

Dynamic load balancing of SCSI WRITE and WRITE SAME commands

Vladimir Tikhomirov

Supervised by Dr. Simo Juvaste

Advised by Antti Vikman and Timo Turunen

Master's thesis.

April 27, 2013

School of Computing

University of Eastern Finland, Joensuu, Finland.

Abstract

Load balancing is computer networking method, which makes systems faster and in the same time use less resources. Almost every system can be improved by less resource usage or acceleration of speed, that is why this methodology is an important part of any computation system. The thesis describes the model of communication of SCSI controller and SCSI disks by different write commands. The application system sends SCSI WRITE and WRITE SAME commands to the disks to make the complete erasure of the disks. The aim of the research is to figure out if the parallel strategy of overwriting the disks is optimal. Moreover, in the thesis we could find the discussions about what parameters can prevent the fast speed of erasure and how is better to avoid these parameters if it is possible. The model, presenting in the thesis, is subjective, which means that it can be different for other systems, but works properly in this case.

Keywords:

Load balancing, RAID, SCSI, pass through commands.

Contents

List of abbreviations	iv
1 Introduction	1
2 Introduction to load balancing and RAID controllers	3
2.1 Load balancing	3
2.2 RAID	5
2.3 SCSI RAID controllers	8
2.4 SCSI WRITE and WRITE SAME commands	9
3 Mathematical problem definition	12
3.1 Model of the SCSI controller communication with the disks . .	12
3.2 Mathematical problem definition	15
3.3 Technical limitations	18
3.4 Theoretical estimations	23
4 Application of load balancing strategies to the HP Smart Array 642	27
4.1 HP Smart Array 642	27
4.2 Load balancing strategies	28

CONTENTS

5 Experiments and results	31
5.1 Results with WRITE SAME 10 command	31
5.2 Results with WRITE 10 command	38
5.3 Analysis of the results	42
Bibliography	44

List of abbreviations

ATA	Serial ATA (AT Attachment)
HP	Hewlett Packard
IDE	Integrated Drive Electronics
IT	Information Technology
JBOD	Just a Bunch of Disks
RAID	Redundant Array of Independent Disks
SAS	Serial attached SCSI
SCA	Single Connector Attachment
SCSI	Small Computer System Interface
SMART	Self-Monitoring, Analysis and Reporting Technology
VHDCI	Very-High-Density Cable Interconnect

Chapter 1

Introduction

Information technology (IT) is growing up so fast nowadays that we can not even compare what was 10 years ago and what is now. Current situation shows that every month, every day, every minute IT-community makes progress and it is so huge that the product, which was so popular 1 or 2 years ago is already old and people even could not use it, because the company already stopped supporting the product. All over the world people start using more and more mobile phones, laptops, read books and so on. And, of course, this is the reason why the amount of information is increasing so fast. From the users side the devices should work without errors and delays. Moreover, user wants that it will be easy to use and the information will keep privacy that nobody can get the access to it. From technical side the systems start to be too complicated, because the processes start to take more memory and the commands start to mix up, that is why for the device it started to be more difficult to handle all these things. Every computation system has the limits of speeding, memory or some other parameters, and even if the program works correct it can be decreased by some reason or can

start taking too much memory. That is why IT-community started to focus on that problem a lot.

The methodology called *Load balancing* was invented not so long time ago and serves for making the computation system faster with less usage of the resources [13]. Mostly this technology is using for the Internet services, for example, one of the most used common applications of load balancing is to provide a single Internet service from multiple servers. In this paper load balancing methodology will be considered on the example of communication between Small Computer System Interface (SCSI) controller and different amount of SCSI disks. From the computer the program sends the SCSI write and SCSI write same commands to the disks to make full erasure. The problem is to figure out what technique is faster and take less memory by load balancing methods. This problem is similar with the Internet one, but has some special issues, which will be opened in the thesis.

Chapter 2

Introduction to load balancing and RAID controllers

2.1 Load balancing

In general load balancing is the methodology which increases the speed by balancing the resources. The main idea is to do as much as possible with least amount of resources. Mostly it depends on the processes, which we can divide by threads and launch them parallel. Sometimes load balancing is divided to the 4 types: static, dynamic, combined and preemptive. In scientific articles combined and preemptive load balancing systems are mentioned very rarely because of their specificity and usefulness only in some certain cases. That is why in this research it will be considered 2 main different types: static and dynamic.

Static load balancing system makes the analysis before the application starts. That is why for creating good load balancing system via static methods the

CHAPTER 2. INTRODUCTION TO LOAD BALANCING AND RAID CONTROLLERS

programmer should know several pieces of information. Firstly, the amount of distributed work should be well known before the program is written. And secondly, the duration of the tasks to be divided should be approximately the same. There are several methods to hand out the tasks such as Round Robin algorithm, Random algorithm, Recursive Bisection algorithm and some others [10]. Frequently the initial parameters are taken from the previous launch results and genetic algorithms are used. An example of load balancing system can be taken from image processing. Let consider 16×16 segmented image and 8 processors. With that image each processor would work on 4 segments. The most important piece of all this division of work is that each processor determines this information and also which segments it will work on.

Dynamic load balancing systems have more possibilities because the balancing process can be done during the application running [5]. This fact gives the advantage of resource movement from busy part to less busy, that all of the parts of the application will be in average busy condition. The disadvantage of that method is that the system should also collect the information about the condition of application, which follows the increasing of the resources. There are 3 general criteria, which programmer should follow during making the dynamic load balancing system:

1. Which process is busy or free
2. Priority of the task to be done
3. Length and/or duration of the task

The program, which realized the dynamic load balancing system, should check the loading of computation units, connection possibility and frequency of sending commands. Dynamic load balancing methodology has a lot of

methods, such as Bidding algorithm, Drafting algorithm, Recursive Coordinate Bisection (RCB) and some others. In scientific papers they can be called as strategies and policies, but in the following content it will be called methods.

2.2 RAID

RAID - Redundant Array of Independent Disks - is a storage technology that combines multiple disk drive components into a logical unit. That means that this technology gives a possibility to keep the information in a specific way. RAID controllers are the devices, which realize this technology. In simple words without a RAID controller computer can see physical disks as a logical disks without any difference. But with RAID controller the computer will see only the logical disks, which will be definitely different in comparison in physical disks. All issues about creation of these logical disks are taken by RAID controller.

There are several types of RAID: RAID 0, RAID 1, RAID 5, RAID 6 and some others [4]. These types of RAID, which were mentioned, are basic types of RAID, thus there are many other types, which also combine them together. For example, RAID 10 is actually combining of RAID 0 and RAID 1. Each type or architecture provides a different balance between the key goals: reliability and availability, performance, and capacity. That means that every RAID keeps the information in a different way.

One of the huge advantages of using RAID is that if one of the disks breaks down others can recover the information on that disk. Not all RAID architectures supports this functionality, but most of them do. Some types of

CHAPTER 2. INTRODUCTION TO LOAD BALANCING AND RAID CONTROLLERS

RAID provide fault tolerance of even two drive failures (RAID 6).

Let consider some types of RAID more detailed. RAID 1 called also as "mirroring" is one of the main fundamental types of RAID, which refers to maintaining duplicate sets of all data on separate disk drives. There must be two disks in the configuration and there is a cost disadvantage as the usable capacity is half the number of available disks. RAID 1 provides cost-effective, high fault tolerance for configurations with two disk drives. The prototype of RAID 1 is presented on the figure 2.1.

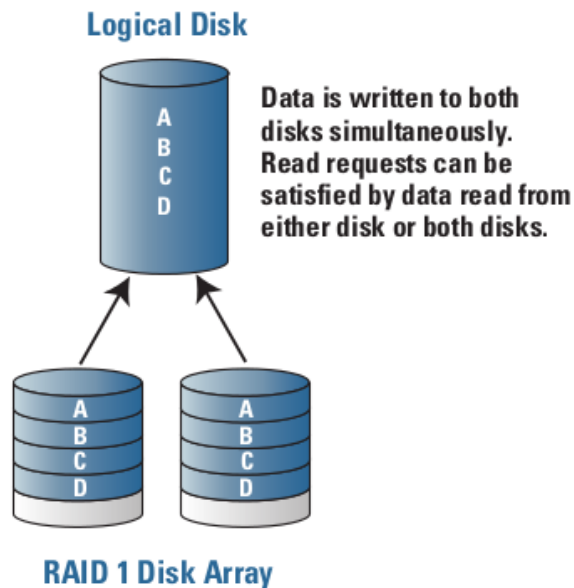


Figure 2.1: Prototype of RAID 1

Another type RAID 5 uses data striping in a technique designed to provide fault-tolerant data storage, but doesn't require duplication of data like RAID 1 [2]. Data is striped across all of the drives in the array, but for each stripe through the array (one stripe unit from each disk) one stripe unit is reserved

CHAPTER 2. INTRODUCTION TO LOAD BALANCING AND RAID CONTROLLERS

to hold parity data calculated from the other stripe units in the same stripe. Read performance is therefore very good, but there is a penalty for writes, since the parity data has to be recalculated and written along with the new data. RAID 5 requires a minimum of three disks and a maximum of 16 disks to be implemented. RAID 5 usable capacity is between 67% - 97% depending on the number of data drives in the RAID set.

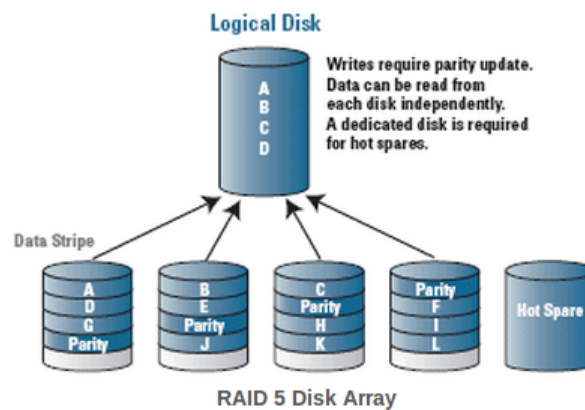


Figure 2.2: Prototype of RAID 5

Moreover, nowadays there are RAID controllers, which already have load balancing system inside and the user can set the mode, which is more comfortable in the some situations. In our case the commands should be sent straight to the disk that is why the load balancing system from controller is not helping at all, because of its own methodology.

On the picture 2.3 one of the SATA RAID controllers is presented. From the name of the device it is clear that this controller can be connected to 8 SATA disks. It supports following types of RAID: 0, 1, 5, 10, 50 and JBOD (Just a Bunch of Disks).



Figure 2.3: 3ware SATA RAID controller 9500S-8

2.3 SCSI RAID controllers

SCSI - Small Computer System Interface - is a set of standards for physically connecting and transferring data between computers and peripheral devices [11]. The SCSI standards define commands, protocols, and electrical and optical interfaces. There are also other interfaces such as SATA (Serial ATA (AT Attachment)), IDE (Integrated Drive Electronics), SAS (Serial attached SCSI). In general, the interfaces are divided by two categories: ATA and SCSI. In the following content it means the according protocol. IDE and SATA interfaces are related to ATA and SAS is related to SCSI. Several types of disks are presented on the figure 2.4. In this paper the topic mostly is about SCSI interface, but sometimes it will be compared with ATA. If SAS and SCSI are compared, it can be considered similar interfaces. However, SAS is a new version of SCSI and it gives a single channel for each disk. In

CHAPTER 2. INTRODUCTION TO LOAD BALANCING AND RAID CONTROLLERS

comparison with SCSI it was only one channel for all disks.



(a) SCSI disk



(b) SATA and IDE disks

Figure 2.4: Three different types of disks

2.4 SCSI WRITE and WRITE SAME commands

This research is made for the Blancco Oy Ltd, which produce the erasure software. The main part of the erasure is performed by SCSI WRITE and WRITE SAME commands [7]. In general, all commands are divided to three groups: Non-Data commands, Data-In and Data-Out. For the last two groups it means that the command is sent including the buffer. Data-Out means that the data is sent to the disk. For the Data-In group it is in opposite way and data is come from the disk. The example of Data-In command is any SCSI READ command. Write commands are related to the Data-Out group because these commands must have Data-Out buffer, containing the data, which will be written to the disk. Nowadays, for SCSI there are 4 different versions of write command and 3 for the write same command.

CHAPTER 2. INTRODUCTION TO LOAD BALANCING AND RAID CONTROLLERS

Lets consider SCSI write commands first. Mostly the commands are different because of the length. Each command has not only operation code, which is the main criteria for the command, but also some other variables such as Logical Block Address (LBA), transfer length, control and some others. These variables do not include Data-Out buffer and the Figure 2.5 shows the example how the commands are defined in specification [7].

Table 108 — WRITE (10) command									
Byte	Bit	7	6	5	4	3	2	1	0
0		OPERATION CODE (2Ah)							
1		WRPROTECT			DPO	FUA	Reserved	FUA_NV	Obsolete
2		LOGICAL BLOCK ADDRESS							
...									
5									
6		Reserved			GROUP NUMBER				
7		TRANSFER LENGTH							
8									
9		CONTROL							

Figure 2.5: Definition of WRITE 10 command

The following commands exist for writing the data to the disk:

- WRITE (10)
- WRITE (12)
- WRITE (16)
- WRITE (32)

The numbers in brackets shows the length of the command in bytes. It is obvious that the WRITE (10) command is the base for others. It seems that by now the best choice is WRITE (16) because in [7] there are some notes, that «Migration from the WRITE (10) command to the WRITE (16)

CHAPTER 2. INTRODUCTION TO LOAD BALANCING AND RAID CONTROLLERS

command is recommended for all implementations». WRITE (32) should be used only in some special cases [7].

SCSI write same commands have the same purposes as SCSI write commands. The difference and biggest advantage is that SCSI write same command can be sent once and the Data-Out buffer can be written to the disk several times. That gives faster speed because the bus is not busy any more - only one command was sending, but the buffer is still writing. There are 3 different types of write same commands, as it was noticed before:

- WRITE SAME(10)
- WRITE SAME(16)
- WRITE SAME(32)

The number in brackets also shows the length of the command in bytes. Number of logical blocks is one of the most important variable in these commands because this number of blocks will be written to the disk with the same buffer. It is obvious that as bigger is the length of the SCSI write same command as more number of blocks we can write to the disk. ATA has the similar command, but it should be send differently and the device should support SMART (Self-Monitoring, Analysis and Reporting Technology) [8]. It is a monitoring system for computer hard disk drives to detect and report on various indicators of reliability.

Chapter 3

Mathematical problem definition

3.1 Model of the SCSI controller communication with the disks

There are a lot of resources, which can be used for different computations. Currently the time costs very much that is why it is important to make programs as fast as possible. Because of that optimization problem is one of the most significant ones. In general it is the problem of finding the best solution from the set of all feasible solutions. It is often done with different algorithms, which analyse the solutions during variable changing. In our case the optimization problem is based on time optimization. That means that the time will be optimized by several variables which can be changed.

Lets consider the model presenting on the picture 3.1, which shows how SCSI controller is connected to the disks and CPU (Central processing unit). This picture can be divided to 5 parts: CPU, where Blancco software is running, bus *a* - from CPU through motherboard to the SCSI controller,

SCSI controller b , buses c - from SCSI controller to the disks d . Let consider that the SCSI controller is connected to N disks.

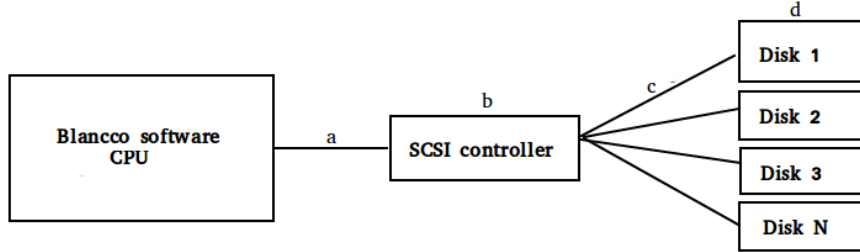


Figure 3.1: Communication model of SCSI controller

This figure presents how SCSI controller is connected in our case.

All the time people want to do everything as fast as possible, but of course there are limitations in different places. The same situation is in this model - each part has the limitation of the speed. Moreover, there are also limitations for the memory.

Lets consider the figure 3.1 more detailed. In this model CPU part will be skipped, because there are so powerful computers in current time that it should not be taken into consideration. The main idea of Blanco software is to write the data to the disks. Lets discuss what *data* means in this context. On the figure there are N disks and each of them has its own capacity d_i , where $i = \overline{1, N}$ shows the number of the disk. Capacity d_i is calculated in megabytes. For each disk i SCSI controller b should send the amount of data, which is at least equal to the capacity of the disk. Moreover, it is needed to send the commands for setting the connection with the disk.

As discussed in 2.4 Blanco software uses two different versions of write

command. Each of these commands has its own capacity. Let define these capacities in a mathematical way. The capacity of the write command is c_1 , the capacity of the write same command is c_2 and the capacity of other commands is c_3 . In these context other commands means the commands for identification the device and setting the connection. All these three variables c_1, c_2, c_3 are constant and are defined by Blancco software. Hence for replac- ing all data on N disks it should be sent F MB from SCSI controller b to the disks d through the buses c . Let F be called as *complete erasure*, is it is calculated by the following formula:

$$F = \sum_{i=1}^N F_i, \quad (3.1)$$

where F_i is the amount of megabytes, that should be sent to the disk i , which is calculated as

$$F_i = c_1 m_{i1} + c_2 m_{i2} + c_3 m_{i3}, \quad (3.2)$$

where m_{ik} is the number of times that command c_k should be sent, where $k = \overline{1, 3}$. Furthermore, complete erasure F should fulfil the following condition:

$$c_1 c_2 = 0, \forall F_i, i = \overline{1, N}. \quad (3.3)$$

This condition means that only 1 command c_1 or c_2 can write to each disk. Thus the software can not erase the same disk with both commands c_1 and c_2 . That means that only software decide which command is better to use and dynamic load balancing system should apply the strategy how to apply it more efficiently.

For clear vision and comprehension it is convenient to present the set of F_i

as a matrix F_m :

$$F_m = \begin{pmatrix} F_1 \\ \vdots \\ F_N \end{pmatrix} = \begin{pmatrix} c_1 m_{11} & c_2 m_{12} & c_3 m_{13} \\ \dots & \dots & \dots \\ c_1 m_{N1} & c_2 m_{N2} & c_3 m_{N3} \end{pmatrix}, \quad (3.4)$$

where F_m is the matrix with dimension $N \times 3$.

The capacity of each disk i should be less than amount of data, which is sent for writing. That means that the following condition should be fulfilled:

$$d_i < F_i, i = \overline{1, N}. \quad (3.5)$$

The main aim of the research is to find the way how to send the data faster to the disks. Of course, that it depends on the devices and on the strategy how the software sends the commands. Let define the variable S_{max} for maximum speed and it is obvious that this variable is constant. According to the formula 3.1 and maximum speed S_{max} it is possible to calculate the minimum time T_{min} for realization complete erasure F by the following formula:

$$T_{min} = \sum_{i=1}^N t_i = \frac{\sum_{i=1}^N F_i}{S_{max}} = F/S_{max}, \quad (3.6)$$

where t_i is the time for sending F_i MB to i disk with the maximum speed S_{max} .

3.2 Mathematical problem definition

Let consider that in the same model there is dynamic load balancing system and it calculates the application state $m = M/s$ times, where M is the amount of megabytes to complete the application task and s is the step. In

this context step s means the amount of commands and m means how many times the dynamic load balancing system should change the strategy of the input data for the next step.

There is matrix $H = \{h_{ij}\}$, where $i = \overline{1, N}$ is the disk number, $j = \overline{1, m}$ is the step number of the dynamic load balancing system and h_{ij} is the number of commands, which is sent to the i disk at j step. That means that elements of matrix H should fulfil the following condition:

$$M = \sum_{i=1}^N \sum_{j=1}^m h_{ij} = \sum_{i=1}^N F_i = F. \quad (3.7)$$

That means that for the complete erasure F it is needed M commands. Moreover, each of h_{ij} can also be described by F_i . As an example of matrix H it is convenient to consider 4 disks, which can be erased by 100 commands with the condition that we can not send more than 100 commands per step. Thus the matrix H can be presented as

$$H_1 = \begin{pmatrix} 100 & 0 & 0 & 0 \\ 0 & 100 & 0 & 0 \\ 0 & 0 & 100 & 0 \\ 0 & 0 & 0 & 100 \end{pmatrix} \quad (3.8)$$

or a bit more complicate:

$$H_2 = \begin{pmatrix} 50 & 10 & 20 & 20 \\ 20 & 50 & 10 & 20 \\ 20 & 20 & 50 & 10 \\ 10 & 20 & 20 & 50 \end{pmatrix}. \quad (3.9)$$

Matrix H_1 and H_2 are practical examples of the matrix H . The first matrix shows the trivial strategy - to send as many commands as the program can to fulfil the condition during each step. Matrix H_2 shows that it could be

done in a different way and we can divide the whole amount of commands and send it separately. The strategy should be applied by some dynamic load balancing algorithm and the aim of the research is to find out what is the optimal matrix H_{opt} and could load balancing help in this case or not.

Also there is matrix $Q = q_{ij}$, where q_{ij} is the order number in the queue for sending commands to i disk at j step. In this case it is possible to show to which disk the program should send the command first. The dimension of Q is the same as dimension of H :

$$\dim(Q) = \dim(H) = N \times m, \quad (3.10)$$

where N is the number of disks and m is the number of steps in the load balancing system.

For mathematical problem definition it is convenient to introduce the functional of time. Let the *functional of time* be defined as the functional

$$T = \sum_{j=1}^m t_j, \quad (3.11)$$

where t_j is the time for sending $\sum_{i=1}^N h_{ij}$ commands at j step of dynamical load balancing process.

The main idea of the optimization problem is to minimize the functional of time T

$$\min_{H, Q, s} T. \quad (3.12)$$

According the formula 3.6 current minimization problem can be rewritten as

$$T \xrightarrow{H, Q, s} T_{min}. \quad (3.13)$$

That means that the application of dynamic load balancing algorithm should make the program faster and the value of T should come close to the minimum

time T_{min} . The algorithm should use H , Q , s as parameters for making the program faster.

3.3 Technical limitations

Every device has its own technical limitations. For example, a car has a limitation for the speed, capacity of petrol tank, maximum engine power and so on. Discussing the SCSI controllers every device has special characteristics, which we can compare with characteristics of the car.

Let introduce some parameters, which can influence on the communication between controller and disk. In this paper the word SCSI is used a lot, but there was no clear explanation what is it. As it is written in the beginning of the paper SCSI is Small Computer System Interface. The main word in this abbreviation is interface, which shows *how* the devices can be connected with each other. Lets again consider the figure 3.1 and the situation that we write some data from CPU to the disk. Controller b has 2 interfaces, because it is connected to the CPU and disks d . Let call bus a as *input* interface and buses c as *output* interface. For the disks d there is only input interface, that is why it is possible to skip the word *input* and use only interface. So, if there is a sentence, which includes "the controller has SCSI interface", that means that the controller has the possibility to connect SCSI disks and the output interface is SCSI. Moreover, because controller has 2 connections there is also input interface, which can be, for example, different modifications of PCI (Peripheral component interconnect) such as PCI-X, PCI64, PCI Express or some others. In general words if there is a discussion about SCSI interface it means only the buses c .

Moreover, SCSI is also a data protocol, which shows what commands should be sent through buses c . Depends on previous content interface SCSI connection has different technical limitations. There are several SCSI interfaces such as SCSI-1, Fast-Wide SCSI, Ultra2 SCSI, Ultra-320 SCSI and some others, but in this research Ultra-320 SCSI interface will be considered more, because the testing device is supported that interface. Nowadays there is interface with faster bandwidth - 640 MB/s - but it did not get popular, because it supports maximum 2 devices per cable.

Ultra320 SCSI [6] is the seventh generation of SCSI I/O (Input/Output) technology. The dominant feature is that the speed is increased to 320 megabyte per second (MB/s). In Ultra320 SCSI devices all support packetized protocol and may support Quick Arbitration and Selection (QAS). Expander communications techniques have also been defined. In general Ultra320 shows that only buses c from the figure 3.1 have this speed and not disks d , and of course, not controller b . In the ideal situation to get the maximum speed of the device communication should fulfil the following condition:

$$S_a > S_b > S_c > S_d. \quad (3.14)$$

Otherwise there is a question, how controller is going to manage several disks if it has the same speed. Currently there are different PCI buses, which can have the maximum bandwidth up to 4 GB/s. Moreover, there are also some possibilities to increase the performance of PCI [1]. But even it is not the limit, because there are also some other buses, such as QPI (QuickPath Interconnect) and HyperTransport with different modification, which can easily have the bandwidth around 25 GB/s. By now the maximum bandwidth has HyperTransport 3.1 with value 51,20 GB/s. Of course, that bus depends on the frequency quite a lot. HyperTransport 3.1 has the maximum frequency

3,2 GHz. Depends on the data content and the number of the disks the data for the transport through bus a can increase, but the numbers show that it seems that it will not be a problem, but still it is better to keep this in mind.

All SCSI devices are always backward compatible, which means that when newer SCSI devices are attached to a system with devices from a previous generation, the newer devices will fall back to the maximum operating speed of the older generation when they are talking. When not talking to older devices, the newer devices operate at their normal speeds. That means that if one of SCSI drives from the figure 3.1 works on Ultra160 interface, the controller, which supports Ultra320, will talk to the disk using Ultra160 interface.

Most of the time in the hardware if the device or cable can be connected to some place that means that it will work. The same situation is with SCSI interfaces, but for some generations connectors are the same. For the interfaces Ultra2 Wide, Ultra3, Ultra320 and Ultra640 the connectors can be 68-pin and 80 pin which belongs to the connection type SCA/SCA-2 (Single Connector Attachment). The examples of the connectors are presented on the figure 3.2.

Mostly the discussion about technical limitations contained the information about the buses a and c , but what about controller and disks? Let start from the disks. In this chapter it was already introduced that each disk i has capacity d_i , speed S_d and special interface. But there are some other parameters, which could be also important. It is *cache*, *spindle speed*, average *seek time* for reading and writing, connector type and dimension. The last parameter is not very interesting, because it should be imagined that there is enough space.



Figure 3.2: 68-pin and 80-pin SCA connectors

In general cache is used in many situations, for example, in operation systems, in controllers, in disks and some other places. Cache is the part of memory, which has very fast speed access [9]. For the drives it means that cache is stored before the physical hard disk platter, which gives the possibility to get the data much faster. Currently disks can have cache with capacity from 8 to 64 MB. Of course, that big cache gives higher performance for the disk. Moreover there are different cache algorithms such as LRU (Least Recently Used), MRU (Most Recently Used), LFU (Least-Frequently Used), Direct-mapped cache and some others [14], [12]. The aim of every cache algorithm is to manage the data in the cache and show, which data block to delete and which one to keep. Depends on the disk it can be applied different cache algorithms, which could give the advantage or disadvantage in the current problem. Also, depends on the number of the disks maybe load balancing system should find some tricky steps through disk cache, which can make the process faster.

Another important characteristic of disk is spindle speed, which shows the frequency of rotation. This value is measured in revolutions per minute (RPM). The value of spindle speed shows how many full rotations were completed in one minute around a fixed axis. Usually for ATA disks spindle speed can be 5400 or 7200 rpm (90 or 120 Hz) and in this case SCSI disks win, because their spindle speed can be 10,000 or 15,000 rpm (160 or 250 Hz). This value is constant and is set by manufacture. So, there is no reason to try make it higher, the only possibility is to take the disk with higher spindle speed.

Both of described parameters cache and spindle speed are influenced on the seek time of the disk. Seek time means the time, which needs for the head of the disc to move to the right position. It is clear that this time can be only average, because in different situations there are different values of cache, spindle speed and data task. In general there are two values of seek time - one for reading operations and another for writing, but usually it is only one value for both operations. Most of the time in the research it will be considered only as a writing seek time, but the value of seek time for the reading should not be thrown out. There are two seek measurements called track-to-track and full stroke. The track-to-track measurement is the time required to move from one track to an adjacent track. This is the shortest (fastest) possible seek time. In hard disk drives (HDD) this is typically between 0.2 and 0.8 ms. The full stroke measurement is the time required to move from the outermost track to the innermost track. This is the longest (slowest) possible seek time.

Lets consider the controller now and figure out what parameters this device has. First of all, it is the possibility to connect buses a and c , which was decided to call as input interface and output interface, which also follow the specific parameter, such as speed. In the literature it is also possible

to see that input interface is called as host interface. Secondly, it is the amount of RAIDs, which is supported by this controller. There are also some more parameters, which are still important, such as cache, cache function or algorithm, maximum amount of physical disks, maximum amount of logical disks. It was discussed about cache, but for the controllers cache has more power. Currently cache of the controllers can have values up to 512 MB. Moreover, some controllers have even their own cache algorithms, which make it stronger.

There is one parameter, which should be marked out. It is the bandwidth of the controller, which definitely influences on the performance of the device. But most of the time in the parameters of the controller it is possible to find only the bandwidth of the input and output interfaces. It is like this, because currently the controllers are so powerful, that they can manage all the data, which comes to the device. That means that the speed of the controller is not the bottleneck and is not very important in our case.

3.4 Theoretical estimations

Lets model the real situation and try to estimate how long time does it take to write some data to the disk. Let consider that there is a Compaq SCSI drive BD01864552 with the following parameters:

- Capacity: 18.2 GB
- Spindle speed: 10,000 RPM
- Interface: Ultra3 SCSI
- Seek time: 3 ms

From these parameters mostly we will focus on the interface, which has a speed 160 MB/s. That means that the speed of the bus c is equal to 160 MB/s. Basically all data on disks are divided to the sectors. Most of the time the size of the sector is 512 bytes, but in some cases this value can be different. In this research we consider that the sector size has 512 bytes value. We will focus on the parameters of the commands and the disks because the controller should not be a bottleneck, but we still need to think that there are some limitations also in the controller.

WRITE SAME command Let consider the WRITE SAME 10 command, which comes from the CPU [7]. This command has 32 bits for the LBA (Logical block addressing) and 16 bits for the Number of Logical Blocks. The main huge advantage of that command is that we can send the same buffer to several logical blocks only with one command. It gives us a possibility to write much more data than we really send by the command.

We can calculate the maximum amount of data that we can write with one WRITE SAME 10 command to the disk:

$$\frac{2^{16} * 512}{1024^2} = 32MB. \quad (3.15)$$

Let consider that the model 3.1 consists of 1 disk and we do not need any other commands except WRITE SAME one. That means that in condition 3.3 parameter c_2 will be equal to 0. It gives a possibility to calculate complete erasure F as a multiplication of command c_1 and amount of commands m_{11} . Because we know the capacity of the disk we could compute the number of commands m_{11} that we need to send for erasure 18.2 GB disk.

$$m_{11} = \frac{F_1}{c_1} = \frac{18.2 * 1024}{32} = \frac{18636.8MB}{32MB} = 582.4. \quad (3.16)$$

For performance of real erasure the data should be divided for filling the last part of the disk, because we could not send 0.4 command. The main aim of the research is testing the erasure time that is why we will not take this little part into consideration and will use the approximate value of 582.

The size of WRITE SAME 10 command is 10 bytes. Including the buffer, which is 512, we get that one whole command is 522 B. So, to erase 18.2 GB disk we need to send the following amount of data:

$$(10 + 512) * 582 = 303804B. \quad (3.17)$$

That means that for complete erasure of one disk we need to send only 304 KB. Another command WRITE 10 does not have so huge advantage, but lets compare these two commands, because in some cases WRITE SAME 10 command is not working.

WRITE command Let consider another pass through command WRITE 10 [7]. The structure of this command is completely the same in comparison with WRITE SAME 10 command, but some parameters have different values. For example, WRITE SAME 10 command has operation code 0x41 and WRITE 10 command has 0x2A. The main difference of these 2 commands is that with WRITE 10 command we need to send whole buffer all the time. The buffer length is calculated by the following Formula:

$$Buffer\ length = Sector\ size * Transfer\ length. \quad (3.18)$$

In our situation the disks have sector size 512 and this value depends only on the disk. But buffer length can be varied by transfer length, which can have a value up to 65536. That means that if we decide to set the maximum transfer length we could send 32 MB with one WRITE 10 command. In

comparison with WRITE SAME 10 command, which has 512 B buffer, our buffer can increase to 512×65536 , which is quite big value if we consider that the controller connect to several disks. Usually, in Linux driver of the controller there is a value of maximum buffer, which is 4096×512 in our case. It can be optimal for some amount of disks, but it should be checked, because this is one of the parameter that we can influence. If a lot of data go with one command through the controller it can make the erasure process slower because the cache of the controller will be full of messages, waiting in the queue. We should influence on the erasure speed by this parameter and find the optimal value during testing. By now this value is set to 256, which means that we send 128 kB at once. We can calculate how many commands we need to erase one 18.2 GB disk by WRITE 10 command with 128 KB buffer:

$$m_{11} = \frac{F_1}{c_2} = \frac{18.2GB}{128kB} = \frac{18636.8MB}{0.125MB} = 149094.4 \quad (3.19)$$

The derived value is 256 times bigger than the amount of commands that we need to send for erasure by WRITE SAME. It is obvious, because the current value of buffer length is 256 bigger. In the chapter 4 we will see the results of erasure testing with different transfer lengths.

Chapter 4

Application of load balancing strategies to the HP Smart Array 642

4.1 HP Smart Array 642

In this chapter we consider experiments and results, which we got during testing different load balancing strategies with HP Smart Array 642. This device is a SCSI RAID controller, which supports protocol Ultra-320 SCSI [3]. The maximum bandwidth of this protocol is 320 MB/s. In current time there is also Ultra-640 SCSI protocol, which is faster and has a bandwidth 640 MB/s. However, Ultra-640 did not get the popularity because of some disadvantages that is why we chose a controller which supports Ultra-320 protocol.

HP Smart Array 642 has one internal and one external VHDCI (Very-High-

Density Cable Interconnect) SCSI ports. The internal port gives the possibility to connect up to 6 internal disks and the external one up to 14 disks. Without any high mathematics we get that HP Smart Array 642 controller supports 20 disks. Maximum capacity for all disks in total is 6 TB. Controller has PCI-X bus, which has a 133-MHz frequency. That means that the maximum bandwidth of the bus is 1 GB/s. Moreover, the architecture is 64-bit.

The described controller is presented on the figure 4.1.



Figure 4.1: HP Smart Array 642 controller

4.2 Load balancing strategies

This section presents a description of the methods for solving the problem presented above in the chapter 3. The solutions that we offer are simple, but effectively show the results in practice. We applied them for testing the erasure time of several disks. We remember that the main problem of

CHAPTER 4. APPLICATION OF LOAD BALANCING STRATEGIES TO THE HP SMART ARRAY 642

the research is to minimize the functional of time T or prove that it is not possible because of some limiting factor.

During testing the erasure time of disks with HP Smart Array 642 we applied three load balancing strategies. First two are sequential and the third is parallel, which is working more effectively than the others. Because first two strategies are much slower than the third one, they were tested with the backplane of only 5 disks connected to the internal port of the controller. The third strategy was tested with the enclosure of 14 disks connected to the external port of the controller. In both cases we chose 18.2 GB disks from Compaq. Theoretical expectations about these disks were discussed in the section 3.4.

The first strategy performs the complete erasure disk by disk. That means that until the current disk is not erased the program does not start erasure of the next disk. From the section 3.4 we know that for erasure one 18.2 GB disk we need to send 582 WRITE SAME 10 commands. According to the point that for testing first strategy we took 5 disks we get the following matrix H_{s1} from the model 3.7 for complete erasure:

$$H_{s1} = \begin{pmatrix} 582 & 0 & 0 & 0 & 0 \\ 0 & 582 & 0 & 0 & 0 \\ 0 & 0 & 582 & 0 & 0 \\ 0 & 0 & 0 & 582 & 0 \\ 0 & 0 & 0 & 0 & 582 \end{pmatrix} \quad (4.1)$$

From the matrix H_{s1} we can conclude that first strategy has only 5 steps in the algorithm.

Second strategy sends one command per disk and then changes the disk.

That means that in this case matrix H is completely different:

$$H_{s2} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & \dots & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & & 0 & 0 & 0 & 0 & 1 \end{pmatrix}. \quad (4.2)$$

For the second strategy the dimension of matrix H is 5×2910 . Value 2910 came from multiplication of number of disks and number of commands that we need to send for the erasure. From the logic of sending commands the devices should be changed quite often, which can be a disadvantage of this strategy.

The third strategy is using parallel computing, which makes the process for complete erasure much faster. The application creates separate thread for each disk and start sending WRITE SAME commands to that disk. It is obvious that during increasing the amount of disks the value of functional of time will increase too, but slowly, because threads are working completely separately. But important limitations are the speeds of the disk and the bus. As we can see from Figure 3.1 separate buses c give the possibility to send data from the controller b to the disks d . In our situation, only one bus c exists, which is connected to the enclosure containing several disks. If we write information in parallel to several disks the bus speed can be the limitation for data transfer. The maximum speed of the current bus is 320 MB/s. In this case it is difficult to present matrix H , because we do not know in which order WRITE SAME commands come to the disks.

Chapter 5

Experiments and results

In this chapter we present the test results of the application, which sends WRITE 10 and WRITE SAME 10 commands with different parameters. We made all the tests with WRITE 10 command, when the disks were connected with the enclosure to the controller. But some tests with WRITE SAME 10 command were done with the backplane. Mostly we are interested in decreasing the time of complete erasure, but sometimes, because of the different speed of the disks, the results started to be useless and strange, that is why we started to calculate also the time of erasure for single disks.

5.1 Results with WRITE SAME 10 command

Using WRITE SAME 10 command we tested 3 different strategies that we considered before. The tests for the first and second strategies were done with 5 disks, connected with the backplane to the internal port of the HP Smart Array 642 controller. The third strategy was testing with 14 disks, connected with the enclosure to the external port of the controller.

CHAPTER 5. EXPERIMENTS AND RESULTS

During testing erasure with first strategy we got the results presented in the Table 5.1.

Table 5.1: Results for complete erasure by first strategy

Disks	Capacity	Time	Average time/disk	Average speed
1	18.2 GB	8 min 32 sec	8 min 32 sec	35.5 MB/s
2	36.4 GB	18 min 12 sec	9 min 6 sec	33.3 MB/s
3	54.6 GB	28 min 5 sec	9 min 21 sec	32.4 MB/s
4	72.8 GB	39 min 7 sec	9 min 46 sec	31.0 MB/s
5	91 GB	48 min 45 sec	9 min 45 sec	31.1 MB/s

Figure 5.1 presents the results from the Table 5.1 more clearly.

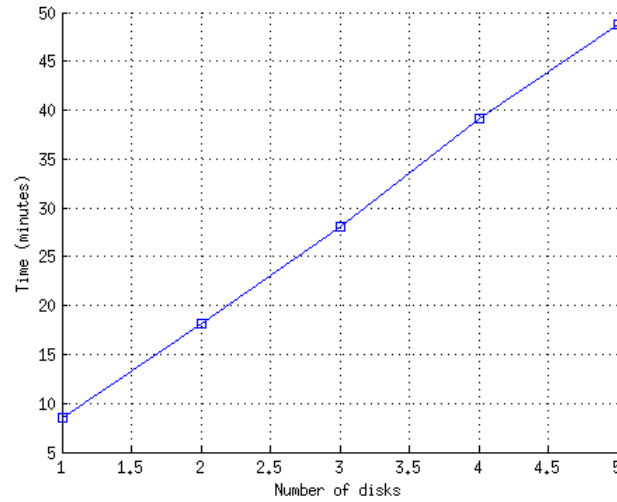


Figure 5.1: Time results for complete erasure by first strategy

CHAPTER 5. EXPERIMENTS AND RESULTS

During testing erasure with second strategy we got the results presented in the Table 5.2.

Table 5.2: Results for complete erasure by second strategy

Disks	Capacity	Time	Average time/disk	Average speed
1	18.2 GB	10 min 54 sec	10 min 54 sec	27.8 MB/s
2	36.4 GB	20 min 42 sec	10 min 21 sec	29.3 MB/s
3	54.6 GB	30 min 41 sec	10 min 13 sec	29.6 MB/s
4	72.8 GB	40 min 25 sec	10 min 6 sec	30.0 MB/s
5	91 GB	50 min 10 sec	10 min 2 sec	30.2 MB/s

Figure 5.2 presents the results from the Table 5.2 more clearly.

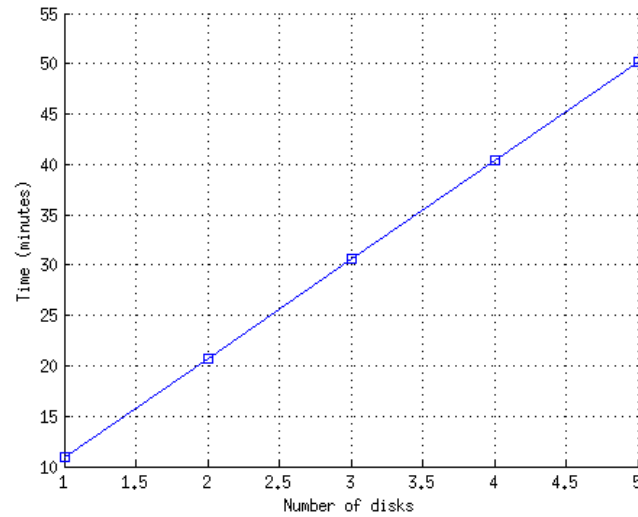


Figure 5.2: Time results for complete erasure by second strategy

CHAPTER 5. EXPERIMENTS AND RESULTS

During testing erasure with third strategy we got the results presented in the Table 5.3.

Table 5.3: Results for complete erasure by third strategy

Disks	Capacity	Time	Average time/disk	Average speed
1	18.2 GB	6 min 5 sec	6 min 5 sec	49.8 MB/s
2	36.4 GB	8 min 29 sec	4 min 15 sec	71.5 MB/s
3	54.6 GB	9 min 45 sec	3 min 15 sec	93.3 MB/s
4	72.8 GB	10 min 57 sec	2 min 44 sec	111.2 MB/s
5	91 GB	10 min 57 sec	2 min 11 sec	138.5 MB/s
6	109.2 GB	10 min 57 sec	1 min 50 sec	166.2 MB/s
7	127.4 GB	10 min 57 sec	1 min 34 sec	194.0 MB/s
8	145.6 GB	10 min 57 sec	1 min 22 sec	221.7 MB/s
9	163.8 GB	10 min 58 sec	1 min 13 sec	249.0 MB/s
10	182 GB	10 min 57 sec	1 min 6 sec	277.1 MB/s
11	200.2 GB	11 min 0 sec	1 min 0 sec	303.4 MB/s
12	218.4 GB	11 min 0 sec	55 sec	331.0 MB/s
13	236.6 GB	13 min 26 sec	1 min 2 sec	293.6 MB/s
14	254.8 GB	13 min 28 sec	58 sec	315.5 MB/s

Figure 5.3 presents the results from the Table 5.3 more clearly.

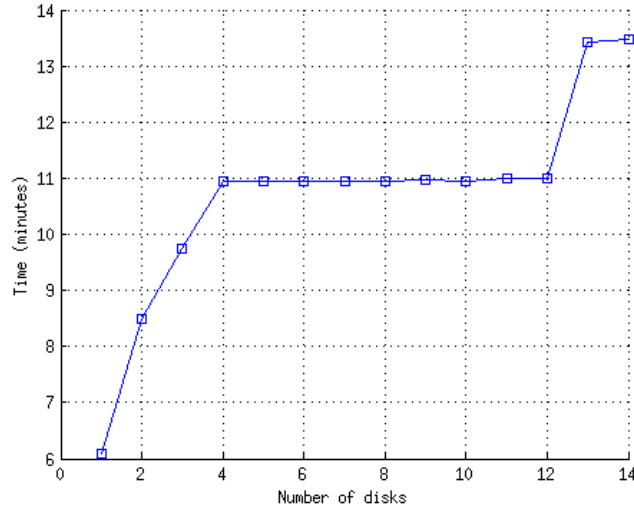


Figure 5.3: Time results for complete erasure by third strategy

From the Figure 5.3 we can see that for disks from 4 to 12 the time of complete erasure is constant. Third strategy sends the commands parallel to several disks. That means if there is one slow disk in the task, the whole time of the erasure will be slow. Therefore, we decided to calculate the time for erasure of single disk.

CHAPTER 5. EXPERIMENTS AND RESULTS

During testing erasure for single disks with third strategy we got the results presented in the Table 5.4. Moreover, in the Table 5.4 we present the time of erasure for single disk during complete erasure.

Table 5.4: Results for erasure of single disks by third strategy

Disks	Capacity	Time	Time during complete erasure
1	18.2 GB	6 min 5 sec	6 min 5 sec
2	36.4 GB	9 min 39 sec	10 min 59 sec
3	54.6 GB	9 min 44 sec	9 min 43 sec
4	72.8 GB	10 min 57 sec	10 min 59 sec
5	91 GB	10 min 55 sec	10 min 53 sec
6	109.2 GB	9 min 42 sec	9 min 42 sec
7	127.4 GB	6 min 37 sec	6 min 37 sec
8	145.6 GB	8 min 29 sec	8 min 29 sec
9	163.8 GB	6 min 38 sec	6 min 37 sec
10	182 GB	6 min 5 sec	6 min 5 sec
11	200.2 GB	10 min 13 sec	10 min 12 sec
12	218.4 GB	9 min 51 sec	10 min 31 sec
13	236.6 GB	9 min 55 sec	9 min 55 sec
14	254.8 GB	6 min 6 sec	6 min 5 sec

CHAPTER 5. EXPERIMENTS AND RESULTS

Figure 5.4 presents the results from the Table 5.4 more clearly.

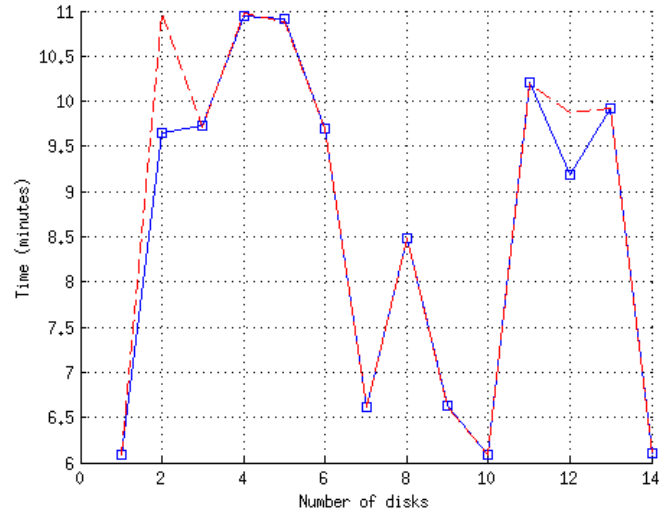


Figure 5.4: Time results for complete erasure by third strategy

By replacing disks 4,5 and 11 we even can improve our results, which are shown on the Figure 5.5.

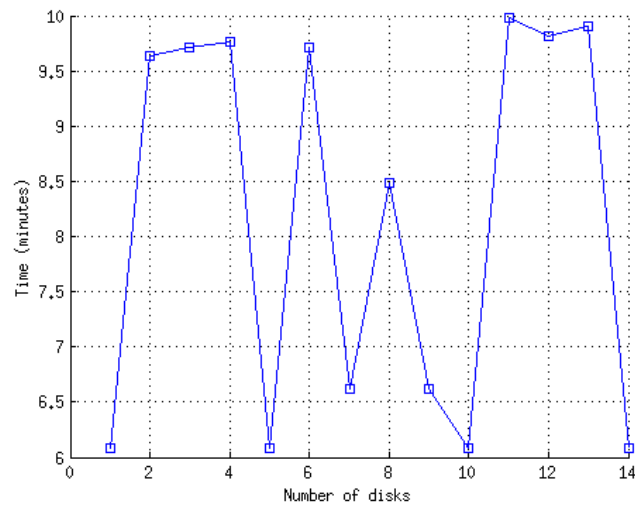


Figure 5.5: Time results for single erasure by third strategy

5.2 Results with WRITE 10 command

During research with WRITE SAME 10 command we understood that the parallel strategy works fine with that controller that is why experiments with WRITE 10 command were done only with the third strategy. Several tests were done with WRITE 10 command with transfer length 256, which is not the maximum, but seems to be optimal. During testing erasure for single disks with third strategy we got the results presented in the Table 5.5.

Table 5.5: Results for erasure of single disks by third strategy

Disks	Capacity	Time	Time during complete erasure
1	18.2 GB	19 min 52 sec	53 min 43 sec (-3)
2	36.4 GB	23 min 24 sec	45 min 6 sec (+4)
3	54.6 GB	23 min 28 sec	49 min 25 sec (+3)
4	72.8 GB	24 min 43 sec	54 min 15 sec (0)
5	91 GB	24 min 38 sec	57 min 39 sec (-1)
6	109.2 GB	23 min 28 sec	60 min 19 sec (+3)
7	127.4 GB	20 min 27 sec	54 min 27 sec (+2)
8	145.6 GB	22 min 15 sec	57 min 47 sec (-3)
9	163.8 GB	20 min 24 sec	42 min 40 sec (+2)
10	182 GB	19 min 50 sec	50 min 37 sec (-5)
11	200.2 GB	24 min 1 sec	33 min 0 sec (-5)
12	218.4 GB	23 min 32 sec	35 min 10 sec (+1)
13	236.6 GB	23 min 39 sec	38 min 14 sec (-6)
14	254.8 GB	19 min 51 sec	52 min 40 sec (-5)

CHAPTER 5. EXPERIMENTS AND RESULTS

Figure 5.6 presents the results from the Table 5.5 more clearly.

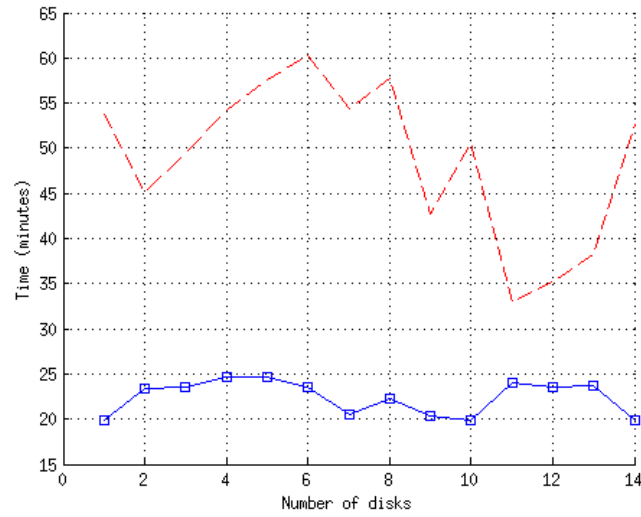


Figure 5.6: Time results for complete erasure by WRITE 10 command

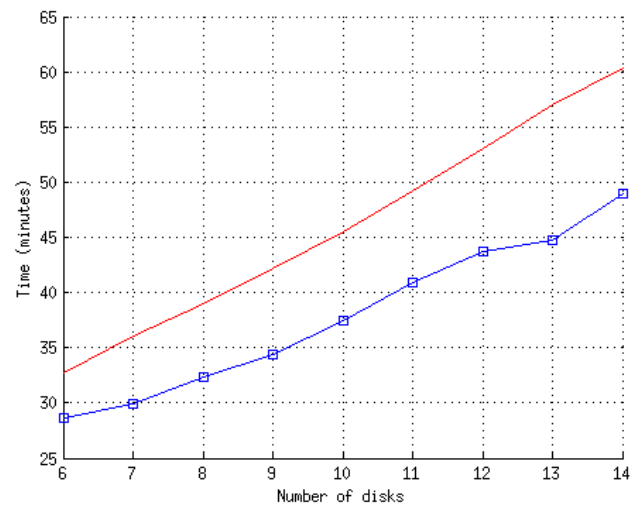


Figure 5.7: Time results for average erasure and complete erasure

CHAPTER 5. EXPERIMENTS AND RESULTS

Figure 5.8 presents the time results for the disk 1 depending on the different transfer length.

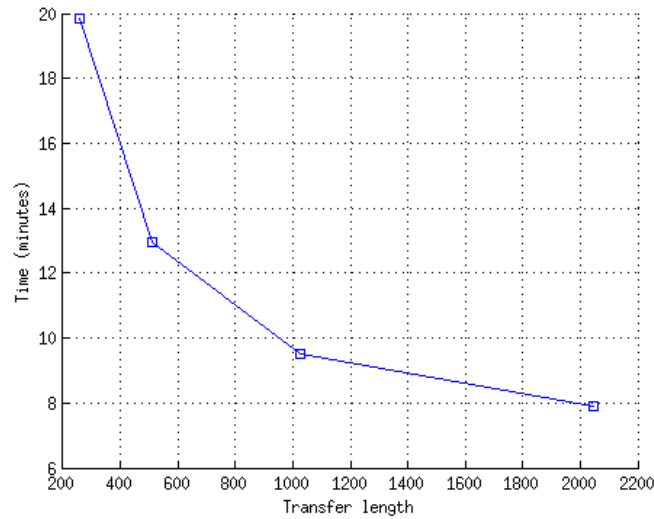


Figure 5.8: Time results for the disk 1

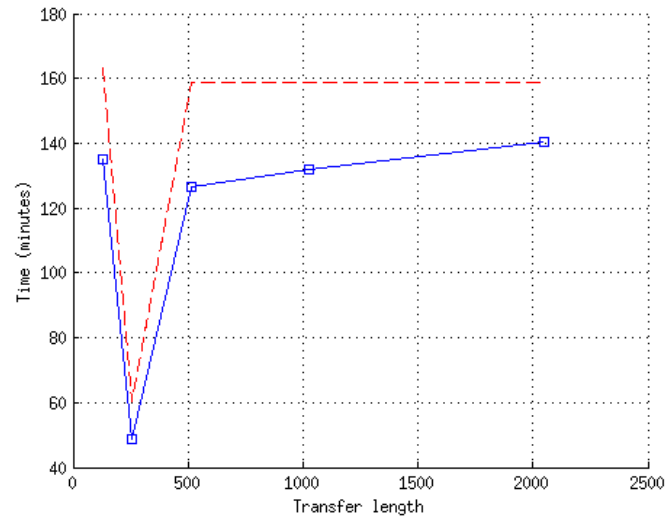


Figure 5.9: Time results for average erasure and complete erasure depending on the transfer length

CHAPTER 5. EXPERIMENTS AND RESULTS

Figure 5.10 presents the time results of complete erasure with transfer length 256 and 512.

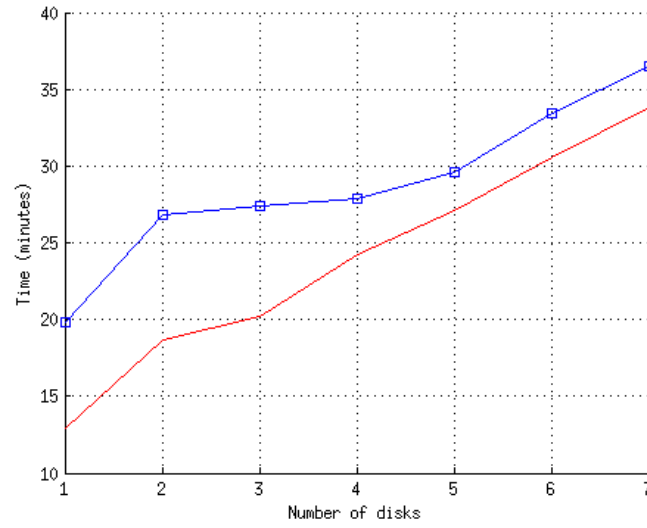


Figure 5.10: Results for complete erasure depending on the transfer length

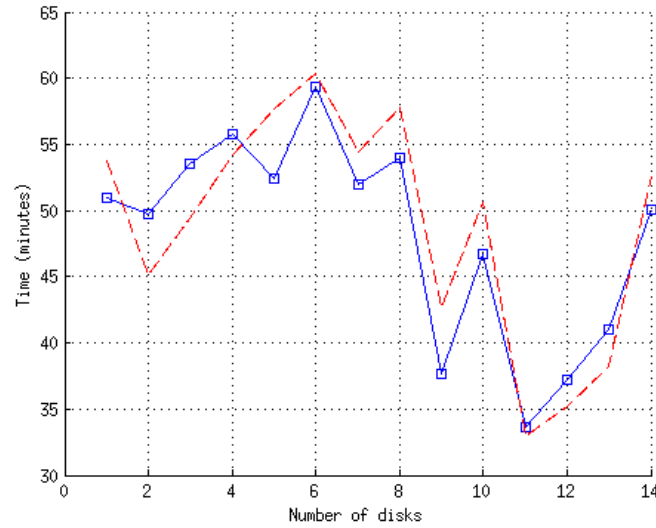


Figure 5.11: Results after replacement of the disks 4,5 and 11.

5.3 Analysis of the results

After testing with WRITE SAME 10 command we can notice that the second strategy is working slower than the first one, which is obvious because it needs to change the devices quite often. The third strategy works without any problems, but the results of several disks are quite different, which gives an idea that we did not consider all the parameters or it depends just on the hardware.

We can see from the Figures 5.1 and 5.2 that the time of first and second strategies is increasing linearly. We could not say the same about the third strategy because the disks have different speed or some other parameters disturb disk work on the maximum speed. Figure 5.4 shows that the difference is very low between the times of erasure for single disk and for single disk during complete erasure. That means that the bus is not a bottleneck in the situation, when we send WRITE SAME 10 command and 14 disks are connected to the controller, and only rotation speed of the disks is a limiting factor. This fact is proved one more time by the Figure 5.5, where we show the results of erasure time for single disk during complete erasure after replacement of 3 slow disks. The results shows that the time of complete erasure decreased from 11 to 10 minutes. But still it is strange that the difference of the time of erasure for several disks can be even 4 minutes, which is almost half time of the complete erasure, but the disk parameters are completely the same.

The testing results do not depend on the application running and the erasure time is not divided randomly to different disks. Table 5.5 proved this fact by results from the same test running several times, which gave the whole difference with -13 seconds for erasure of 14 disks by WRITE 10 command.

That means that even for erasure of the fastest disk 13 seconds will take only 0.01 of whole erasure. That is why we could ignore this difference.

We can notice that during complete erasure by WRITE SAME 10 command disks 2, 4 and 5 showed the worst time. In the same time during complete erasure by WRITE 10 command disks 6 and 8 showed the worst time. That gives an idea that there are some other parameters that also depending on the time of erasure. Moreover, Figure 5.11 shows how changes the results after replacement disks 4,5 and 11. We can conclude that time for complete erasure decreased for 1 minute and the same disks show the maximum and minimum times. It is also strange that during single erasure by WRITE 10 command disks 1 and 14 were the fastest and during complete erasure their results are very closed to very slow ones.

Figures 5.8 and 5.9 show the results of testing different transfer length, which is one of the parameters that we can influence. To summarize these results we can say that for 14 disks transfer length 256 is definitely the optimal value, but it is not correct for 1 disk. Using the value 2048 we could erase one disk much faster. That means that the load balancing application could help with analysis, which transfer length to take, for making the erasure faster.

Bibliography

- [1] *Techniques for Increasing PCI Performance*. Intel Corporation, 1999.
- [2] *RAID Technology Overview*. Hewlett-Packard Development Company, L.P., 2007.
- [3] *HP Smart Array 642 Controller*. Hewlett-Packard, 2009.
- [4] *Which RAID Level is Right for Me?* Adaptec, Inc., 2010.
- [5] Ali M. Alakeel. *A Guide to Dynamic Load Balancing in Distributed Computer Systems*. IJCSNS International Journal of Computer Science and Network Security, VOL.10 No.6, 2010.
- [6] Michael Arellano. *Ultra320 SCSI: New Technology - Still SCSI*. SCSI Trade Association, 2001.
- [7] T10 Technical Committee. *Working Draft. Revision 35. Information technology - SCSI Block Commands - 3 (SBC-3)*. Technical Committee of Accredited Standards Committee INCITS, 2012.
- [8] T13 Technical Committee. *Working Draft. Revision 4c. AT Attachment 8 - ATA/ATAPI Command Set (ATA8-ACS)*. American National Standard of Accredited Standards Committee INCITS, 2007.

BIBLIOGRAPHY

- [9] Carlo Kopp. *An Introduction to SCSI Disk Performance*. Dr Carlo Kopp's Industry Publications, 2005.
- [10] Abraham Silberschatz; James Peterson. *Operating System Concepts, 8th edition*. Addison-Wesley, 1982.
- [11] Gary Field; Peter M. Ridge. *The Book of SCSI, 2nd Edition*. Publishers Group West, 2000.
- [12] Ramakrishna Karedla; J. Spencer Love; Bradley G. Wherry. *Caching Strategies to Improve Disk System Performance*. Computer, 1994.
- [13] B. Litow; S.H. Hosseini; K. Vairavan; G.S. Wollfe. *Performance Characteristics of a Load Balancing Algorithm*. Journal of Parallel and Distributed Computing, 1995.
- [14] Andrew S. Tanenbaum; Albert S. Woodhull. *Operating Systems: Design and Implementation (Third Edition)*. Pearson Education, Inc., 2006.