluster and heterogeneous Hadoop cluster. Firstly, do experiments on the isomorphic Hadoop cluster. The results show that the average allocation of all images to each node is more efficient. As the heterogeneous cluster, this paper presents a task allocation strategy using GA to optimize the allocation of large-number images processing. Experimental results show that the optimization can significantly speed up the image processing.

## References

[1] T. White, Hadoop: The Definitive Guide. O'Reilly Media, Yahoo! Press, pp. 9-11, 2009.

[2] F.N. Afrati, J.D. Ullman. Optimizing multiway joins in a Map-Reduce environment. IEEE Transactions on Knowledge and Data Engineering, Vol.23, No.9, pp. 1282–1298, 2011.

[3] P.J. LI, G.J. CHEN, W.M. GUO. A distributed storage architecture for regional medical image sharing and cooperation based on HDFS. J South Med Univ, Vol.31, No.3, pp. 495-498, 2011.

[4] M. WANG, X.Z. ZHU, J.M. ZHAO, C.F. HUANG. Massive Images Retrieval System Based on Hadoop. Computer Technology and Development. Vol.23, No.1, pp. 204-208, 2013.

[5] HOLLAND J H. Adaptation in natural and artificial systems [M]. Ann Arbor: University of Michigan Press, 1975.

[6] M. Zhou, S.D. Sun, Genetic Algorithm Theory and Applications [M]. Beijing. National Defence Industry Press, 1999.

[7] YOGESWARAN M, PONNAMBALAM S G, TIWARI M K. An efficient hybrid evolutionary heuristic using genetic algorithm and simulated annealing algorithm to solve machine loading problem in FMS. International Journal of Production Research. 2009.

[8] J.F. GONCALVES. et al, A Hybrid Genetic Algorithm for Assembly Line Balancing [J]. Journal of Heuristics, pp. 629-642, 2002.

[9] FILLIAT D, MEYER J A. Map-based navigation in mobile robots: I. A review of localization strategies [J]. Cognitive Systems Research, 4(4), pp. 243-282, 2003.

[10] G.R. Zhang, Q.X. Wu, Z.Q Zhuo. A Large-scale Images Processing Model based on Hadoop Platform. International Conference on Innovative Computing and Cloud Computing (ICCC), pp. 51-54, 2013.