Journal of Network and Computer Applications ■ (■■■) ■■■-■■■

Contents lists available at ScienceDirect



2 3

4 5 6

12

13 14

17

18 19 20

21 22

23

24

25

26

27

28

29

30

31

32

33

35

36

37

38

39

40

41

42

43

45

46

47

48

49

50

51

52

53

55

56

57

58

59

60

61

62

63

64

65

66

## Journal of Network and Computer Applications



journal homepage: www.elsevier.com/locate/jnca

## A new energy-aware task scheduling method for data-intensive applications in the cloud

Qing Zhao a,\*, Congcong Xiong a, Ce Yu b, Chuanlei Zhang a, Xi Zhao a

School of Computer Science and Information Technology, Tianjin University of Science and Technology, 300222 Tianjin, China

#### ARTICLE INFO

#### Keywords: Energy aware scheduling Data-intensive application SLA violation rate Data correlation

#### ABSTRACT

Maximizing energy efficiency while ensuring the user's Service-Level Agreement (SLA) is very important for the purpose of environmental protection and profit maximization for the cloud service providers. In this paper, an energy and deadline aware task scheduling method for data-intensive applications is proposed. In this method, first, the datasets and tasks are modeled as a binary tree by a data correlation clustering algorithm, in which both the data correlations generated from the initial datasets and that from the intermediate datasets have been considered. Hence, the amount of global data transmission can be reduced greatly, which are beneficial to the reduction of SLA violation rate. Second, a "Tree-to-Tree" task scheduling approach based on the calculation of Task Requirement Degree (TRD) is proposed, which can improve energy efficiency of the whole cloud system by optimizing the utilization of its computing resources and network bandwidth. Experiment results show that the power consumption of the cloud system can be reduced efficiently while maintaining a low-level SLA violation rate.

© 2015 Published by Elsevier Ltd.

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

### 1. Introduction

With the arrival of Big Data, many applications with a large amount of data have been abstracted as scientific workflows and run on a cloud platform. Cloud computing has been envisioned as the next-generation computing paradigm because of its advantages in powerful computing capacity and low application cost (Buyya, 2008; Armbrust et al., 2010; Pedram, 2012). However, the growing quantity of cloud data centers has greatly increased the total energy consumption in the world, which has become a critical environmental issue because of high carbon emissions. On the other hand, high power consumption is also a big problem in terms of economic cost from the perspective of cloud service providers. Researchers (Qureshi et al., 2009) have found, a 3% reduction in energy cost for a large company like Google can translate into over a million dollars in cost savings. High energy consumption not only translates to high cost but also leads to high carbon emissions, which are not environmentally friendly. Therefore, the problem of energy-aware performance optimization has attracted significant attention.

In literature, many existing works (Gorbenko and Popov, 2012; Fard et al., 2012) have shown that a task scheduling strategy is crucial for the overall performance of cloud workflow systems' energy

http://dx.doi.org/10.1016/j.jnca.2015.05.001 1084-8045/© 2015 Published by Elsevier Ltd. efficiency. The high ratio of under-loaded machines is the main reason for low energy efficiency. Researchers (Chen et al., 2008) have found that even during idle periods, most of today's servers consume up to 50% of their peak power. Virtualization technology (Barham et al., 2003) is a significant technology for improving the utilization of resources, and is also the technical foundation of cloud computing. Therefore, the virtual machine (VM) is the basic deployment unit in this paper. A reasonable task-scheduling strategy based on VM placement can make hosts work in a proper load so as to achieve the objective of high energy efficiency.

With the development of precision instrument technologies, many scientific research fields have accumulated vast amounts of scientific data, such as astronomy, meteorology, and bioinformatics. To solve the energy problem of these data-intensive applications, a special-task scheduling method will be presented in this paper. The scheduling principles can be derived according to the characteristics of the data-intensive application as follows:

1. Decrease network traffic through rational data layout and task scheduling.

I/O operations are the most time consuming part for dataintensive applications in the cloud. A bad task scheduling strategy can increase the amount of data transmission greatly and would directly result in an increased task response time. Then, to satisfy the user's service-level agreement (SLA), the cloud system would have to allocate more computing resources

<sup>&</sup>lt;sup>b</sup> School of Computer Science and Technology, Tianjin University, 300072 Tianjin, China

<sup>\*</sup> Corresponding author. E-mail address: zhaoqing@tust.edu.cn (Q. Zhao).

to applications to decrease their time consumptions on computing. An indirect result is that the energy consumption on the server side would be increased.

On the other hand, frequent data transmission will also lead to a large amount of power consumption. As shown in (Cavdar and Alagoz, 2012), network devices consume up to 1/3 of the total energy consumption (excluding cooling equipments), and network conflicts are the main reason for this high consumption. The frequent and large amounts of data movement will inevitably lead to the prolonging of the total network transmission time consumption and an increase in the risk of online conflict.

Therefore, we believe it is important to reduce the amount of data transmission data intensive applications in the cloud. Fortunately, there are data dependencies between tasks, and reasonable data placement and task scheduling based on these dependencies can decrease the number and time consumption of data transfers. This is the first optimization objective of this paper.

2. Improve the utilization of servers in the cloud, and reduce the generation of inefficient energy.

As mentioned above, under-loading of a machine is the main cause of low energy efficiency. Therefore, if the utilization of one machine is low, two strategies can be performed. One is to allocate more tasks to this host in order to improve its resource utilization situation. The other is to migrate the tasks on it to other machines so as to make it closed.

On the other hand, the overloaded status also needs to be changed. This is because the error rate will increase greatly under the condition of overload, hence leading to a high increase of power consumption. In the literature (Srikantaiah et al., 2009), the optimal CPU utilization of today's servers is about 70% in terms of energy efficiency. Therefore, making the active servers work at a balanced energy efficient utilization rate, and turning the underloaded server off should be a intelligent strategy. This is another optimization objective of this paper.

For the conveniences of the readers, the symbols defined in this paper are illustrated in Table 1.

The remainder of the paper is organized as follows. Section 2 presents related works. Section 3 builds the user workflow model. Section 4 builds the cloud environment model. Section 5 shows our energy consumption model. Section 6 gives the detail of the task scheduling method. Section 7 presents and analyzes the simulation results. Finally, Section 8 addresses conclusions and future work.

#### 2. Related works

The great amounts of energy consumed by supercomputers and computing centers have been a major resource and environmental concern facing today's society. In data centers, large amounts of energy are wasted by leaving computing and networking devices such as servers, switches, and routers-powered on in a low utilization state. A survey of the energy utilization state of 188 data centers mostly located in the United States points out that on average 10% of servers are never utilized (Grid, 2010). A proportion of this power could be saved if these servers were powered off or switched to low-power mode while idle. Therefore, improving the

Table 1

D	The number of global data items
$l_i (1 \le i \le  D )$	The <i>i</i> th data item
$O_{in}(t)$	The set of input data items of the task $t$
ize <sub>i</sub>	The size of the data item $d_i$
5	The number of physical servers in the cloud environment
$(1 \le j \le  S )$	The <i>j</i> th physical server
rj storage	The storage capacity of server $s_j$
i	A set of data items stored on server s <sub>j</sub>
$M_X(1 \le x \le m)$	The x th type of VMs in the cloud platform
cpu	The type of VM $VM_x$ .
$C_{mem}^{x}$	The memory c the CPU capacity of apacity of the type of VM $VM_x$ .
, crark	The $k$ th task The Worst-Case execution time of the task $t_k$ runs on a specific type of VM
/CET <sup>k</sup>	The deadline restrict of the whole application the task $t_k$ belongs to
k deadline	
j oi	The network bandwidths between server $s_i$ and servers $_j$ The CPU capacity of the physical server $s_i$
-j · cpu · :	
rj mem	The memory capacity of the physical servers <sub>j</sub>
$M_{\rm x}$	The number of the type of VM $VM_x$ that are allocated into the server $s_j$ .
til <sup>j</sup>	The CPU utilization rate of s <sub>j</sub> at time t
$ot_j$	The optimal utilization level in terms of performance-per-watt for the server $\mathbf{s}_j$
fecUtil <sup>j</sup>	The effective CPU utilization of the server $s_j$ at the time t
$M_{x\_runningTask}(t)$	The number of $VM_x$ that is running task on the server $s_j$ at the time $t$
$O_{in}(t_k)$	The number of data items in the input set of the task $t_k$
nitialData i	The data correlation between $d_i$ and $d_j$ derived from the initial data items.
nt ermediateData i	The data correlation between $d_i$ and $d_j$ derived from the intermediate data items.
j i	The integrated correlation between $d_i$ and $d_i$ .
j i	The set of tasks which need the data item $d_i$ as its input data
RD	Task Requirement Degree
OED	On-off expectation degree at time t
isk margin <sub>j</sub> (t)	How much additional workload is needed to increase the machine's utilization to $Opt_j$ in order to improve the energy efficiency of this server
veUtil <sub>p</sub>	The average CPU utilization in the period t
ower <sup>j</sup>	The overall power consumption of the server $s_j$
Power <sup>J</sup> (t)	The power consumption of the server $s_j$ at time $t$ .
otal_Power	The total power consumption the cloud system

## Download English Version:

# https://daneshyari.com/en/article/6884979

Download Persian Version:

https://daneshyari.com/article/6884979

<u>Daneshyari.com</u>