

## RESEARCH ARTICLE

# A Speculative Approach to Spatial-Temporal Efficiency with Multi-Objective Optimisation in a Heterogeneous Cloud Environment

Qi Liu<sup>1</sup>, Weidong Cai<sup>1</sup>, Jian Shen<sup>2</sup>, Zhangjie Fu<sup>3\*</sup>, Xiaodong Liu<sup>4</sup>, and Nigel Linge<sup>5</sup>

<sup>1</sup>School of Computer and Software, Nanjing University of Information Science and Technology, Nanjing, China

<sup>2</sup>Jiangsu Engineering Centre for Network Monitoring, Nanjing University of Information Science and Technology, Nanjing, China

<sup>3</sup>Jiangsu Collaborative Innovation Center of Atmospheric Environment and Equipment Technology (CICAEET), Nanjing University of Information Science and Technology, Nanjing, China

<sup>4</sup>School of Computing, Edinburgh Napier University, 10 Colinton Road, Edinburgh EH10 5DT, UK

<sup>5</sup>School of Computing Science and Engineering, University of Salford, Salford, UK

## ABSTRACT

A heterogeneous cloud system, e.g. a Hadoop 2.6.0 platform provides distributed but cohesive services with rich features on large-scale management, reliability and error tolerance. As big data processing is concerned, newly built cloud clusters meet the challenges of performance optimisation focusing on faster task execution and more efficient usage of computing resources. Present proposed approaches concentrate on temporal improvement, i.e. shortening MapReduce (MR) time, but seldom focus on storage occupation; however, unbalanced cloud storage strategies could exhaust those nodes with heavy Map/Reduce cycles, and further challenge the security and stability of the entire cluster. In this paper, an adaptive method is presented aiming at spatial-temporal efficiency in a heterogeneous cloud environment. A prediction model based on an optimised K-ELM algorithm is proposed for faster forecast of job execution duration and space occupation, which consequently facilitates the process of task scheduling through a multi-objective algorithm called TS-NSGA-II. Experiment results have shown that compared to the original load-balancing scheme, our approach can save approximate 47-55 seconds averagely on each task execution. Simultaneously, 1.254% of differences on hard disk occupation were made among all scheduled reducers, which achieves 26.6% improvement over the original scheme. Copyright © 2015 John Wiley & Sons, Ltd.

## KEYWORDS

MapReduce; Cloud Storage; Load Balancing; Multi-Objective Optimisation; Prediction Model

### \*Correspondence

Zhangjie Fu, School of Computer and Software, Nanjing University of Information Science and Technology Nanjing, China. E-mail: [wwwfzj@126.com](mailto:wwwfzj@126.com)

Received . . .

## 1. INTRODUCTION

In recent years, distributed computing has been widely investigated and deployed in both academic and industrial fields due to its features of large-scale, virtualization, failure control among connected components, and asynchronised communication. Cloud computing as one of the successful commercial distributed systems provides users

with on-demand services by allocating rational computing and storage resources transparently [1, 2].

MapReduce paradigm proposed by Google is being exploited by a fast growing number of companies and research institutes [3]. Hadoop, as a type of open-source implementation provided by Apache, gives them a good chance to conduct efficient big data processing and discover potential and valuable information in a non-traditional way. Enterprises and companies therefore benefit from analysing and dealing with real-time data. At the moment, data analysis applications in a cloud have shown different complexity, resource requirements and data delivery deadlines; such diversity has created new

† Please ensure that you use the most up to date class file, available from the SEC Home Page at

<http://www3.interscience.wiley.com/journal/114299116/home>

requirements of job scheduling, workload management and program design in a cloud. Several projects have been launched to reduce challenges on writing complex programs for data analysis and/or data mining, e.g. Pig [4] built upon the MapReduce engine in the Hadoop environment. In addition, HBase [5] and Hive [6], implemented by Apache, are widely used in a cloud environment to achieve better performance. In these applications, however, low-level improvement based on MapReduce is still required due to its direct interaction with HDFS (Hadoop Distributed File System) [7]. An outstanding strategy that improves the security and stability of a cloud system is necessary.

While the optimization of job scheduling in MapReduce has been widely conducted in recent activities [8-20], current Hadoop systems still suffer from poor load-scheduling strategies due to their lack of consideration on the usage of cloud storage, which would bring heavy loads on certain data nodes and therefore cause a long delay on total execution. Although theoretically infinite computing resources can be provided in a cloud system, unreasonable increment of mappers/reducers cannot achieve processing efficiency, and even waste more storage to complete.

A scheme is therefore presented in this paper to achieve process efficiency and load balance in a cloud system both spatially and temporally. Our contributions are three folds as follows:

- (1) A prediction model called PMK-ELM is firstly built providing prediction on the number of reducers needed for newly coming tasks, as well as possible execution duration and storage size they may take.
- (2) An optimized algorithm based on NSGA-II [21] called TS-NSGA-II is then designed to maintain such an equalized status that the total time completing the job distributed in each reducer is almost same while keeping the ratio of hard disk space similar.
- (3) A practical Hadoop environment is constructed to verify the feasibility and performance of the scheme.

The remainder of this paper is organized into five sections. Related work on load balancing is reviewed in Section II. In Section III, preliminaries of core algorithms manipulated in our approach are introduced. Section IV explains the adaptive method to achieve fair loads during map and reduce processes. Results are presented and evaluated in Section V with a comparison of corresponding algorithms. Finally, Section VI concludes the paper and identifies potential future work.

## 2. RELATED WORK

A balanced load is hard to be achieved due to the imbalanced input data of Reduce phase. [9] proposed an optimization method, by repartitioning the inputting data

of map and reduce tasks, all available data node can complete its task at the same time. This method can handle all kinds of load deflections, but it is too difficult to be implemented, and it has greatly changed Hadoop. Also, extra reassigning tasks would produce additional network overhead. Partition methods are also research hot pots, such as those methods based on historical data [10] and the sampling results [11], which could allocate input data to different nodes more flexibly. Though these methods can achieve dynamic load balancing, their performance system was not verified in an actual Hadoop system. Through offline and online analysis, resource requirements can be predicted by using a benchmark or real application workloads, for example, [12] proposed a prediction model based on SVM in a heterogeneous environment. Combined with an adaptive algorithm HAP, it can be used for predicting the amount of data assigned to different tasks node. However, the reduce tasks required repeated cutting and consolidation of data blocks, which can lead to extra time cost. In addition, the training phase of SVM would require much time. A prediction model focused on resource consumption of MapReduce processes, based on a classification and regression tree, was presented by Jing et al. [16].

The efficiency of virtualization deployment has been extensively studied. [13] proposed a general method for estimating resource requirements when running applications in a virtual environment. [14] studied the resource requirement of starting a new virtual instance. Through a resource prediction model, dynamic resource provision was achieved in a cloud environment. Metrics for performance and load efficiency assessment in cloud systems have also gained much attention. [15] described a method for more accurate assessment of distributed cloud application performance.

Besides the above methods, some researchers are studying optimizing the speculative execution strategy in MapReduce. A key advantage of MapReduce is its automatic processing failure. Its high fault tolerance makes it easier for a programmer to use. If a node collapses, MapReduce will restart the task on different machines. Some speculative execution strategies have been proposed in some literature. Google only started backup tasks when a job was close to completion; their experiments showed that proposed speculate execution can reduce the execution time of operation 44% [3]. In order to improve the performance of the cluster, Hadoop and Microsoft Dryad [31] also provided an implementation speculative execution strategy. At first, their strategy was roughly the same as that proposed by Google. However, an optimized speculative execution called Longest Approximate Time to End (LATE) algorithm was proposed in which a different metric was defined to start tasks for speculative execution. The remaining time was estimated, rather than considering the progress of the current task. LATE gave a more clear assessment of struggling tasks' impacts on the overall job running time. But the time that

every stage occupies was not stable while the std representing standard deviation used in LATE cannot be applied in all applications. Qi et al. therefore proposed MCP to overcome the disadvantages in LATE. MCP identified slow tasks based on average progress rate of a whole cluster though in reality, the progress rate can be unstable. Struggles can be appropriately judged in homogenous environments. However, there are still a lot of disadvantages in MCP, including average progress rate and its mediocre performance in heterogeneous environments.

Data placement schemes have also been researched. To address this problem, a new Data-gRouping-Aware data placement scheme was proposed in [19]. It extracts optimal data groupings and re-organizes data layouts to achieve load balancing in per group. CoHadoop was proposed in [20]; it permits applications to decide where data should be stored. However, these schemes are aimed at the data placement when storing the data and not fit for MapReduce. Furthermore, they cannot be applied when data have been stored.

Comprehensive load and usage efficiency have achieved large improvement in a distributed environment. However, it is still challenging to achieve spatial-temporal efficiency in a cloud system, especially in a heterogeneous one.

### 3. PRELIMINARIES

A detailed introduction to some advanced techniques used in this paper is given in this section.

#### 3.1. MapReduce

In MapReduce, computation works are implemented through map tasks and reduce tasks. Map tasks put different pairs of data into multiple lists grouped by different keys. So, data having the same key are distributed to the same list. Then, results generated by map tasks, as intermediate data, are pulled by reduce tasks to process further and get the final results [22].

MapReduce jobs are divided into multiple tasks, then, these generated tasks are distributed to nodes and executed in the cluster. Map tasks are partitioned into different datanodes according to a logical split of input data that generally resides on HDFS [23]. Reduce tasks are produced according to an equation in reduce stage. The map task reads the data from HDFS as input data, map functions designed by the user are then applied and put the results into buffers. This data are written to the memory of the node executing the map task when it is less than the threshold user set. Otherwise, this data will be spilled into the hard disk of the nodes. There are three phases in reduce tasks, called shuffle (copy), sort (merge), and reduce. In the shuffle phase, the reduce tasks pull the intermediate data files generated by the map tasks. Then, the intermediate files from all the map tasks are sorted in the following phase. After all the intermediate data are shuffled and transferred, the reduce phase starts working.

Job scheduling in Hadoop is performed by the namenode, which manages a number of datanodes in the cluster. In MapReduce 2.x, each datanode will prepare containers for map tasks and reduce tasks, which can be seen as an abstraction of resource and used to execute the task. The number of map and reduce container is calculated the configuration file. Application Master periodically checks the heartbeats coming from datanodes and calculates the reported state of free resources and current progress of tasks that they are currently executed.

#### 3.2. Basic ELM

Recently, Artificial Neural Networks (ANNs) have been widely applied in applications involving classification or function approximation [24]. However, they also suffer from low learning speed, which has become the main bottleneck when applying an ANN algorithm to practical applications. In order to overcome this drawback, many researchers explore the approximation capability of feedforward neural networks, especially in a limited training set, from the point of view of mathematics. A novel machine learning algorithm called Extreme Learning Machine (ELM) [25, 26] was therefore designed based on Single-hidden Layer Feedforward Neural networks (SLFNs) [27].

Let  $X = \{x_1, x_2, \dots, x_N | x_i \in R^D, i = 1, 2, \dots, N\}$  denote the training set with  $N$  samples,  $D$  represent dimension. Let  $Y = \{y_1, y_2, \dots, y_N | y_i \in R\}$  denote the vectorised label where column  $j$  ( $\{j = 1, 2, \dots, P\}$ ) set by 1 for class  $j$  while other columns set by 0, and  $P$  is the number of classes. Then, the model of a single layer hidden layer neural network having  $L$  hidden neurons and an activation function  $g(x)$  can be expressed as

$$\sum_{j=1}^L \beta_j \cdot g(\langle w_j, x_i \rangle + b_j) = y_i \quad (1)$$

where  $i = 1, 2, \dots, N$ ,  $w_j$  and  $\beta_j$  represents the weight vectors from inputs to hidden layer and from hidden layer to output layer, respectively,  $b_j$  is the bias of  $j$ th hidden neuron,  $g(\langle w_j, x_i \rangle + b_j)$  is the output of the  $j$ th hidden neuron with respect to the input sample  $x_i$ . Note that (1) can be rewritten in a compact form as

$$H \cdot \beta = Y' \quad (2)$$

where  $H$  is the hidden layer output matrix of SLFNs and  $\beta$  is the output weight matrix,  $Y'$  is the transpose of  $Y$ . Optimal weights and bias of SLFNs can be found by using back propagation learning algorithms, which requires users to specify learning rates and momentum. However, there is no guarantee that the global minimum error rate can be found. Thus, the learning algorithm suffers from local minima and over-training. In exploration of the approximation capability of feedforward neural networks in a finite training set, it is found that SLFNs can reach the approximate capacity at a specified error

$\varepsilon(\varepsilon > 0)$  level with the hidden layer neurons is much less than the number of training samples. And based on the minimum norm least-squares function, the weight matrix  $\beta$  in (2) can be solved by

$$\beta = H^+ \cdot Y \quad (3)$$

Where  $H^+$  is a MooreCPenrose matrix generalized inverse of matrix  $H$ .

### 3.3. K-ELM

K-ELM(Kernel-ELM) has simplified the complexity of the ELM algorithm, with the improvement of the operation speed. Meanwhile it improves the simulation precision of the algorithm and the fitting ability based on the original ELM algorithm. In K-ELM, a positive number is added to the diagonal of  $H^T H$  or  $HH^T$ , which makes the ELM algorithm more stable and present a better generalization performance [28, 29]. The prediction model established based on the training set can be described as:

Minimum value:

$$L_{PELM} = \frac{1}{2} \|\beta\|^2 + \frac{1}{2} C \sum_{i=1}^N \|\xi_i\|^2 \quad (4)$$

Constraint:

$$h(x_i)\beta = y_i^T - \xi_i^T, i = 1, 2, \dots, N \quad (5)$$

where  $\beta = [\beta_1, \beta_2, \dots, \beta_L]$  is the weight of the hidden layer outputs.  $C$  is the ridge regression parameter.  $\xi_i$  is the error vector between expected outputs and training outputs,  $h(x_i)$  is output vector of hidden neurons corresponding to the training sample  $x_i$ . Finally, the output function of ELM regression can be expressed as

$$f(x) = h(x)H^T \left( \frac{I}{C} + HH^T \right)^{-1} T = \begin{pmatrix} K(x, x_1) \\ \vdots \\ K(x, x_N) \end{pmatrix} \left( \frac{I}{C} + \Omega_{ELM} \right)^{-1} T \quad (6)$$

Similar to SVM, nuclear ELM (or kernel-based ELM, K-ELM) is not required to set the number of neurons in the hidden layer and the activation function types. Common kernel functions are shown as followed.

Linear:  $K(x_i, x_j) = x_i \cdot x_j$

Polynomial:  $K(x_i, x_j) = (x_i \cdot x_j + b)^d, b \geq 0$

RBF:  $K(x_i, x_j) = \exp(-\sigma \|x_i - x_j\|^2), \sigma > 0$

Sigmoid:  $K(x_i, x_j) = \tan(ax_i \cdot x_j + b), a > 0, b < 0$

### 3.4. NSGA-2

NSGA-II as one of the multi-objects optimization algorithms has lots of operations that are the same as those

in GA. For example in NSGA-II, the population undergoes initialization, crossover and mutation as usual. However there are three main differences:

- (1) each chromosome is sorted based on non-domination sorting into a front to obtain a fitness value;
- (2) crowding distance used to measure the diversity of the population is employed to decide the distance between individuals;
- (3) the population with the current population and current offspring (obtained by crossover and mutation) is sorted again based on the rank and the crowding distance.

After that, the best N (population size) individuals are selected to be the next generation. The main consideration in the design of the NSGA-II algorithm consists of six aspects, involving code generation, determination of the initial population, fitness evaluation, selection, crossover and mutation. Detailed procedure is shown in Figure 1.

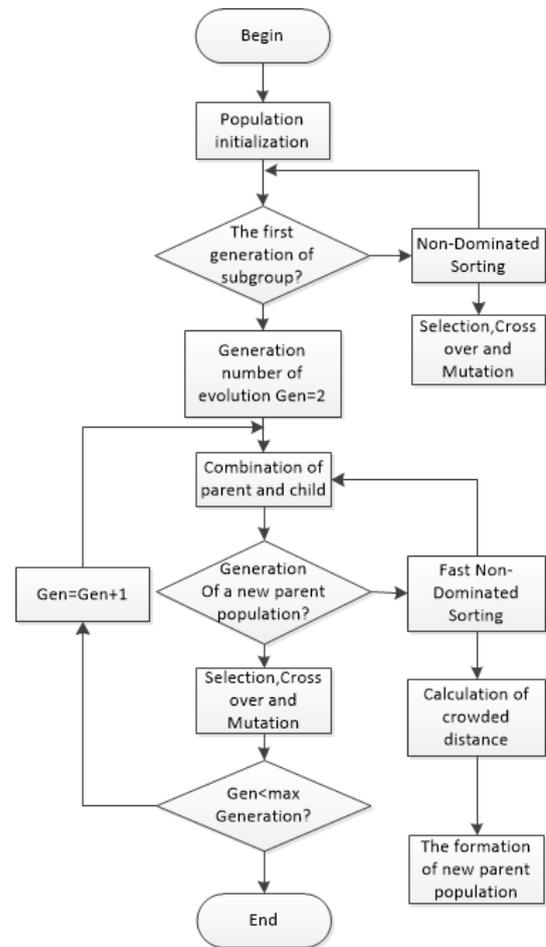


Figure 1. Flow chart of NSGA-2

## 4. APPROACH TO LOAD BALANCING

### 4.1. A method for partition reconstruction

MapReduce uses a hash function as the original partition function, where splits are generated and distributed to different reducers. The Original hash function may lead to sever load skew, especially in a heterogeneous environment, which will decrease the speed of some node. However, the overall job finishing time is decided by the node that finishes the task at last according to wooden barrel effect. Algorithm1 depicts the way that fairly equal size of splits is ensured for distribution, which helps the system dispense different volume of data to the different node having different computing capacity. Before starting the work, we run the WordCount application on each node separately to get the approximate capacity of each node. Then, the volume of data is given away according to a different capacity. According to algorithm1, the list of the partition has a relatively balanced data amount according to different capacities.

---

#### Algorithm 1 Partition Reconstruction

---

##### Input:

The input size of reduce stage, *size*;  
The number of data chunks, *number*;

##### Output: *partition\_list*

```

Get the list of capacity of each server Lc
Set iterator = 1
for iterator < number, iterator ++ do
  Get the ratio_list according to ratio_list =
  capacity/avg_capacity
  Maxr = Max(ratio_list)
  Minr = Min(ratio_list)
  if Maxr/Minr > 1.5 then
    Maxr = minr = (Maxr + Minr)/2
    Add Minr and Maxr to ratio_list
  else
    Break
  end if
end for
partition_list = size * ratio_list
return partition_list

```

---

### 4.2. A prediction model for load balancing based on K-ELM

In this section, the training set is set as:  $TS = \{time, reducer\_no, datanode\_no, input\_size, shuffle\_size\}$ , where *reducer\_no* represents the reducer number, actually, it also indicates the sequence when reducers run. *datanode\_no* represents the number of a datanode. Generally, a datanode can be mapped to several reducers. Here *input\_size* does not represent the input size of the whole task, but the input size of reducers at the reduce stage. *shuffle\_size* denotes the data size of a reducer that needs to shuffle when map processes have finished.

In details, the building progress of prediction model for execution time based on K-ELM (PMK-ELM) is as follows:

Step 1: Data pre-processing. First, samples that contain great network congestion are removed. Then the trimmed datasets are divided into training samples and test samples. The training samples are used for training the prediction model, whereas the test ones are for checking if the prediction model has been well trained.

Step 2: Model training. To build the K-ELM prediction model (PMK-ELM), training parameters of the model are obtained by using the training set sample generated by Step 1. The specific processes are as follows:

- (1) Randomly generated weights between the input layer and the hidden layer, and between the hidden layer neurons  $w$  and the threshold value  $b$ ;
- (2) Use the hidden layer neuron activation function to calculate the hidden layer output matrix  $H$ ;
- (3) Work out output layer weights.

Step 3: Data validation. Datasets generated by Step 1 are used to validate the PMK-ELM algorithm. According to the parameters trained in Step 2, the predictive values of test sets can be retrieved, which are then compared with the actual values to verify the prediction performance of the model.

### 4.3. TS-NSGA-II

#### 4.3.1. Mathematical model

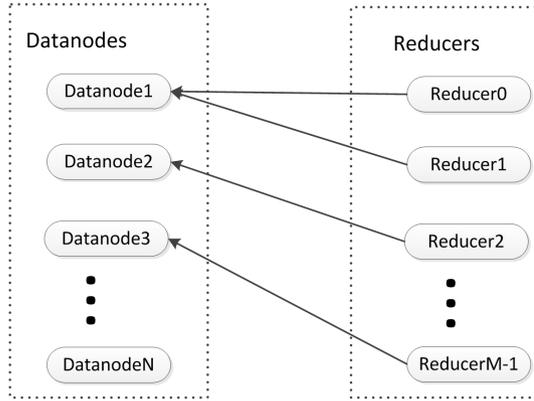
When a map task is completed, the data will be shuffled and merged, and then assigned to different reducers; however, the amount of data assigned to each reducer is not equal, which consequently causes uneven allocation of reducers to datanodes. In order to make reduce tasks consume less time and hard disk space occupation, following conditions should be satisfied:

- (1) The data amount handled by a reducer assigned to a datanode cannot be more than disk usage of the datanode;
- (2) A reducer can only be assigned to a datanode, but a datanode can handle multiple reducers, as in Figure 2.

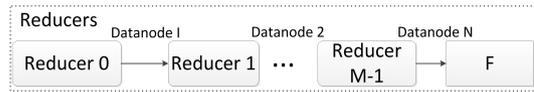
Although an actual reduce process is parallel, it is assumed in a virtual serialization line. A datanode called F is further abstracted so that when the procedure arrives at F, the reduce task is completed, as shown in Figure 3.

Assuming that the output of map tasks can be randomly divided into  $m$  data chunks and there are  $n$  datanodes in the clusters. If  $t_{mn}$  represents the execution time that each reducer needs, then the execution time of each split can be noted as a matrix  $M_t$ , as shown below:

$$M_t = \begin{pmatrix} t_{11} & \cdots & t_{1n} \\ \vdots & \ddots & \vdots \\ t_{m1} & \cdots & t_{mn} \end{pmatrix}$$



**Figure 2.** Relationship between datanodes and reducers



**Figure 3.** Virtual serialization

In order to evaluate the usage of storage space, the percentage of input size  $S_{mn}$  from total unused size  $sl_{mn}$  are calculated and noted as  $ps_{mn}$ .

$$ps_{mn} = s_{mn}/sl_{mn} \quad (7)$$

Then the hard disk space ratio of each split can be described as  $M_s$ :

$$M_s = \begin{pmatrix} ps_{11} & \cdots & ps_{1n} \\ \vdots & \ddots & \vdots \\ ps_{m1} & \cdots & ps_{mn} \end{pmatrix}$$

Finally, the elements of  $M_t$  and  $M_s$  are combined to format a new matrix  $M$  with new elements expressed as  $(t, ps)_{mn}$ , as shown below:

$$M = \begin{pmatrix} (t, ps)_{11} & \cdots & (t, ps)_{1n} \\ \vdots & \ddots & \vdots \\ (t, ps)_{m1} & \cdots & (t, ps)_{mn} \end{pmatrix}$$

The real execution time of datanode  $i$  can be described as  $t_i$ , whereas the split size can be represented as  $S_i$ . Accordingly, the real processing results list  $L$  can be calculated as:

$$L = \{(t, ps)_1, (t, ps)_2, \dots, (t, ps)_n\}$$

Here, two objective functions can be formatted as shown in (8) and (9); whereas the constraints are shown as (10) and (11), where in (10),  $InSum$  represents the total sum of reduce Input size.

$$\min T = \sum_{i=1}^n \left| \frac{\bar{t} - t_i}{\bar{t}} \right| \quad (8)$$

$$\min S = \sum_{i=1}^n |ps_i - \bar{ps}| \quad (9)$$

$$InSum = \sum_{i=1}^n s_i \quad (10)$$

$$t_i > 0, ps_i > 0. \quad (11)$$

#### 4.3.2. Design of TS-NSGA-II

The design of algorithm consists of six aspects, including determination of the initial population, fitness evaluation, selection, mutation, code generation and crossover. Major changes have been made on the latter two.

##### (1) Code generation

Non-negative integers are used as the index of reducers, i.e.  $0, 1, 2, \dots, M-1$  for  $M$  reducers, however. On the other hand,  $N$  datanodes are indexed using positive integers, i.e.  $1, 2, \dots, N$ . In this case, distribution of  $M$  reducers to  $N$  datanodes may generate  $N^M$  possible combinations.

##### (2) Crossover

The original NSGA-II algorithm uses Simulated Binary Crossover (SBX) [19] in this stage; however, in our scheme, crossover probability called  $pc$  is used for better grouping after being selected. The Crossover stage in this scheme consists of two steps:

- 1) Randomly match a group of chromosomes;
- 2) During matching chromosomes, randomly set intersections to make matched individual chromosomes exchange their information.

Chromosome should always be kept permutations, so the procedure of crossover is: after randomly selecting paired chromosomes, two crossover positions are randomly generated; the cross section of elements on the other side of the parent is also removed. Then, the new cross section is added to the sequence of the parent that has cut out some of the elements. Taking two pairs of chromosomes as an example, where chromosome A=2313|1122|32 and chromosome B = 3123|2213|12. The cross section is divided by a vertical bar. First, the element corresponding to |1122| of A is removed from B, so B' = 312312; then a gene fragment of A is added to B, so the offspring B'' is 3123|1211|22. Similarly, the offspring A'' is 2313|3222|13. For new produce offspring A'' and B'', it needs to be decided whether the total data size is bigger than the storage quota. If not, they are regarded as effective; otherwise, iteration will be operated. The complete procedure of the algorithm is shown in Algorithm 2 as follow:

**Algorithm 2** Crossover**Input:**

The list of chromosomes,  $L_i$ ;  
 Crossover probability,  $p_c$ ;  
 The hard disk space ratio of each split,  $PS_{mn}$ ;

**Output:** New list of chromosomes,  $NewLi$ 

Randomly match a group of individual in  $Li$  according to  $p_c$  noted as  $A$  and  $B$

**while true do**

Randomly generate two number not larger than the length of  $A$ , described as  $m$ ,  $n(m \leq n)$

Divide  $A$  into 3 parts:  $SeqAm, SeqAc, SeqAn$

Do the same operator to  $B$

Get  $SeqBm, SeqBc, SeqBn$

$A' = SeqAm \cup SeqAn, B' = SeqAm \cup SeqAn$

$A'' = A' \cup SeqBc, B'' = B' \cup SeqAc$

Get the  $ps$  according to  $ps_{mn}$

**if  $ps$  is smaller than 1 then**

Break

**else**

Continue

**end if****end while**

Replace  $A$  with  $A''$  and  $B$  with  $B''$  in  $NewLi$

**return**  $NewLi$

## 5. EXPERIMENT AND ANALYSIS

In order to test the performance and benefits of the load balancing scheme, a practical heterologous cloud testing environment was implemented, which consists of a desktop computer and a server. The server has 288 GB of memory and 10 TB of SATA hard disks. The desktop contains 12GB of memory, a single 500GB disk and a Core 2 Quad processor. Eight virtual machines were created in the server with different amounts of memory and number of shared processors. The detailed information is shown in Table I.

**Table I.** The detailed information of each virtual machine

NodeId	Memory(GB)	Core processors
Node1	10	8
Node2	8	4
Node3	8	1
Node4	8	8
Node5	4	8
Node6	4	4
Node7	18	4
Node8	12	8

K-means (KM) and WordCount algorithms were manipulated to evaluate the performance of load scheme. The Purdue MapReduce Benchmarks Suite provides us with the K-means clustering workload, where 26 GB of free datasets, and a free datasets of 50GB in WordCount clustering workload [30] were selected as the inputs.

All our test applications were built based on Hadoop 2.6.0. According to the Apache Hadoop documents, `mapreduce.tasktracker.reduce.tasks.maximum` has been set as 1.

Overall testing processes were conducted in three stages.

- (1) Dataset Collection. A Hadoop analysis tool was implemented to get historical data.
- (2) Execution Time Prediction. The PMK-ELM was enabled to predict the execution time of next reduce tasks.
- (3) Load balancing. The core MRContainerAllocator class was modified in the Hadoop system to apply the results generated by TS-NSGA-II.

### 5.1. Evaluation of PMK-ELM

To evaluate the performance of PMK-ELM, different input size and different numbers of reducers were tested during experiments, as depicted in Table II. SVM (PM-SVM) proposed in [12] was also replicated in the testing environment for comparison purposes. A log analysis tool was developed to collect training and test sets.

**Table II.** Experiment parameters

	Dataset size (pieces)	Training dataset size (pieces)	Testing dataset size (pieces)
K-Means	910	800	110
WordCount	800	700	100

A Generic Algorithm (GA) was employed to generate the parameters that PM-SVM and PMK-ELM need. In the experiments,  $max\_gen$  was set as 200 and the range of  $C$  and  $b$  was from 0 to 1000.  $\sigma$  and  $p$  were set between 0 and 100. The size of the population was set as 50. The results generated by GA are shown in Table III.  $MAPE$  was used to evaluate the results, same as the method mentioned in [12].

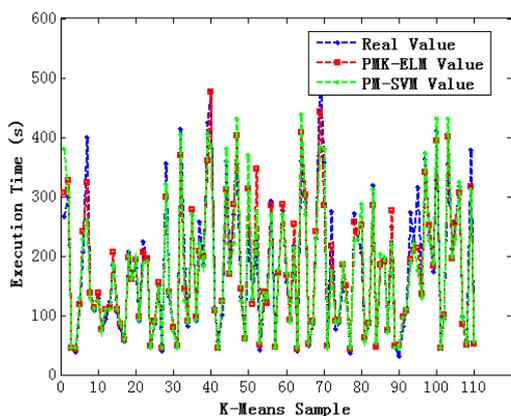
**Table III.** The best parameters generated by GA

	K-Means		WordCount	
	PMK-ELM	PM-SVM	PMK-ELM	PM-SVM
$C$	15.838	-	20.521	-
$\sigma$	0.069	-	0.867	-
$b$	-	2.285	-	6.961
$p$	-	41.967	-	16.583
$MAPE$	10.05%	10.60%	12.64%	13.42%

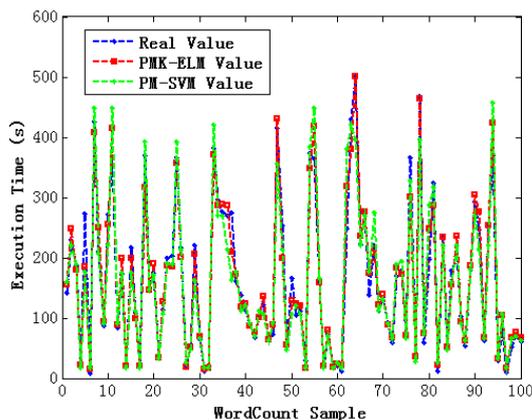
In Table IV, the results are the average value after having run for 50 times. The training time of PMK-ELM is almost 80 times shorter than PM-SVM. Moreover, for both group, the test time of PMK-ELM is about 80 times shorter than PM-SVM. Besides, the accuracy of PMK-ELM is higher than PM-SVM, too.

**Table IV.** The performance comparison between PMK-ELM and PM-SVM

		Training Time(sec)	Testing Time(sec)
K-Means	PMK-ELM	0.055	0.004
	PM-SVM	4.462	0.250
WordCount	PMK-ELM	0.043	0.03
	PM-SVM	3.324	0.307



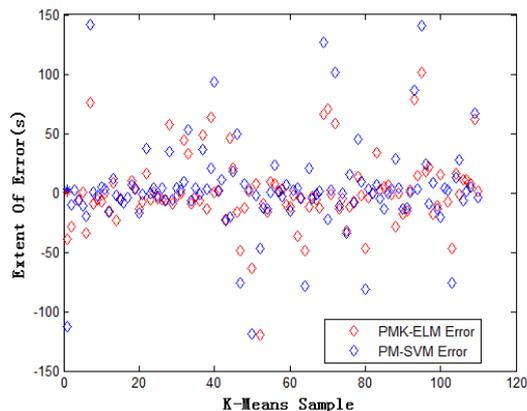
**Figure 4.** Comparison between PMK-ELM and PM-SVM in execution time of K-Means



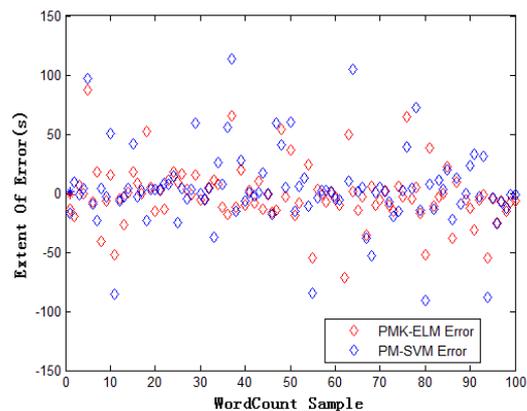
**Figure 5.** Comparison between PMK-ELM and PM-SVM in execution time of WordCount

In Figure 4, Figure 5, Figure 6 and Figure 7 the detailed results of PMK-ELM and PM-SVM are depicted. In Figure 4 and Figure 5, the line of PMK-ELM lays more closely to the real value than that of PM-SVM in two groups. On the peaks, this phenomenon is more apparent in both pictures. Although values predicted by PMK-ELM are not very accurate under some circumstance, accuracy of PMK-ELM is relatively higher compared with PM-SVM. In the Figure 6 and Figure 7, the errors of PMK-ELM are distributed near 0 intensively, while PM-SVM shows separate distribution. Trend shown in these pictures in

consistent with that shown in Figure 4 and Figure 5, which shows the performance of the PMK-ELM is better than PM-SVM. Furthermore, when the training time and test time are taken into consideration, PMK-ELM is obviously a better choice.



**Figure 6.** Distribution of error of K-Means



**Figure 7.** Distribution of error of WordCount

**5.2. The performance of proposed load balancing scheme**

In this section, the K-Means experiment is firstly run once with its execution time and hard disk space recorded. Corresponding results are shown in Table V and VI.

From Table V and Table VI, we can see that Reducer3 and Reducer6 consumed when executing the task, so the overall execution time is decided by the longest time. In Table VI, Node1 did not take part in the task, which has a better performance and may help the overall task finish earlier.

Then, we deleted the results generated by the application and applied PMK-ELM and TS-NSGA-II to this application and we got a better performance. The points shown in Figure 8 and Figure 9 are all

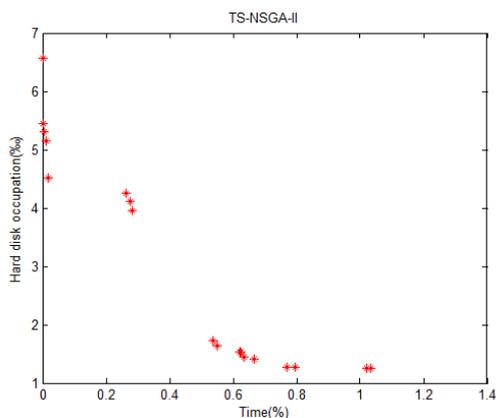
**Table V.** Hard disk space change with original Hadoop settings

NodeId	Before Execution(GB)	After Execution(GB)
Node1	405.16	403.08
Node2	406.79	404.69
Node3	404.82	402.75
Node4	412.36	410.23
Node5	405.09	402.83
Node6	413.44	411.32
Node7	404.71	404.71
Node8	404.51	402.11

the feasible solutions created by our scheme in two groups of experiments. Our scheme randomly chooses a group of solutions from each group, one is group A={1,4,6,2,8,5,3}, which represents assigning reducer0 to datanode1, reducer1 to datanode4 and so on, the other group is B={1,5,6,4,7,8,3}. The benefits we got are shown in Figure 8, Figure 9, Table VII, Table VIII and Table IX.

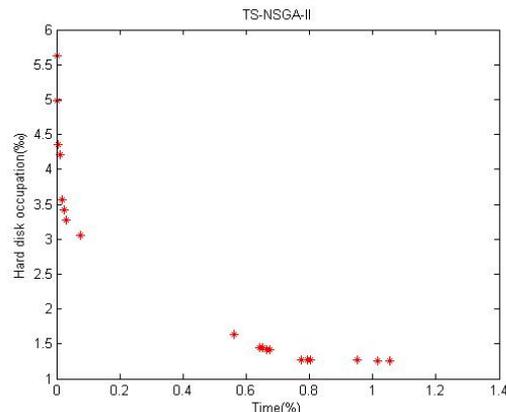
**Table VI.** Execution time of different reducers

NodeId	Reducer Group	Reducer Execution Time(sec)
Node1	Reducer0	196
Node2	Reducer5	199
Node3	Reducer1	227
Node4	Reducer4	226
Node5	Reducer3	240
Node6	Reducer6	269
Node7	-	-
Node8	Reducer2	181



**Figure 8.** Results of Group A

As shown in Figure 10, the maxim reducer execution time of Group A and B is shorter than the original Group, which determines the group A and B finish the reduce stage faster than the original. The results shown in Table IX also prove it. Not only does our load balancing scheme make the application run faster, but also helps the hard

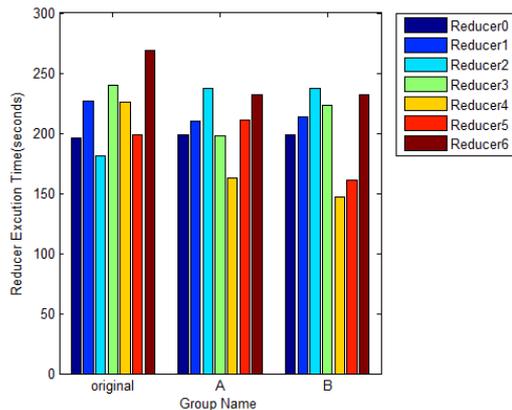


**Figure 9.** Results of Group B

**Table VII.** Hard disk space change with original Hadoop settings

NodeId	Before Execution(GB)	After Execution(GB)	
		A	B
Node1	405.16	403.08	403.08
Node2	406.79	404.53	406.79
Node3	404.82	402.70	402.70
Node4	412.36	410.29	408.03
Node5	405.09	402.99	403.03
Node6	413.44	411.05	411.05
Node7	404.71	404.71	402.58
Node8	404.51	402.38	402.41

disk occupation more reasonable. Table VIII shows the hard disk occupation when PMK-ELM and TS-NSGA-II are applied.  $S$  in Table IX is an evaluation parameter that has described in Eq.(9) in Section IV, which also shows our scheme has a better performance in job execution time.



**Figure 10.** Comparison between original and optimized schemes in reducer execution time

**Table VIII.** Comparison between original and optimize schemes in disk balancing( $S$ )

	Original	A	B
$S(\%)$	1.709	1.415	1.125

**Table IX.** The overall execution time change with PMK-ELM and TS-NSGA-II

	Original(sec)	A(sec)	B(sec)
Overall Job Execution Time	615	560	568

## 6. CONCLUSIONS

In this paper, an adaptive approach is proposed combined with a prediction model, PMK-ELM and a multi-object selective algorithm, TS-NSGA-II. The PMK-ELM can help facilitate the prediction of the execution time of tasks; whereas the TS-NSGA-II is designed to facilitate the selection of a suitable number of reducers. The experiment results have shown that both models achieve a good performance. About 47-55 seconds have been saved during experiments. In terms of storage efficiency, only 1.254% of differences on hard disk occupation were made among all scheduled reducers, which achieves 26.6% improvement than the original scheme. In the future, we would like to optimize the speculative strategy in MapReduce and try to improve the performance of the strategy.

## ACKNOWLEDGEMENTS

This work is supported by the NSFC (61300238, 61300237, 61232016, 1405254, 61373133), Marie Curie Fellowship (701697-CAR-MSCA-IFEF-ST), the 2014 Project of six personnel in Jiangsu Province under Grant No. 2014-WLW-013, the 2015 Project of six personnel in Jiangsu Province under Grant No. R2015L06, Basic Research Programs (Natural Science Foundation) of Jiangsu Province (BK20131004) and the PAPD fund.

## REFERENCES

1. Armbrust M, Fox A, Griffith R, Joseph A, Katz R, Konwinski A, Zaharia M. A view of cloud computing. *Communications of the ACM* 2010; **53**(4): 50-58.
2. Fu Z, Sun X, Liu Q, Zhou L, Shu J. Achieving Efficient Cloud Search Services: Multi-keyword Ranked Search over Encrypted Cloud Data Supporting Parallel Computing. *IEICE Transactions on Communications* 2015; **E98B**(4): 190-200.
3. Sandholm T, Lai K. MapReduce optimization using regulated dynamic prioritization. *Acm Sigmetrics Performance Evaluation Review* 2015; **37**(1): 299-310.
4. Anyanwu K, Kim H S, Ravindra P. Algebraic optimization for processing graph pattern queries in the cloud. *IEEE Internet Computing* 2013; **17**(2): 52-61.
5. Shamsi J, Khojaye M A, Qasmi M A. Data-intensive cloud computing: requirements, expectations, challenges, and solutions. *Journal of grid computing* 2013; **11**(2): 281-310.
6. Lee K H, Lee Y J, Choi H, Chung Y D, Moon B. Parallel data processing with MapReduce: a survey. *Acm Sigmod Record* 2012; **40**(4): 11-20.
7. Wu Y, Ye F, Chen K, Zheng W. Modeling of distributed file systems for practical performance analysis. *IEEE Transactions on Parallel and Distributed Systems* 2014; **25**(1): 156-166.
8. Kwon Y, Balazinska M, Howe B, Rolia J. Cost-effective resource provisioning for mapreduce in a cloud. *IEEE Transactions on Parallel and Distributed Systems* 2015; **26**(5): 1265-1279.
9. Palanisamy B, Singh A, Liu L. CSkewTune: Mitigating Skew in MapReduce Applications. *ACM SIGMOD International Conference on Management of Data* 2012; 25-36.
10. Fu J, Du Z. Load Balancing Strategy on Periodical MapReduce Job. *Computer Science* 2013; **40**(30): 38-40 (In Chinese).
11. Gufler B, Augsten N, Reiser A, Kemper A. Handling data skew in MapReduce. *International Conference on Cloud Computing and Services Science* 2011; 574-583.
12. Yuanquan F, Weiguo W, Yunlong X, Heng C. Improving MapReduce performance by balancing skewed loads. *China Communications* 2014; **11**(8): 85-108.
13. Mei Y, Liu L, Pu X, Sivathanu S, Dong X. Performance analysis of network I/O workloads in virtualized data centers. *IEEE Transactions on Services Computing* 2013; **6**(1): 48-63.
14. Islam S, Keung J, Lee K, Liu A. IEmpirical prediction models for adaptive resource provisioning in the cloud. *Future Generation Computer Systems* 2012; **28**(1): 155-162.
15. Breb S, Beier F, Rauhe H, Sattler K U, Schallehn E, Saake G. Efficient co-processor utilization in database query processing. *Information Systems* 2013; **38**(8): 1084-1096.
16. Piao J T, Yan J. Computing Resource Prediction for MapReduce Applications Using Decision Tree. *In Web Technologies and Applications* 2012; 570-577.
17. Xiao Z, Song W, Chen Q. Dynamic resource allocation using virtual machines for cloud computing environment. *IEEE Transactions on Parallel and Distributed Systems* 2013; **24**(6): 1107-1117.

18. Chen Q, Liu C, Xiao Z. Improving MapReduce performance using smart speculative execution strategy. *IEEE Transactions on Computers* 2014; **63**(4): 954-967.
19. Lam Thu Bui, Abbass H, Barlow M, Bender A. Robustness Against the Decision-Maker's Attitude to Risk in Problems With Conflicting Objectives. *IEEE Transactions on Evolutionary Computation* 2012; **16**(1): 1-19.
20. Wang J, Shang P, Yin J. DRAW: A New Data-Rouping-Aware Data Placement Scheme for Data Intensive Applications with Interest Locality. *IEEE Transactions on Magnetics* 2013; **49**(6): 2514-2520.
21. Eltabakh M Y, Tian Y, Zcan F, Gemulla R, Krettek A, Mcpherson J. CoHadoop: flexible data placement and its exploitation in Hadoop. *Proceedings of the Vldb Endowment* 2011; **4**(9): 575-585.
22. Aji A, Wang F, Vo H, Lee R, Liu Q, Zhang X, Saltz J. Hadoop GIS: a high performance spatial data warehousing system over mapreduce. *Proceedings of the VLDB Endowment* 2013; **6**(11): 1009-1020.
23. Lu X, Islam S, Wasi-ur-Rahman M, Jose J, Subramoni H, Wang H, Panda D. Improving Mapreduce Performance in Heterogeneous Environments. *42nd International Conference on Parallel Processing (ICPP)* 2013; 641-650.
24. Huang G, Huang G B, Son S, You K. Trends in extreme learning machines: a review. *Neural Networks* 2015; **61**: 32-48.
25. Samat A, Du P, Liu S, Li J, Cheng L. Ensemble Extreme Learning Machines for Hyperspectral Image Classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 2014; **77**(4): 1060-1069.
26. Huang G, Zhou H, Ding X, Zhang R. Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics* 2012; **42**(2): 513-529.
27. Savitha R, Suresh S, Kim H J). A meta-cognitive learning algorithm for an extreme learning machine classifier. *Cognitive Computation* 2014; **6**(2): 253-263.
28. Jemai J, Manel Z, Khaled M. An NSGA-II algorithm for the green vehicle routing problem. *Evolutionary computation in combinatorial optimization* 2012; 37-48.
29. Ahmad F, Chakradhar S, Raghunathan A, Vijaykumar T. K-means clustering in the cloud C a Mahout test. *IEEE Workshops of International Conference on Advanced Information Networking and Applications (WAINA)* 2011; 514-519.
30. Esteves R, Pais R, Rong C. Tarazu: optimizing MapReduce on heterogeneous clusters. *ACM SIGARCH Computer Architecture News* 2012; 61-74.
31. Isard M, Budiu M, Yu Y, Birrell A, Fetterly D: Dryad: distributed data-parallel programs from sequential building blocks. *ACM SIGOPS Operating Systems Review* 2007; 59-72.