



To Detect Terror-Related Activities on the Web using Data Mining Techniques

Ishwarya. M¹, Mozibur Raheman Khan.²

Research Scholar, Department of Computer Science, Jamal Mohamed College, Tiruchirapalli, Tamilnadu, India ¹

Assistant Professor, Department of Computer Science, Jamal Mohamed College, Tiruchirapalli, Tamilnadu, India ²

Abstract: *An innovative knowledge-based methodology for terrorist detection by using Web traffic content as the audit information is presented. The proposed methodology learns the typical behavior ('profile') of terrorists by applying a data mining algorithm to the textual content of terror-related Web sites. The resulting profile is used by the system to perform real-time detection of users suspected of being engaged in terrorist activities. The Receiver-Operator Characteristic (ROC) analysis shows that this methodology can outperform a command-based intrusion detection system.*

Keywords: *Data Mining, User Modeling, Terrorist Trend Detection, Anomaly Detection, Activity Monitoring.*

I. INTRODUCTION

Terrorist cells are using the Internet infrastructure to exchange information and recruit new members and supporters (Lemos 2002; Kelley 2002). For example, high-speed Internet connections were used intensively by members of the infamous 'Hamburg Cell' that was largely responsible for the preparation of the September 11 attacks against the United States (Corbin 2002). This is one reason for the major effort made by law enforcement agencies around the world in gathering information from the Web about terror-related activities. It is believed that the detection of terrorists on the Web might prevent further terrorist attacks (Kelley 2002).

One way to detect terrorist activity on the Web is to eavesdrop on all traffic of Web sites associated with terrorist organizations in order to detect the accessing users based on their IP address.

Unfortunately it is difficult to monitor terrorist sites (such as 'Azzam Publications' (Corbin 2002)) since they do not use fixed IP addresses and URLs. The geographical locations of Web servers hosting those sites also change frequently in order to prevent successful eavesdropping. To overcome this problem, law enforcement agencies are trying to detect terrorists by monitoring all ISPs traffic (Ingram 2001), though privacy issues raised still prevent relevant laws from being enforced. In this paper a new methodology to detect users accessing terrorist related information by processing all ISPs traffic is suggested. The main design criteria for the proposed methodology are:

1. Training the detection algorithm should be based on the content of existing terrorist sites and known terrorist traffic on the Web.
2. Detection should be carried out in real-time. This goal can
3. be achieved only if terrorist information interests are presented in a compact manner for efficient processing.
4. The detection sensitivity should be controlled by user-defined parameters to enable calibration of the desired detection performance.

The paper is organized as follows. In the second section a brief review of intrusion detection systems, cluster analysis, and the vector space model which form the theoretical foundation behind the new methodology are presented.

In the third section the new content-based detection methodology is described in detail. The fourth section illustrates the methodology through an initial case study designed to test its feasibility.

The fifth section elaborates on the ways a system based on the new methodology can be deployed. Finally, the sixth section outlines directions for the next stages of the research.

II. BACKGROUND

This research integrates issues from the research fields of computer security (Intrusion Detection Systems), information retrieval (the vector-space model), and data mining (cluster analysis). The following subsections include a brief overview of these topics and their relation to the newly proposed methodology.

A. Intrusion Detection System

An *Intrusion Detection System* (IDS) constantly monitors actions in a certain environment and decides whether

they are part of a possible hostile attack or a legitimate use of the environment (Debar *et al.* 1999). The environment may be a computer, several computers connected in a network or the network itself. The IDS analyzes various kinds of information about actions emanating from the environment and evaluates the probability that they are symptoms of intrusions. Such information includes, for example, configuration information about the current state of the system, audit information describing the events that occur in the system (e.g., event log in Windows XP), or network traffic.

Several measures for evaluating an IDS have been suggested (Debar *et al.* 1999; Richards 1999; Spafford and Zamboni 2000; Balasubramaniyan *et al.* 1998). These measures include accuracy, completeness, performance, efficiency, fault tolerance, timeliness, and adaptivity. The more widely used measures are the True Positive (TP) rate, that is, the percentage of intrusive actions (e.g., terror-related pages) detected by the system, False Positive (FP) rate which is the percentage of normal actions (e.g., pages viewed by normal users) the system incorrectly identifies as intrusive, and Accuracy which is the percentage of alarms found to represent abnormal behavior out of the total number of alarms. In the current research TP, FP and Accuracy measures were adopted to evaluate the performance of the new methodology

B. Vector-Space Model

One major issue in this research is the representation of textual content of Web pages. More specifically, there is a need to represent the content of terror-related pages as against the content of a currently accessed page in order to efficiently compute the similarity between them. This study will use the vector-space model commonly used in Information Retrieval applications (Salton 1989; Salton *et al.* 1975) for representing terrorists' interests and each accessed Web page. In the vector-space model, a document d is represented by an n -dimensional vector $d = (w_1, w_2, \dots, w_n)$, where w_i represents the frequency-based weight of term i in document d . The similarity between two documents represented as vectors may be computed by using one of the known vector distance measuring methods such as Euclidian distance or Cosine (Boger, *et al.* 2001; Pierrea, *et al.* 2000). In this study each Web page is considered as a document and is represented as a vector. The terrorists' interests are represented by several vectors where each vector relates to a different topic of interest. The cosine similarity measure is commonly used to estimate the similarity between an accessed Web page and a given set of terrorists' topics of interests.

C. Clustering Techniques

Cluster analysis is the process of partitioning data objects (records, documents, etc.) into meaningful groups or clusters so that objects within a cluster have similar characteristics but are dissimilar to objects in other clusters (Han and Kamber 2001). Clustering can be viewed as unsupervised classification of unlabelled patterns (observations, data items or feature vectors), since no pre-defined category labels are associated with the objects in the training set. Clustering results in a compact representation of large data sets (e.g. collections of visited Web pages) by a small number of cluster centroids. Applications of clustering include data mining, document retrieval, image segmentation, and pattern classification (Jain *et al.* 1999). Thus, clustering of Web documents viewed by Internet users can reveal collections of documents belonging to the same topic. As shown by Sequeira and Zaki (2002), clustering can also be used for anomaly detection: normality of a new object can be evaluated by its distance from the most similar cluster under the assumption that all clusters are based on 'normal' data only. In this study clustering of Web pages retrieved from terrorist-related sites is used to find collections of Web pages belonging to the same terrorists' topic of interest. For each collection a centroid is computed and represented by the vector space model.

III. CONTENT-BASED DETECTION OF TERROR-RELATED ACTIVITY

A. Detection Environment

This study suggests a new type of knowledge-based detection methodology that uses the content of Web pages browsed by terrorists and their supporters as an input to a detection process. In this study, refers only to the textual content of Web pages, excluding images, music, video clips, and other complex data types. It is assumed that terror-related content usually viewed by terrorists and their supporters can be used as training data for a learning process to obtain a 'Typical-Terrorist-Behavior'. This typical behavior will be used to detect further terrorists and their supporters. A 'Terrorist-Typical-Behavior' is defined as an access to information relevant to terrorists and their supporters. A general description of a system based on the suggested methodology is presented in Figure 1. Each user under surveillance is identified as a 'user's computer' having a unique IP address rather than by his or her name. In the case of a real-time alarm, the detected IP can be used to locate the computer and hopefully the suspected terrorist who may still be logged on to the same computer.

The suggested methodology has two modes of operation:

Learning typical terrorist behavior: in this mode, a collection of Web pages from terror-related sites is downloaded and represented as a set of vectors using the vector space model. The collected data is used to derive and represent the typical behavior ('profile') of the terrorists and their supporters by applying techniques of unsupervised clustering. Since the IP addresses of downloaded pages are ignored, moving the same or similar contents to a new address, as frequently carried out by terror-related sites, will not affect the detection accuracy of the new method.

Monitoring users: this mode is aimed at detecting terrorist users by comparing the content of information accessed by users to the Typical-Terrorist-Behavior. The textual content of the information accessed by a user on the Web is converted into a vector called 'access vector'. An alarm is issued when the similarity between the 'access vector' and the

Typical -Terrorist-Behavior is above a predefined threshold. The privacy of regular users is preserved, since the system does not need to store either the visited IP addresses or the actual content of the viewed pages. Due to extensive dimensionality reduction procedures, the access vectors do not hold sufficient information to restore the actual content of the page. Other types of viewed content (e.g., images) are ignored by the system.

Apparently the task of making a distinction between a legitimate user and a terrorist is a typical classification problem (see Han and Kamber 2001) with two alternative classes: terrorists vs. non-terrorists. However, the most popular classification algorithms (such as Naïve Bayes and C4.5) are probabilistic in their nature, i.e., they assume a stable and relatively balanced probability distribution of classes over the entire population. Moreover, they usually ignore any differences between the misclassification costs of objects belonging to different classes. All these assumptions are totally wrong when dealing with terrorist detection on the Web. The monitored population is completely unbalanced so that the actual percentage of terrorists in the entire population of Web users is usually close to zero. It is not expected that the amount of terrorist activities on the Web to be stable either – in fact, this study is interested in any fluctuations in these activities. And, of course, missing one real terrorist in a haystack of legitimate users may be more costly than suspecting several legitimate users of being active terrorists. Consequently, in this study it was decided to follow the more flexible clustering approach to terrorist detection, while leaving the investigation of classification methods to future research.

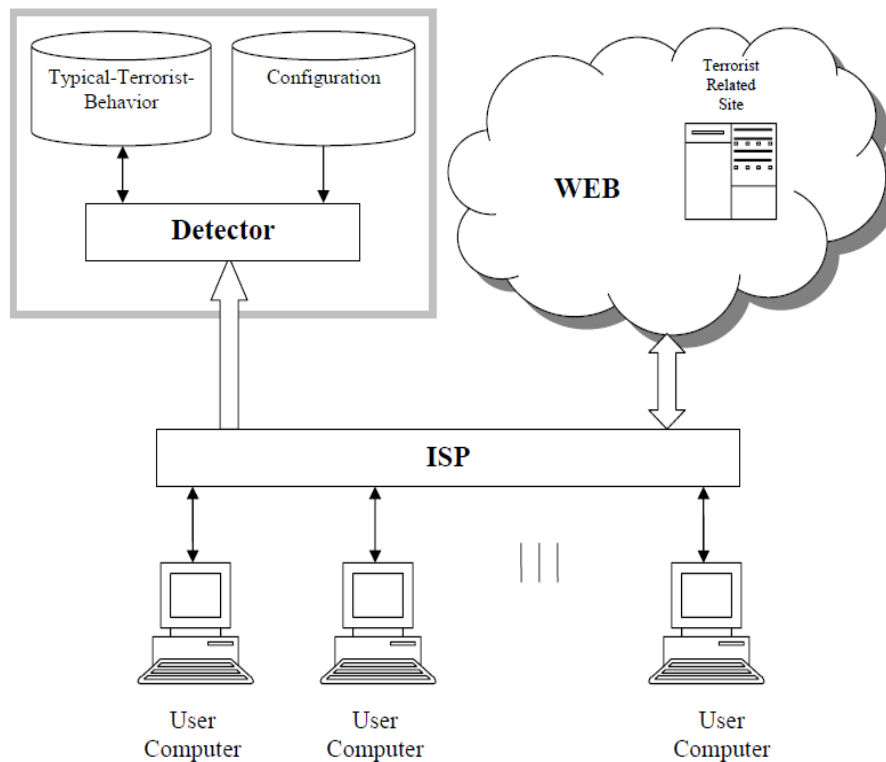


Figure 1: Detection environment

The following subsections describe in detail the Typical-Terrorist-Behavior learning module and the detection algorithm.

Learning the Typical-Terrorist-Behavior

The learning Typical-Terrorist-Behavior part of the methodology defines and represents the typical behavior of terrorist users based on the content of their Web activities. Figure 2 describes the learning module. It is assumed that it is possible to collect Web pages from terror-related sites. The content of the collected pages is the input to the Vector Generator module that converts the pages into vectors of weighted terms (each page is converted to one vector). The vectors are stored for future processing in the Vector of Terrorists Transactions DB.

The Clustering module accesses the collected vectors and performs unsupervised clustering resulting in n clusters representing the typical topics viewed by terrorist users. For each cluster, the Terrorist-Representor module computes the centroid vector (denoted by C_{vi}) which represents a topic typically accessed by terrorists.

As a result, a set of centroid vectors represent a set of terrorists interests referred to as the 'Typical-Terrorist-Behavior'.

The Typical -Terrorist- Behavior is based on a set of Web pages that were downloaded from terrorist-related sites and is the main input of the detection algorithm. In order to make the detection algorithm more accurate, the process of generating the Typical-Terrorist-Behavior has to be repeated periodically due to changes in the content of terrorist related site. Typical-Terrorist-Behavior depends on the number of clusters. When the number of clusters is higher, the Typical-Terrorist-Behavior includes more topics of interest by terrorists where each topic is based on fewer pages. It is hard to hypothesize what the optimal number of clusters is. In the case study presented in the next section detection performance for two settings of the number of clusters are presented.

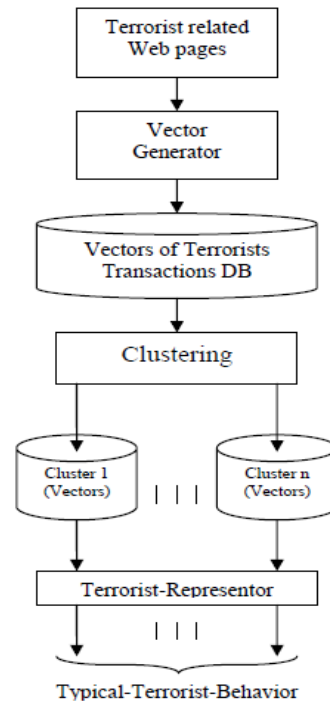


Figure 2: Learning the Typical-Terrorist- Behavior

Detecting Typical Terrorist Behavior

In the Monitoring module (Figure 3), the Vector-Generator converts the content of each page accessed by a user into a vector representation (referred to as the ‘access vector’). The Detector uses the access vector and the Typical-Terrorist-Behavior and tries to determine whether the access vector belongs to a terrorist. This is done by computing similarity between the access vector and all centroid vectors of the Typical-Terrorist-Behavior. The cosine measure is used (Boger et al. 2001; Pierrea et al. 2000) to compute the similarity.

The detector issues an alarm when the similarity between the access vector and the nearest centroid is higher than the predefined threshold denoted by tr:

$$Max \left(\frac{\sum_{i=1}^m (tCv_{i1} \cdot tAv_i)}{\sqrt{\sum_{i=1}^m tCv_{i1}^2 \cdot \sum_{i=1}^m tAv_i^2}}, \dots, \frac{\sum_{i=1}^m (tCv_{im} \cdot tAv_i)}{\sqrt{\sum_{i=1}^m tCv_{im}^2 \cdot \sum_{i=1}^m tAv_i^2}} \right) > tr$$

where Cvi is the i th centroid vector, Av - the access vector, $tCvi1$ - the i th term in the vector Cvi , $tAvi$ - the i th term in the vector Av , and m - the number of unique terms in each vector.

The threshold parameter tr controls the sensitivity of the detection. Higher value of tr will decrease the sensitivity of the detection process, decrease the number of alarms, increase the accuracy and decrease the number of false alarms. Lower value of tr will increase the detection process sensitivity, increase the number of alarms and false alarms and decrease the accuracy. The optimal value of tr depends on the preferences of the system user. In the next section, the feasibility of the new methodology is explored using a case study.

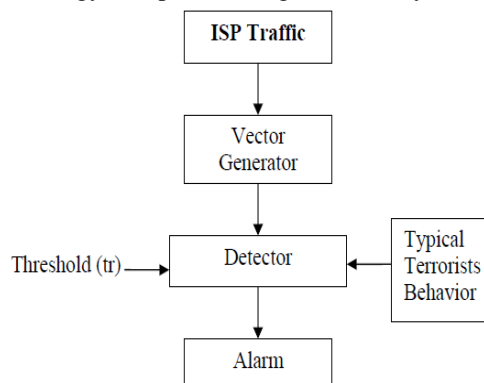


Figure 3: Detecting Terrorists – Monitoring Module

IV. CASE STUDY

Experimental Settings

An initial evaluation of the proposed knowledge-based detection methodology is conducted through a prototype system. The experimental environment included a small network of nine computers, each computer having a constant IP address and a proxy server through which all the computers accessed the Web. In the experiment, the proxy server represented an ISP.

Eight students in information systems engineering were instructed to access Web sites related to general topics and generated about 800 transactions. Several other users were requested to access terrorist related information (mainly the 'Azzam Publications' sites visited by one of the September 11 terrorists) and generated about 214 transactions. In this experiment, the users accessed only pages in English, though the methodology is readily applicable to other languages as well. The Vector Generator, Clustering and Detector modules described above were implemented and installed inside the server. Vcluster program from the Cluto Clustering Tool (Karypis 2002) is used to implement the clustering module where the clustering algorithm used was 'k-way' and the similarity between the objects computed using the cosine measure (Boger et al. 2001; Pierrea et al. 2000).

The 'k-way' clustering algorithm is a variation of the popular K-Means clustering algorithm. One problem with these algorithms is that it is hard to find the optimal k (number of clusters) that will achieve the best clustering performance. The experiments were done with different k's and compared results.

Evaluation Measures

To evaluate the system performance the following measures (based on Sequeira and Zaki 2002) were used:

True Positive Rate (TP) (also known as Detection Rate or Completeness): the percentage of terrorist pages receiving a rating above the threshold (referred to as tr in the model). In the experiments, terrorist pages will be obtained from the users simulating terrorists.

False Positive Rate (FP): the percentage of regular Internet access pages that the system incorrectly determined as related to terrorist activities, i.e., the percentage of non-terrorist pages receiving a rating above threshold and suspected falsely as terrorists.

Accuracy – percentage of alarms related to terrorist behavior out of the total number of alarms.

Since no benchmark data on content-based intrusion detection is currently available, the results are compared to the best numbers achieved with ADMIT which is a command-level method using the K-means clustering algorithm to detect intruders (Sequeira and Zaki 2002).

Summary of Results

As mentioned above, 800 vectors representing pages accessed by non-terrorist users and 214 vectors representing pages related to terrorist activities were used. The latter pages were collected from various terrorist related Web sites and used to train the system in the learning mode. The experiment included the following steps:

1. Set the initial value of the threshold parameter tr to zero and the initial number of clusters k to 9.
2. Randomly select 43 (out of 800) vectors from the non-terrorist set of vectors. These vectors are used to test the system's ability to ignore non-terrorist users. The same number of terror-related pages was used to evaluate the detection capability of the system (see the next step).
3. Randomly select 43 (~20%, out of 214) vectors from the terrorist set of vectors as a validation set. These vectors are used to test the system's ability to detect terrorist users.
4. Train the system (learning phase) using the remaining 171 vectors (~80%, out of 214) representing terrorist related pages. Apply the clustering algorithm to produce a set of clusters representing terrorist topics of interests and calculate the centroid for each topic (cluster). This step results in a set of k vectors forming the 'Typical Terrorist Behavior'.
5. Use the 43 vectors representing pages accessed by a terrorist user as an input to the detector and observe the percentage of terrorist related vectors that raised an alarm (True Positive Rate measure).
6. Use the 43 vectors representing pages accessed by regular users as an input to the detector and observe the percentage of vectors that mistakenly raised an alarm (False Positive Rate measure). Finally, calculate the percentage of alarms that are terrorist-related out of the total number of alarms (Accuracy measure).
7. Repeat steps 5-6 for different values of threshold tr between 0 and 1.
8. Repeat steps 3-7 five times (each time selecting a different set of terrorists vectors at step 3 for cross validation). The average results of these five runs are shown on the ROC graph (Provost and Fawcett 2001) in Figure 4.

9. Repeat the whole process for a different value of k (number of clusters) to evaluate the sensitivity of the system performance to this parameter.

The ROC (Receiver-Operator Characteristic) curves in Figure 4 describe the entire evaluation process for 5 and 9 clusters ($k=5, k=9$). The X-axis represents the FP (False Positive), and the Y-axis represents the TP (True Positive). Every point in the ROC curves represents an average result of five cross-validation runs where in each run different 43 terror-related vectors were selected for validation. The graph does not reveal a significant change in the performance as a result of reducing the number of clusters from 9 to 5.

The graph in Figure 5 describes the accuracy as a function of the threshold parameter. The results show that the accuracy is a monotonically increasing function of the threshold: for a sufficiently high value of Tr the FP rate approaches zero. However this graph cannot suggest the best threshold as its choice depends on user-specific preferences (such as the false alarm and the non-detection cost functions).

The initial case study clearly suggests that the proposed methodology is feasible and a system implementing this methodology might reliably detect terrorists accessing their sites on the Internet based on the content of monitored Web traffic. In the experiments, the prototype system reached TP=93% and FP=11.7% on average (see Table 1, Threshold = 0.2), compared to TP=70% and FP=15% obtained by the ADMIT system (Sequeira and Zaki 2002), which utilized user command-level data.

Table 1- results for 9 clusters

Average Roc			Accuracy
Threshold	Fp	Tp	
0	1	1	0.5
0.05	0.916	0.991	0.520
0.1	0.538	0.977	0.645
0.14	0.285	0.963	0.772
0.15	0.225	0.958	0.812
0.16	0.197	0.944	0.829
0.17	0.173	0.939	0.845
0.18	0.168	0.935	0.848
0.19	0.135	0.930	0.873
0.2	0.117	0.930	0.889
0.21	0.103	0.911	0.898
0.22	0.093	0.897	0.906
0.23	0.093	0.893	0.905
0.24	0.079	0.874	0.917
0.25	0.075	0.855	0.920
0.26	0.051	0.827	0.941
0.27	0.042	0.823	0.951
0.28	0.033	0.799	0.960
0.29	0.028	0.766	0.965
0.3	0.023	0.734	0.969

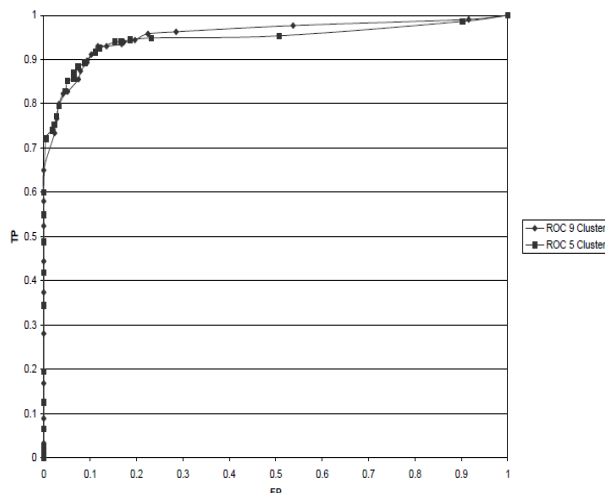


Figure 4: True Positive and False Positive Rate for 9 clusters

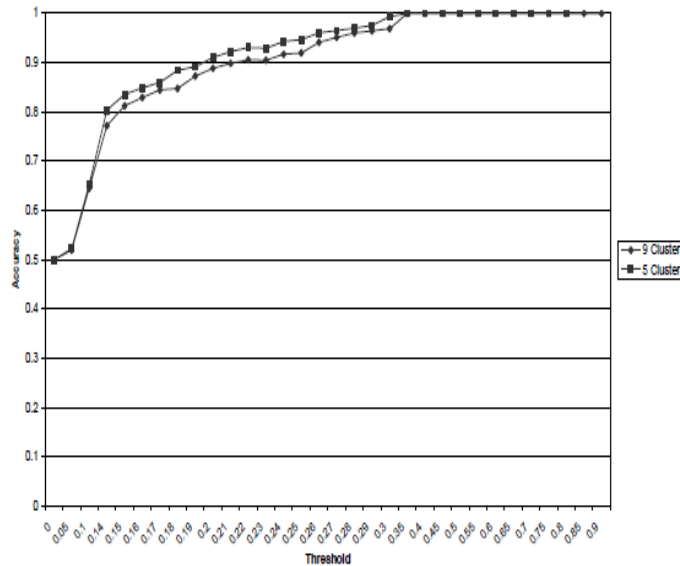


Figure 5: Accuracy as a Function of Threshold

V. DEPLOYMENT ISSUES

A system implementing the new methodology can be deployed by law enforcement agencies in two different ways where each way has its advantages and disadvantages.

ISP-based System: The system implementing the new methodology may be deployed within the ISP infrastructure. The major advantage of such a deployment is that the ISP is able to provide the exact identity of a suspicious user detected by the system since the IP is allocated to the user by the ISP. The disadvantages of such a deployment are that it requires ISP awareness and cooperation and the privacy of the ISP subscribers is violated.

Network-based System: The system implementing the new methodology is eavesdropping on communication lines connecting the ISPs to the Internet backbone. In such a deployment the major advantage is that the ISP cooperation is not required and the privacy of the ISP subscribers is protected, since most ISPs are allocating a temporary IP address for every user. The major disadvantage of such a deployment is that the exact identity of a subscriber using a given IP address is unknown.

It is difficult to effectively deploy a system based on the new methodology on sites that provide internet access to casual users such as an Internet Café or Hot Spots since users are not required to identify themselves to the service operator.

VI. CONCLUSION

In this paper, an innovative, knowledge-based methodology for terrorist activity detection on the Web is presented. The results of an initial case study suggest that the methodology can be useful for detecting terrorists and their supporters using a legitimate ways of Internet access to view terror-related content at a series of evasive web sites.

The ongoing research includes the following five main issues:

1. Document representation: It is planned to perform a comparison between the common vector-space model and a novel graph-based model (Schenker *et al.* 2003) that can capture relationships between terms in a textual document. This issue is important as it is believed that the success of the methodology depends on accuracy of content representation.
2. Similarity Measures: In the monitoring phase, the classification of content in every accessed page depends on the calculation of similarity between the access vector and the cluster centroids. There is a need to compare the results of the cosine measure to other measures such as Euclidean distance.
3. Detection methodology: In a real-world environment, more accurate results may be obtained by monitoring sequences of page views rather than raising an alarm after every suspect page. Developing an anomaly detection system for detecting abnormal content, which may be an indication of terrorist or other criminal activities, is another important research direction which was initiated (Last *et al.* 2001). Applying classification methods to the terrorist detection problem is another interesting direction.
4. Computational complexity: A system based on the new methodology has to process every HTML page that is being accessed by any subscriber of an ISP where it is deployed. There is a need to work on reducing the

computational complexity of the proposed methodology. One way to achieve this goal is to reduce the size of each access vector (dimensionality reduction) without significantly reducing the system detection performance.

5. Optimal settings: Further analysis is required to determine the system settings such as the number of clusters (k), and the detection threshold.
6. Analyzing picture and binaries: Detection may also be based on MD5 hashes of pictures or binaries that are downloadable from terror-related sites. This would again provide data for clustering algorithms. The important features may be extracted manually (e.g., by image processing techniques).

The detection methodology presented here can be applied to detecting other types of criminals surfing the Web such as pedophiles accessing child pornography sites.

ACKNOWLEDGMENT

I am using this opportunity to express my gratitude to Mr. Mozibur Raheman Khan M.Sc.,M.Phil., Assistant Professor, Department of Computer Science, Jamal Mohamed College, Trichy. for his support and guidance.

REFERENCES

- [1] Han, J., Kamber, M. (2001) *Data Mining: Concepts and Techniques*, Morgan Kaufmann.
- [2] Balasubramaniyan, J.S., Garcia-Fernandez, J.O., Isacoff, D., Spafford, E., Zamboni, D.(1998) An architecture for intrusion detection using autonomous agents, *Proceedings 14th Annual Computer Security Applications Conference*, IEEE Comput. Soc, Los Alamitos, CA, USA, xiii+365, pp. 13-24.
- [3] Boger, Z., Kuflik, T., Shoal, P., Shapira, B.(2001) Automatic keyword identification by artificial neural networks compared to manual identification by users of filtering systems, *Information Processing and Management*, **37**:187-198.
- [4] Corbin, J. (2002) *Al-Qaeda: In Search of the Terror Network that Threatens the World*, Thunder's Mouth Press / Nation Books, New York.
- [5] Debar, H., Dacier, H., Dacier, M., Wespi, A. (1999) Towards a taxonomy of intrusion-detection systems, *Computer Networks*, **31**, pp. 805–822.
- [6] Ingram, M. (2001) Internet privacy threatened following terrorist attacks on US, URL: <http://www.wsws.org/articles/2001/sep2001/isps24.shtml>
- [7] Jain, A.K., Murty, M.N., Flynn, P.J. (1999) Data Clustering: A Review, *ACM Computing Surveys*, **31**, 3:264-323.
- [8] Karypis, G. (2002) CLUTO - A Clustering Toolkit, Release 2.0, University of Minnesota, URL: <http://www.users.cs.umn.edu/~karypis/cluto/download.html>.
- [9] Kelley, J. (2002) Terror Groups behind Web encryption, USA Today, URL: http://www.apfn.org/apfn/WTC_why.htm
- [10] Last, M., Shapira, B., Elovici, Y. Zaafrany, O., Kandel, A. (2001) Content-Based Methodology for Anomaly Detection on the Web, submitted to AWIC'03, Atlantic Web Intelligence Conference, Madrid, Spain.
- [11] Lemos, R. (2002) What are the real risks of cyberterrorism?, *ZDNet*, URL: <http://zdnet.com.com/2100-1105-955293.html>.
- [12] Pierrea, S., Kacanb, C., Probstc, W. (2000) An agent-based approach for integrating user profile into a knowledge management process, *Knowledge-Based Systems*, **13**:307-314.
- [13] Provost, F., Fawcett, T. (2001). Robust Classification for Imprecise Environments. *Machine Learning***42**,3:203-231.
- [14] Richards, K. (1999) Network Based Intrusion Detection: A Review of Technologies, *Computers & Security*, **18**:671-682.
- [15] Salton, G. (1989) *Automatic Text Processing: the Transformation, Analysis, and Retrieval of Information by Computer*, Addison-Wesley, Reading.

BIOGRAPHY



currently doing M. Phil degree in JAMAL MOHAMED COLLEGE (Autonomous) Affiliated to BHARATHIDASAN UNIVERSITY, Trichy. Received M.Sc. (Information Technology) degree from ST. JOSEPH'S COLLEGE(Autonomous), Trichy-2. (Affiliated to BHARATHIDASAN UNIVERSITY).