



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 3

Issue: 1

Month of publication: January 2015

DOI:

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Enhanced Data Forensics Technique Based On Parallel Processing and Multilevel File Binning

Neha Razdan¹, Sachin Majithia²

¹Student, ²Assistant Professor, Department of Information Technology, CEC, Landran

Abstract: Data forensic comes from forensic science. Forensic science is the scientific method of gathering and examining information about the past. This is especially important in law enforcement where forensics is done in relation to criminal or civil law, but forensics are also carried out in other fields, such as astronomy, archaeology, biology and geology to investigate ancient times. Data forensics, often used interchangeably with computer forensics, is essentially the study of digital data and how it is created and used for the purpose of an investigation. Data forensics is part of the greater discipline of forensics, in which various types of evidence are studied to investigate an alleged crime. In this paper the main focus on optimizing overall data size using binning and de-duplication or filtering of files. All the files are arranged with their constant size irrespective of their size. Our main work is using parallel processing because it helps to take less time to process, also response the time over the network

Keywords: Data forensics, parallel processing, sequential processing, binning, filtering, kmp algorithm

I. INTRODUCTION

A. Data Forensics

Forensic science is the scientific method of gathering and examining information about the past. This is especially important in law enforcement where forensics is done in relation to criminal or civil law, but forensics are also carried out in other fields, such as astronomy, archaeology, biology and geology to investigate ancient times. In modern use, the term forensics in the place of forensic science can be considered correct, as the term forensic is effectively a synonym for legal or related to courts. However, the term is now so closely associated with the scientific field that many dictionaries include the meaning that equates the word forensics with forensic science. Data forensics can involve many different tasks, including data recovery or data tracking. Data forensics might focus on recovering information on the use of a mobile device, computer or other device. It might cover the tracking of phone calls, texts or emails through a network. Digital forensics investigators may also use various methodologies to pursue data forensics, such as decryption, advanced system searches, reverse engineering, or other high-level data analyses. NIST defines digital forensics as an applied science for “the identification, collection, examination, and analysis of data while preserving the integrity of the information and maintaining a strict chain of custody for the data”. Figure 1 illustrates the process flow of digital forensics.



Figure1: Process flow of Digital or Data Forensics

B. De-duplication

Data deduplication or Single Instancing essentially refers to the elimination of redundant data. In the deduplication process, duplicate data is deleted, leaving only one copy (single instance) of the data to be stored. However, indexing of all data is still retained should that data ever be required. Data de-duplication is a term used to describe an algorithm or technique that eliminates duplicate copies of data from storage. Data deduplication is commonly performed on secondary storage systems such as archival and backup storage. One necessary condition for the data deduplication is that data de-duplication should be scalable. It means that de-duplication should be elastic. It doesn't affect to the overall storage structure.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

The two major categories of data de-duplication:

- 1) *Offline Data de-duplication*: In an offline deduplication state, first data is written to the storage disk and deduplication process take place at a later time.
- 2) *Online Data de-duplication*: In an online deduplication state, replicate data is deleted before being written to the storage disk.

C. Binning

Binning is a way to group a number of more or less continuous values into a smaller number of "bins". Binning is also a combination of two or more ccd images sensor pixels to form a new "super-pixel" prior to readout and digitizing. It is a preprocessing technique used to reduce the effects of minor observation errors

D. Knuth Morris Pratt Algorithm

In computer science, the Knuth–Morris–Pratt string searching algorithm (or KMP algorithm) searches for occurrences of a "word" W within a main "text string" S by employing the observation that when a mismatch occurs, the word itself embodies sufficient information to determine where the next match could begin, thus bypassing re-examination of previously matched characters. The algorithm was conceived in 1974 by Donald Knuth and Vaughan Pratt, and independently by James H. Morris. The three published it jointly in 1977. In studying the worst-case performance of the brute-force and boyer-moore pattern matching algorithms on specific instances of the problem; we should notice a major inefficiency. Specifically, we may perform many comparisons while testing a potential placement of the pattern against the text, yet if we discover a pattern character that does not match in the text, then we throwaway all the information gained by these comparisons and start over again from scratch with the next incremental placement of the pattern. The Knuth-Morris-Pratt (or "KMP") algorithm avoids this waste of information and, in so doing, it achieves a running time of $O(n + m)$, which is optimal in the worst case. That is, in the worst case any pattern matching algorithm will have to examine all the characters of the text and all the characters of the pattern at least once.

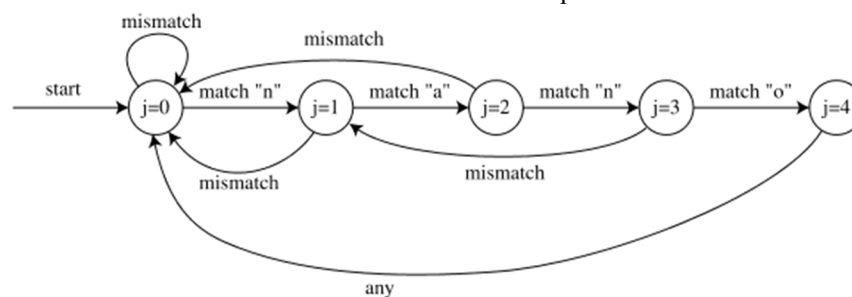


Figure1.1: Implementation of KMP algorithm

E. Example of Knuth Morris Pratt Algorithm

To illustrate the ideas of the algorithm, we consider the following example:

$$T = \text{xyxxyxyxyxyxyxyxyxyxyxy}$$

and

$$P = \text{xyxyxyxyxx}$$

At a high level, the KMP algorithm is similar to the naive algorithm: it considers shifts in order from l to $n-m$, and determines if the pattern matches at that shift. The difference is that the KMP algorithm uses information gleaned from partial matches of the pattern and text to skip over shifts that are guaranteed not to result in a match. Suppose that, starting with the pattern aligned underneath the text at the leftmost end, we repeatedly "slide" the pattern to the right and attempt to match it with the text. Let's look at some examples of how sliding can be done. With numbering, to make it easier to follow.

1. Consider the situation when $P[1, \dots, 3]$ is successfully matched with $T[1, \dots, 3]$. We then find a mismatch: $P[4] \neq T[4]$. Based on

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

our knowledge that $P[1, \dots, 3] = T[1, \dots, 3]$, and ignoring symbols of the pattern and text after position 3, what can we deduce about where a potential match might be? In this case, the algorithm slides the pattern 2 positions to the right so that $P[1]$ is lined up with $T[3]$. The next comparison is between $P[2]$ and $T[4]$.

2. Since $P[2] = T[4]$, the pattern slides to the right again, so that the next comparison is between $P[1]$ and $T[4]$.

3. At a later point, $P[1, \dots, 10]$ is matched with $T[6, \dots, 15]$. Then a mismatch is discovered: $P[11] \neq T[16]$. Based on the fact that we know $T[6, \dots, 15] = P[1, \dots, 10]$ (and ignoring symbols of the pattern after position 10 and symbols of the text after position 15), we can tell that the first possible shift that might result in a match is 12. Therefore, we will slide the pattern right, and next ask whether $P[1, \dots, 11] = T[13, \dots, 23]$. Thus, the next comparisons done are $P[4] = T[16]$, $P[5] = T[17]$, and $P[6] = T[18]$ and so on, as long as matches are found. In studying the worst-case performance of the brute-force and boyer-moore pattern matching algorithms on specific instances of the problem, we should notice a major inefficiency. Specifically, we may perform many comparisons while testing a potential placement of the pattern against the text, yet if we discover a pattern character that does not match in the text, then we throw away all the information gained by these comparisons and start over again from scratch with the next incremental placement of the pattern. The Knuth-Morris-Pratt (or "KMP") algorithm avoids this waste of information and, in so doing, it achieves a running time of $O(n + m)$, which is optimal in the worst case. That is, in the worst case any pattern matching algorithm will have to examine all the characters of the text and all the characters of the pattern at least once. The KMP pattern matching algorithm incrementally processes the text string T comparing it to the pattern string P . Each time there is a match, we increment the current indices. On the other hand, if there is a mismatch and we have previously made progress in P , then we consult the failure function to determine the new index in P where we need to continue checking P against T . Otherwise (there was a mismatch and we are at the beginning of P), we simply increment the index for T (and keep the index variable for P at its beginning). We repeat this process until we find a match of P in T or the index for T reaches n , the length of T (indicating that we did not find the pattern P in T). The main part of the KMP algorithm is the while-loop, which performs a comparison between a character in T and a character in P for each iteration. Depending upon the outcome of this comparison, the algorithm either moves on to the next characters in T and P , consults the failure function for a new candidate character in P , or starts over with the next index in T . The correctness of this algorithm follows from the definition of the failure function. The skipped comparisons are actually unnecessary, for the failure function guarantees that all the ignored comparisons are redundant—they would involve comparing characters that are already known to match.

II. LITERATURE SURVEY

Raju Bhukya, DVLN Somayajulu The current research in life science area is producing large amount of genetic data. DNA related pattern matching in a given sequence is one of the fundamental problems in computational biology. Biologists often search DNA sequences to identify use full information available from protein and genes for the functional and structural behavior. Many algorithms have been proposed but more efficient and robust methods are needed for the exact pattern matching algorithms. In the proposed work an index based sequential multiple pattern matching using pair indexing algorithm gives very good performance when compared with some of the existing algorithms. The current technique avoids unnecessary DNA comparisons as a result the number of comparisons and CPC ratio gradually decreases and overall performance increases.

Zhike Zhang, Deepavali Bhagwat, Witold Litwin, Darrell Long, Thomas Schwarz, S.J. Many modern storage systems use deduplication in order to compress data by avoiding storing the same data twice. Deduplication needs to use data stored in the past, but accessing information about all data stored can cause a severe bottleneck. Similarity based deduplication only accesses information on past data that is likely to be similar and thus more likely to yield good deduplication. We present an adaptive deduplication strategy that extends Extreme Binning and investigate theoretically and experimentally the effects of the additional bin accesses.

Aniello castiglione (member, ieee), giuseppe cattaneo, giancarlo de maio, and alfredo de santis (member, ieee) Digital evidence is increasingly used in juridical proceedings. In some recent legal cases, the verdict has been strongly influenced by the digital evidence proffered by the defense. Digital traces can be left on computers, phones, digital cameras, and also on remote machines belonging to ISPs, telephone providers, companies that provide services via Internet such as YouTube, Face book, Gmail, and so on. This paper presents a methodology for the automated production of predetermined digital evidence, which can be leveraged to forge a digital alibi. It is based on the use of automation, a program meant to simulate any common user activity. In

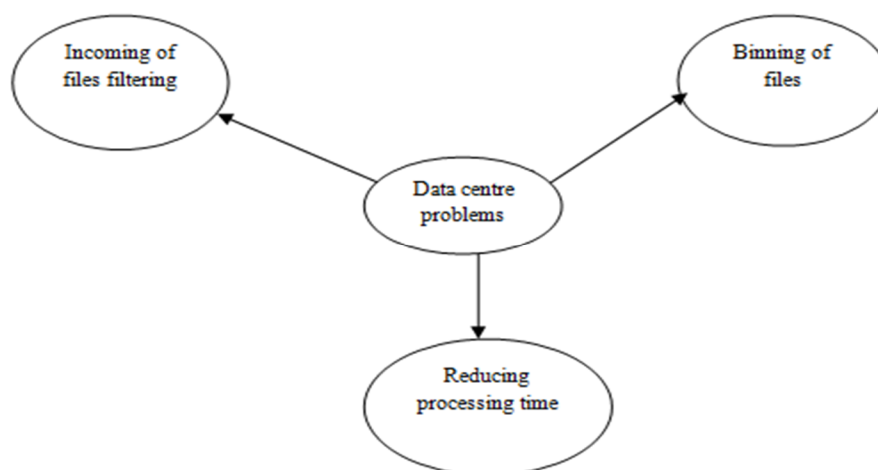
International Journal for Research in Applied Science & Engineering Technology (IJRASET)

addition to wanted traces, the automation may produce a number of unwanted traces, which may be disclosed upon a digital forensic analysis. These include data reminiscent of suspicious files, as well as any kind of logs generated by the operating system modules and services. The proposed methodology describes a process to design, implement, and execute the automation on a target system, and to properly handle both wanted and unwanted evidence. Many experiments with different combinations of automation tools and operating systems are conducted. This paper presents an implementation of the methodology through VBScript on Windows 7. A forensic analysis on the target system is not sufficient to reveal that the alibi is forged by automation. These considerations emphasize the difference between digital and traditional evidence. Digital evidence is always circumstantial, and therefore it should be considered relevant only if supported by stronger evidence collected through traditional investigation techniques. Thus, a Court verdict should not be based solely on digital evidence.

Sampath S1, Bharat Bhushan Sagar2, Nanjesh B R3 Parallel computing operates on the principle that large problems can often be divided into smaller ones, which are then solved concurrently to save time (wall clock time) by taking advantage of non-local resources and overcoming memory constraints. The main aim is to form common single node architecture for both MPI and PVM, which demonstrates the performance gain and losses achieved through parallel processing using MPI and PVM as separate cases. We also demonstrate the performance dependency of parallel applications on RAM of the nodes (desktop PCs) used in parallel computing. This can be realized by implementing the parallel applications like solving matrix multiplication problem, using MPI and PVM separately. The single node architecture consists of a client, a master, capable of handling requests from the client, and a slave, capable of accepting problems from the master and sending the solution back. The master and the slave communicate with each other using MPICH-2 and PVM3.4.6 when computation is under MPI and PVM respectively. The master will monitor the progress and be able to compute and report the time taken to solve the problem, taking into account the time spent in assigning the problem into slave and sending the results along with the communication delays. We aim to evaluate and compare these statistics of both the cases to decide which among MPI and PVM gives faster performance and also compare with the time taken to solve the same problem in serial execution to demonstrate communication overhead involved in parallel computation. We aim to compare and evaluate the statistics obtained for different sizes of RAM under parallel execution in a single node involving only two cores, where one acts as master and other as slave. We also show the dependency of serial execution on RAM for the same problem by executing its serial version under different sizes of RAM.

III. PROBLEM FORMULATION

From the literature survey, it has been concluded that there is issue of investigation of phishing data from the user access pattern to handle the situation when incoming files are of dynamic file sizes. Their separation needs to be done on the basis of file sizes just in two categories i.e. heavy weight and light weight files with the help of binning. In the process of binning the heavy weight files are handled and are managed more systematically more than 3 to 4 categories of different file sizes. The amount of data at the data centre is very huge it needs to be processed keeping the time parameter in mind so the data is to be processed parallel using some string matching algorithm.



International Journal for Research in Applied Science & Engineering Technology (IJRASET)

Figure3: Data centre problems

IV. PROPOSED WORK

This chapter describes the implementation setups used for implementing data security using KMP algorithm and enhanced parallel processing techniques in data forensics.

A. Implementation Procedure

As discussed in literature survey, there is a problem of processing of data in serial, which is the main challenge in data forensics. So there is a need for improvement in the data security techniques for better efficiency and more efficient algorithm for enhancing security.

In this, a kmp- algorithm is implemented on the MATLAB. For implementation, we have done the following steps:

STEP 1: Data arrives at the data centre

STEP 2: Segmentation has been done for the data because it reduces data files into smaller chunks which are easier to transfer over Internet.

STEP 3: After segmentation data compression has done using de-duplication of files which focuses on optimizing overall data files size. This technique reduces cost and increased storage efficiency, while also improving the user experience.

STEP 4: The files are then segregated into different bins depending on their size. When binning is done with constant size for all files irrespective of their size, it either leads to a large number of segments or very few segments. The files are also processed by using parallel processing because it take less time to process and save the segments, also response the time over the network.



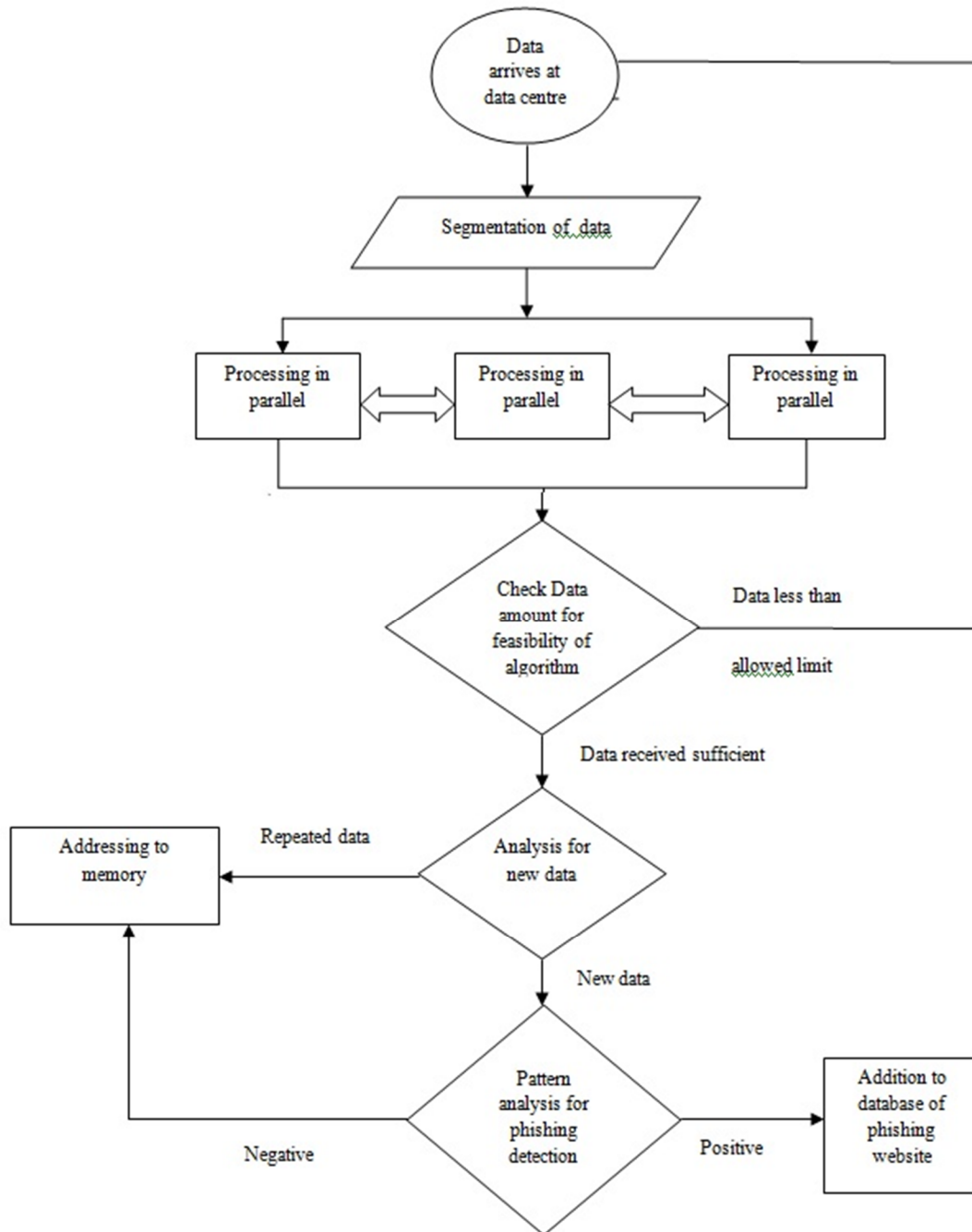
Screen shot4: Shows the parallel processing of files.

STEP 5: Now check the amount of data whether it is feasible for algorithm if the given data is not feasible or less than for the algorithm then again the new data will enter at the data centre till the amount of data is feasible for the algorithm.

STEP 6: If the data which we received is sufficient, then we analyze a new data. If this new data was also repeated then this data will addressing to the memory.

STEP 7: Then the new data will analyzing for the phishing detection. If this data is negative it will addressing to the memory and if it is positive it will goes in the database of phishing websites in which other phishing websites is present. By using the KMP (Knuth Morris Pratt) algorithm the string of data is matched whether it is phishing data or not.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)



Flowchart4.1: Shows the processing of data

B. Implementation Setup

The implementation of our work is done in the MATLAB. The software used to implement our proposed algorithm will be MATLAB. MATLAB is a high-level language and interactive environment for numerical computation, visualization, and programming. Using MATLAB, we can analyze data, develop algorithms, and create models and applications. The language, tools, and built-in math functions enable us to explore multiple approaches and reach a solution faster than with spreadsheets or traditional programming languages, such as C/C++ or Java™.

MATLAB can be used for a range of applications, including signal processing and communications, image and video processing,

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

control systems, test and measurement, computational finance, and computational biology.

We chose MATLAB to implement our algorithm because of the following reasons:

- 1) High-level language for numerical computation, visualization, and application development
- 2) Interactive environment for iterative exploration, design, and problem solving
- 3) Built-in graphics for visualizing data and tools for creating custom plots
- 4) Development tools for improving code quality and maintainability and maximizing performance
- 5) Tools for building applications with custom graphical interfaces

Functions for integrating MATLAB based algorithms with external applications and languages such as C, Java,.N Various toolboxes used in MATLAB are as follows:

- 6) Communications System Toolbox
- 7) Data Acquisition Toolbox
- 8) Database Toolbox
- 9) Simulink

V. METHODOLOGY

A. Methodology Used

This section describes a methodology and its related parameters, experiment factors.

- 1) *System Parameters*- The experiments are conducted using Intel Dual Core 64-bit processor with a minimum of 256 MB RAM or higher and hard disk 10 GB or higher. The implementation is done in the MATLAB. In this report, the KMP algorithm has been implemented parallel processing of the above configuration.
- 2) *Experimental Parameters*- In order to evaluate the effectiveness of the proposed technique in the data forensics. The two parameters are must be determined.
 - a) *Time processing*:-The time processing is defined as the two sets of functions for measuring absolute elapsed time: clock and etime, and tic and toc. The clock and etime functions use dates to track elapsed time, and are useful if you need accuracy to only .01 second. Originally, the TIC and TOC functions relied on clock for the time and had the same accuracy. We use a tic toc function to calculate the time processing.
 - b) *Complexity*: - It is used to characterize something with many parts where those parts interact with each other in multiple ways. Complexity of an object or system is a relative property. For instance, for many functions (problems), such a computational complexity as time of computation is smaller when multitape Turing machines are used than when Turing machines with one tape are used. Random Access Machines allow one to even more decrease time complexity (Green law and Hoover 1998: 226), while inductive Turing machines can decrease even the complexity class of a function, language or set (Burgin 2005). This shows that tools of activity can be an important factor of complexity.

VI. RESULTS ANALYSIS

In this section we evaluate the effectiveness of proposed technique that is by using the parameters. We have generated the graphs of our proposed method by using two parameters first is time processing, second is complexity.

In these graphs we have a number of files which is to be processed .The table is also given below in which number of files and the time taken to complete its processing. The number of files is in KB and the time is given in seconds.

A. Scenerio1 (Sequential Processing)

In this scenario, files are processed in sequential manner. In sequential processing files take much more time to process and the time processing of all the files is very large.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

Table No6.1 Shows the time taken to process the different files

No of files(in kb)	5000	4000	3000	2000	1000	500	250	100
Time (in seconds)	95.55	80.53	76.49	73.63	63.54	50.44	45.42	40.39

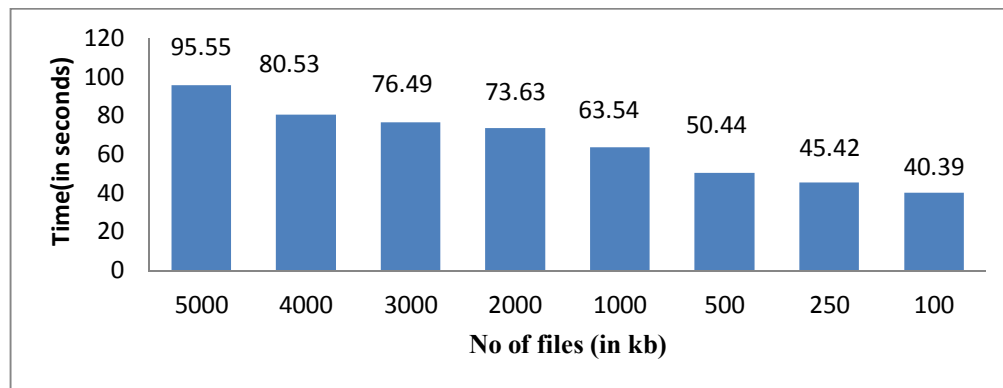


Figure6.2: Represents the graph in sequential processing

B. Scenerio2. (Parallel Processing)

In this scenario, files are processed in parallel manner. In parallel processing, files take small amount of time to process.

Table No6.3: Shows the time taken to process the different files

No of files(in kb)	5000	4000	3000	2000	1000	500	250	100
Time (in seconds)	75.55	66.64	59.54	45.44	39.36	35.32	30.15	27.22

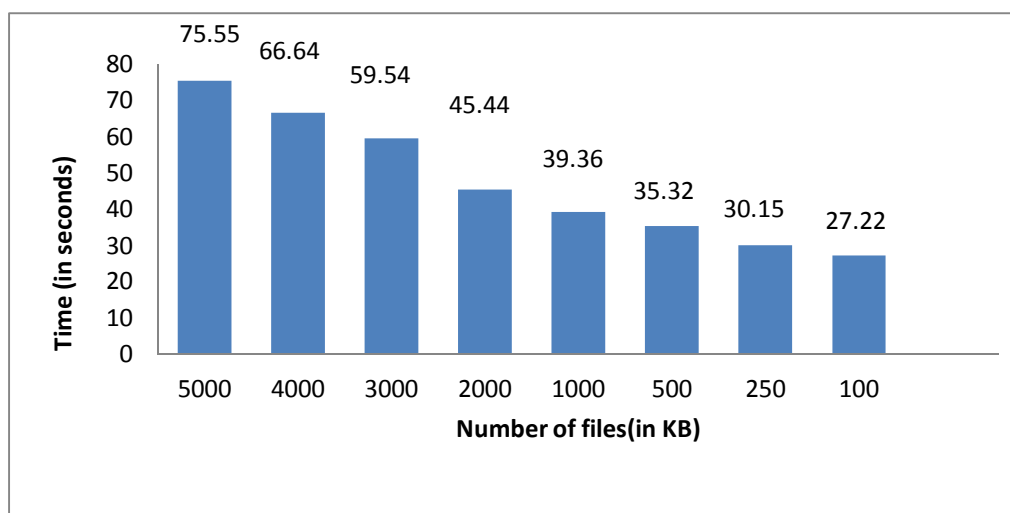


Figure6.4 Represents graph in parallel

With these two graphs we can show that time taken in parallel is much faster than the sequential. Now days we have a number of

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

files in the data that is to be accessed. To make our parallel processing better we use the idea of binning and filtering. With the help of this our time become more less. We have taken a same number of files for processing in parallel by use the concept of binning and filtering.

C. Scenerio3. (Parallel Processing)

In this scenario, files are processed in parallel manner using binning and filtering. A lot of time must be saved with the help of binning and filtering.

Table No6.5 Shows the time taken to process the different files

No of files(in kb)	5000	4000	3000	2000	1000	500	250	100
Time (in seconds)	50.15	40.44	32.12	23.21	19.14	15.13	13.11	10.9

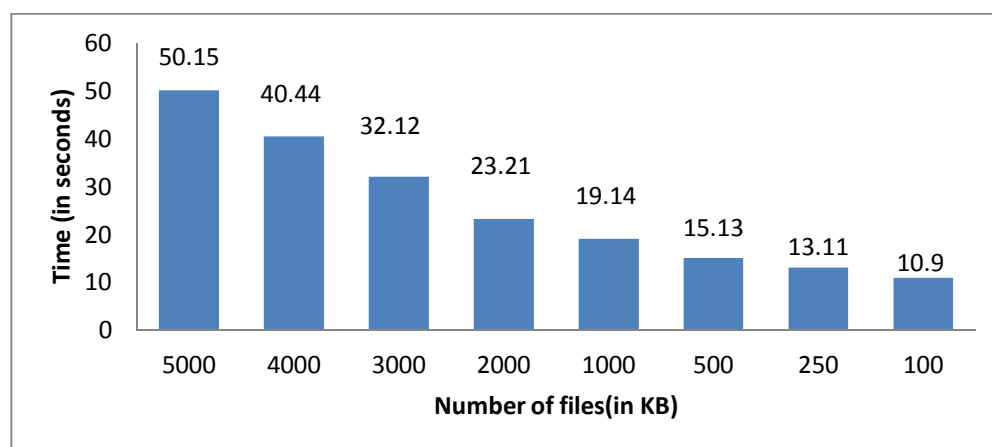


Figure6.6 Represents graph in parallel by using binning and filtering

VII. CONCLUSION

In this paragraph, the result of the simulations that we have mentioned here we conclude that the simulations have proved to improve upon the time factor that was utterly needed for the kind of data processing ability that is needed in this task. Earlier only binning and filtering of files were done to enhance the capacity and efficiency of the system. But as an enhanced technology we have used the idea of parallel processing instead of sequential processing ultimately reducing the overall time required to process the data. Results of the parallel processing have been much better than the previous results.

VIII. FUTURE SCOPE

The field of cloud computing forensics and incident response (Computer security incident management which involves monitoring and proactive detection of security threats) is new field of research that attracts more and more scientists. It poses a lot of challenges due to exponential rise in dynamic being generated and accessed data around the world and also distributed nature of cloud. In our approach we have presented a solution to real time data access patterns on static database. In future more innovative techniques would be required for dynamic database as it spreads over multiple data centers.

REFERENCES

- [1] Irfan Khan. "Prefix Span Algorithm Based on Multiple Constraints for Mining Sequential Patterns", International Journal of Computer Science and Management Research, Vol. 1, Issue 5, December 2012.
- [2] B.Mallik and D.garg, "CFM Prefix Span: A pattern growth algorithm incorporating monetary and compactness", pp 4509-4555, July 2012.
- [3] S. Zawoad and R. Hasan, "Towards building proofs of past data possession in cloud forensics," ASE Science Journal, 2012
- [4] K. Ruan, J. Carthy, T. Kechadi, and M. Crosby, "Cloud forensics: An overview," in proceedings of the 7th IFIP International Conference on Digital Forensics, 2011.

International Journal for Research in Applied Science & Engineering Technology (IJRASET)

- [5] S. Zawoad, A.K. Dutta and R. Hasan, "Salas: Secure Logging-as-a-Service for Cloud Forensics", Symposium on Information, Computer and Communications Security (ASIACCS), 2013
- [6] R. Marty, "Cloud Application Logging for Forensics", Proceedings of the 2011 ACM Symposium on Applied Computing, 2011
- [7] G. Sibiya, H. Venter and T. Fogwill, "Digital forensic framework for a cloud environment", Proceedings of the 2012 Africa Conference, 2012
- [8] Fabio Marturana ,Simone Tacconi ,Gianluigi Me, "A case study on digital forensics in the cloud", International Conference on Cyber-Enabled Distributed Computing and Knowledge Discover,2012
- [9] Shahrzad Zargari ,David Benford , "Cloud Forensics: Concepts, Issues, and Challenges", Third International Conference on Emerging Intelligent Data and Web Technologies,2012
- [10] Guangxuan Chen', Yanhui nu', Panke Qin', Jin Du2 , "Suggestion To Digital Forensics In Cloud Computing Era", Proceedings of IC-NIDC2012
- [11] AmitChhabra, Gurvinder Singh "A Cluster Based Parallel Computing Framework (CBPCF) for Performance Evaluation of Parallel Applications", International Journal of Computer Theory and Engineering, Vol. 2, No. 2 April, 2010.
- [12] Rafiqul Zaman Khan, MdFiroj Ali, "A Comparative Study on Parallel Programming Tools in Parallel Distributed Computing System: MPI and PVM", Proceedings of the 5th National Conference; INDIACom-2011.
- [13] Sampath S, Sudeepa K.B, Nanjesh B R "Performance Analysis and Evaluation of Parallel Applications using a Cluster Based Parallel Computing Framework", International Journal of Computer Science and Information Technology Research Excellence (IJCSITRE), Vol.2, Issue 1, Jan-Feb 2012.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)