# Novel Surveillance System for Instant Messengers using Data Mining Techniques

Punitha M, 4th sem, M.Tech,
Dept of Computer Science and Engineering,
SJBIT,
Bangalore,India.

Mrs. Bindiya M K , Asso. Prof.,
Dept of Computer Science and Engineering,
SJBIT,
Bangalore,India.

*Abstract*: **Instant messengers have now become a common means of communication in today's world. The suspicious messages are communicated through the instant messages which are untraceable via cyber crime. We implement a surveillance system that discovers and predict the suspicious messages that are communicated via instant messages. The instant messenger surveillance system uses data mining and ontology to identify suspicious messengers with the personal details like IP address ,email id etc. Framework is developed using Ontology based Information Extraction technique (OBIE), Association rule mining (ARM) a data mining technique with set of pre-defined Knowledge-based rules (logical), for decision making process that are learned from domain experts and past learning experiences of suspicious dataset like *GTD* (Global Terrorist Database). The experimental results obtained will aid to take prompt decision for eradicating cyber crimes.**

*Keywords: Instant Messenger(IM), Ontology based information extraction, Data Mining techniques.*

## I. INTRODUCTION

**Internet** evolutions led to the growth of innumerable cybercrimes[1]. Criminals have adapted to send suspicious messages via mobile, Instant Messengers applications and Social Networking Sites, which makes it difficult to trace their criminal activities dynamically. The E-crime department must now be improvised with the development of technology to find criminals of cybercrime activity. Many of the Instant Messaging Systems (IMS) developed restricted their limit for sending messages, video, image and audio conferencing through the applications. They are not well equipped to detect online suspicious messages.

To avoid this kind of activity we improve the existing IMS using data mining techniques of Associative rules. And also uses WordNet a lexical database contains huge information on Ontology based on information techniques. This frame work uses the Suspicious pattern Detection (SPD) algorithm initiates the steps to capture the instant messages that are communicated between the users and stored into the database for identifying suspicious messages and predicts the type of cyber threat activity. And also e-crime monitoring system program to trace the culprit details for E-crime department.



Fig 1. Popular instant messaging applications

**Instant messaging** is the real-time transmission of text, audio, image and video using the internet. Fig 1 shows the popular instant messaging application. Instant messaging is also known as the online chat which is communicated between two or more people. In email chat only the person online can communicate messages but in the instant messengers and social networking sites offline messages can be sent which is then sent to the person who was offline. There are many instant messenger and social networking sites applications available today for example, Facebook, Tweeter, and LinkedIn.

**Web chatting** is a system that allows users to communicate in real time using easily available online interfaces on web. It is a type of internet online chat differentiated by its simplicity of using and accessibility to users who do not wish to take the time to install and learn to use specialized chat software. In today's world chatting online has become the popular way to communicate or connect with people individually or a group chat. The online chatting is meant to be short in which it enables other people to respond as quickly as possible. Now is not only left with the short text message but we can send the image, video clips, audio clip and files also.

**Spam** is the messages that are irrelevant or unsolicited, sent over the Internet, mainly to large amount of people, for the purposes of advertising, phishing, spreading malware, etc. E-mail *spam*, also known as unsolicited bulk e-mail (UBE), junk mail, or unsolicited commercial e-mail (UCE) or electronic junk mail. There are mainly two main types of spam, and they have significant effects on Internet users. The first type of spam is Cancellable Usenet spam is a single message sent to 20 or more Usenet newsgroups. The second type of spam is

Email spam targets individual users with direct mail messages.

Any Internet-enabled application is a potential carrier for worms and other malware. Instant messaging is no exception. Currently, there are more than 30 worms spread from instant messaging networks and their clients. Threats in instant messaging, now days it has been a new attack vector in instant messaging. The threats would be sent with the attachment or the links so need to be careful with the instant messaging application.

In the instant messaging the clients are hijacked by reading the user's friend list of the contact with some messages sent with worm attached to it. Instant messaging is an upcoming threat as and a carrier for malware. More peoples are now using instant messaging for all purpose. Instant messaging networks provide the ability to transfers files. Concurrently, instant messengers can transfer worms with the files.

**Phishing** is a trick that makes use Internet users from visiting fake websites. Phishing websites are designed as they look like the popular website's login page. The page look like actual website, people would log in to the website, by accidentally giving access to criminals with user accounts details. The fake website will then take the details and then send spam messages to perpetuate the phishing websites and promote service's or promote product's. When a phished account is used by a spammer, more and more Wall posts are spam or links to phishing sites are sent to the user, and the cycle continues.

Phishing is a continual threat that keeps growing from day to day. The risk grows even more in social media such as Facebook, Twitter etc. Hackers commonly use the phishing sites to attack people using the social media sites in their workplace, homes, or public in order to take the personal information and security information that will affect the user and the company (if in a workplace environment). In internet phishing is used to show trust in the user since the user may not be able to tell that the site being visited or program being used is not actual website, and when this occurs is the hacker has the chance to access the personal information such as email passwords, security codes, and credit card numbers among other things.

This section gives an over view on cybercrimes performed in IMS and other threats related internet enabled activity. The remainder of this paper is organized as follows: In Section II, we reviewed recent research advances in identifying criminals from cyberspace. The section III the literature survey on the work implemented. The section IV the basic system overview that shows the operational phases and its implementation. The section V shows the comparison results with existing IMS. The section VI concludes the paper.

## II. PROBLEM STATEMENT

Nowadays, it's difficult to survive without IMS as users are addicted to. Trillions of messages are sent each day through emails and IMS. Popular IMS such as AOL, MSN, ICQ, Yahoo, Google Talk, Skype, Facebook, Twitter, and LinkedIn have changed the way of communication with friends, acquaintances, and business colleagues. Once limited to desktops, popular instant messaging systems are finding their way onto handheld devices and cell phones, allowing users to chat virtually from anywhere.

## III. LITERATURE SURVEY

Instant messaging (IM) is a type of online chat which offers real-time text transmission over the Internet[1]. Some IM applications can use push technology to provide real-time text, which transmits messages character by character, as they are composed. More advanced instant messaging can add file transfer, clickable hyperlinks, Voice over IP, or video chat. The instant messages in instant messenger and social networking would be monitored of its behavior, activities, or other changing information, usually of people for the purpose of influencing, managing, directing, or protecting them. The monitoring would be done with the help of data mining and ontology. Ontology means the nature of being, becoming, existence, or reality, as well as the basic categories of being and their relations.

Recently Daya.C.et al. 2010 studied Ontology-Based Information Extraction[OBIE] current approaches[2]. Information extraction basically aims to retrieve certain types of information from natural language text by processing them. OBIE is the sun-field of Information Extraction. The major drawback of the work is in all the before OBIE the type of source of natural language is from Wiki, documents or set email or a web page. The future work that can done on this work is (a) improving the effectiveness of the Information Extraction process, (b) integrating OBIE systems with the semantic web and (c) improving the use of ontologies.

To check the multilingual contents from emails and do a framework was the work of Mohd Mahmoodali et al. in 2013[3]. The work also says that the new communication trend for spammers and criminals is email via image embedded with multilanguage text. The work also provides anti-spam and suspicious filtering framework for multilingual content. The future work of proposed work is doing a new framework for SMS in mobile phones and in the social networking sites.

In 2012 the support for Ontology grew and the Suggested Upper Merged Ontology [SUMO] was introduced. SUMO is the only formal ontology that has been mapped to all of the WordNet lexicon[4]. The major features of SUMO are: (a) Mapping to all WordNet, (b) language generation templates which is a software

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**NCRTS-2015 Conference Proceedings**

component that is designed to combine one or more language template with a data model to produce one or more result document.

Long back David.C. et al. in 1996 proposed work on maintence of discovered association rules in large database which led to an incremental updating technique[5]. The work also uses Association Rule Mining [ARM] is basically used for OBIE model. The ARM rules are applied on SSWDB to find support and confidence. The paper is discovering an incremental updating technique for association rules in the large databases. The main drawbacks of the proposed work is the design is not efficient algorithm for mining different types of rules and patterns, the design of efficient algorithm to update, maintain and manage the rules discovered. The future work is to implement a fast update algorithm[FUP] and studies on finding multiple level or generalized associates rules in large transaction database.

The Words that are not considered as the main thing in the sentence is listed in 2012 is named as stop-word in webconfs website[6]. Most search engines do not consider extremely common words in order to save disk space or to speed up search results. These filtered words are known as 'stop words'. Example words are able, am, allow, back, can't, directly etc.

To help social networks to detect malicious web content Michael Robertson et al. in 2010 proposed a security framework[7]. The importance of provenance information as a means to trust and validate the authenticity of available data cannot be stressed enough in today's web-enabled world. The large amount of data now accessible due to the Internet explosion brings with it the related issue of determining how much of it is trustworthy. Proved information, such as who is responsible for the data or how the data came to be, assists in the process of verification of the authenticity of the data.

To avoid the email threatening Appava et al. in 2009 proposed an approach to intelligent Data Mining system[8]. An approach to designing very fast algorithm for approximate string matching in the dictionary is proposed. The algorithm uses multiple spelling errors that corresponding to the insertion, deletion, changes, and transposes operations are done on the character strings that are considered in the fault model. The design of the very fast approximate string matching algorithms is done vai a four-step reduction procedure is described in the work. The most effective step uses hashing techniques to avoid comparing the given word with words at large distances.

To support the phishing detection in instant messengers in 2012 M. Mahmoodali proposed work using Data Mining[9]. Our survey shows that the techniques used in data extraction from deep webs need to be improved to achieve the efficiency and accuracy of automatic wrapper. The further investigations on deep web data extraction indicates that the development of a lightweight ontological

technique using existing lexical database for English (WordNet) is able to check the similarity of data records and detect the correct data region with higher precision using the semantic properties of the data record. The advantages of the method are that it can extract three types of data records, namely, single-sectioned data records, multiple-sectioned data records, and loosely sectioned data records, and the data records will also provides options for aligning iterative and disjunctive data items.

In 2009 Sunitha et al. proposed a RDF Reification for relational wrapper[10]. The paper focus on the different kinds of symmetric key encryption techniques those exist in present world for securing the data communications. The RDF will also frames all the techniques related to variety of encryption like image encryptions, information encryptions, double encryptions. RDF also aims to explain the performance parameters that are used in encryption processes and analyzing on the security issues. RDF helps people to protect their sensitive information from thefts so considered as one of the best tool when it is transmitted via insecure communication channels.

## IV. SYSTEM OVERVIEW

In Existing System it is difficult to trace criminal activities dynamically in Mobile Messages, Instant Messengers and Social Networking Sites. In fast growing country peoples are finding many techniques for entertainment and message transferring purpose also its may be a sites or social networks, in this social networks having many more specialty like message chatting, video calling audio calling and they restricted their sending message count, chat count video calling timing also but they are not well equipped to detect online suspicious messages.

The existing system as the following disadvantages(1) Difficult to trace the correct position of the system.(2)Difficult to find code word of terrorist.(3)Non-Securable data system.(4)At the same time cyber crimes also grownup rapidly, But the social networks are don't have mechanism to restrict this kind of activity.(5)By using this advantages terrors are easily passing their messages throw the internet. (6)Threaten calls also sending from one person to another person's but we couldn't find like that persons.

In the implemented system, we surveyed various architectures of Mobile Phones, Instant messengers and Social Networking sites. These studies helped us to develop a new Framework. WordNet is used as features for classification of words from unstructured text. Similarly, WordNet Ontology based on information extraction technique is discussed. Our Contribution includes improving the existing IMS using data mining technique of Associative rules, Ontology based information retrieval technique (probabilistic models), which is guided with pre-defined Knowledge based rules and ARM. Early detection of suspicious messages from instant messaging systems (Mobile Phone, IM and SNS) is possible with our proposed

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**NCRTS-2015 Conference Proceedings**

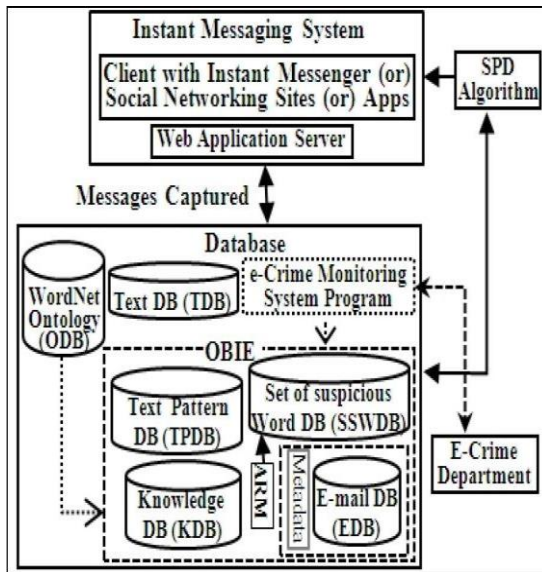Framework to identify and predict the type of cyber threat activity and trace the criminal details.



Fig 2. Implemented system architecture

The Fig 2 shows the operational phases that is implemented. The Suspicious Pattern Detection(SPD) algorithm initiates the steps to capture the instant messages that are communicated between the users and then stores them onto database for detection of suspicious messages. The main algorithm implemented is SPD algorithm which is apart from that, it also instigates the e-crime monitoring system program to trace the culprit details for E-crime department. The implementation makes use of database which stores dynamic messages and Ontology Based Information Extraction technique to detect suspicious words from messages with a Pre-defined Knowledge based rules. The three major tasks performed are as follows: (1) Word Extraction from Unstructured Text.(2) E-Crime Monitoring System Program.(3) SPD Algorithm.

In the word extraction from unstructured text module, filtering of unnecessary words from messages (TDB) is done. During this process, the suspicious words are identified using algorithms. The detected suspicious words are stored in TPDB for further processing. After finding suspicious words, the messages are marked as suspicious in SSWDB. KDB maintains the detected stem words along with the domain.

In the E-Crime Monitoring System Program module, the metadata is checked, for identifying from and to which Email-id the suspicious words belong and other relevant information. It will monitor the personal details of Email-id from EDB, which are provided during the creation of Email account. The suspicious messages that are sent through email id account using which computer details are tracked by E-Crime monitoring system program.
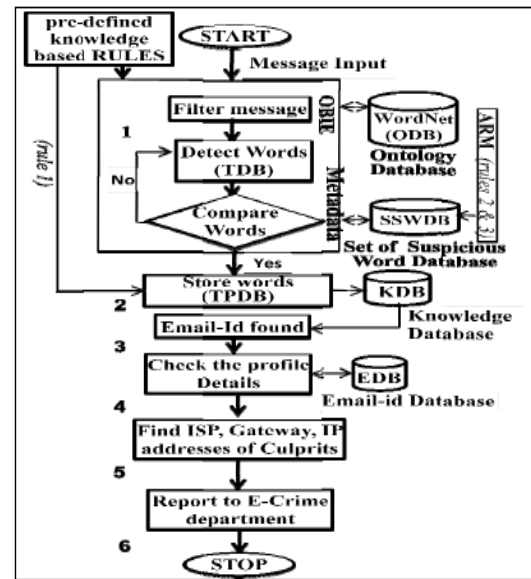


Fig 3. Schematic representation of SPD algorithm.

In the SPD Algorithm, if the detected words are suspicious, the e-crime monitoring system program is initiated by SPD (Suspicious Pattern Detection) algorithm. It is the backbone of our framework. It has initiated the overall progress starting from storing text messages in TDB till finding the culprits by providing a detailed report from KDB and EDB databases to E-Crime department when suspicious messages are found. The Fig 3 shows the SPD algorithm step by step.

## V. RESULTS

Comparing the implemented framework with the existing IMS. Currently none of the Instant Messengers has the ability to detect suspicious messages during online chat. The table below

Table 1. comparison of implemented framework with existing IMS.

| Features | IMS | Implemented Framework |
|---|---|---|
| Cyber threat activity detection | Static detection (time consumed) | Dynamic detection |
| Ontology support | No | Yes |
| Efficiency | Very Good | Moderate |
| Database & data Mining support | No | Yes |
| System architecture | Easy to Design | Complex to desgin |

The implemented framework will not only find the English simple words but can also find the suspicious words sent in short-form, code words. Example the words "kill, murder" can also be represented has "picnic". The messages sent would also be encrypted and only admin can see it in decrypted from.

## VI. CONCULSION

Framework aids to identify suspicious words from cyber messages and trace the suspected culprits. The currently

**Special Issue - 2015**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**NCRTS-2015 Conference Proceedings**

existing Instant Messengers lack of these features of capturing significant suspicious patterns of threat activity from dynamic messages and finds relationships among people, places and things during online chat as criminals have adapted to it. All the communication can be kept track and be safe. The Admin can also delete the suspicious users and can also obtain the IP address and E-Mail ID and other personal details that are being given during the registration.

## REFERENCES

[1] Mohammed Mahmood Ali, Khaja Moizuddin Mohammed, Lakshmi Rajamani, "framework for surveillance of instant Messages in IM and SNS using data mining and ontology", Proceeding of the IEEE Students Technology Sympusium,2014.

[2] Daya C. Wimalasuriya, and Dejing Dou,"Ontology-Based Information Extraction: An Introduction and a Survey of Current Approaches," Journal of Information Science,Volume 36, No. 3, pp. 306-323, 2010.

[3] M. Mahmood Ali, and L. Rajamani, "Framework for surveillance of instant messages, " published by inderscience in IJITST, vol. 5, 2013.

[4] (2012). [Online]. Available: http://www.ontologyportal.org/

[5] David W. Cheung, and et al., "Maintenance of discovered association rules in largedatabases: an incremental updating technique," published by IEEE in 1996.

[6] (2012). [Online]. Available: http://www.webconfs.com/stop-words.php.

[7] Michael Robertson, Yin Pan, and Bo Yuan, "A Social Approach to Security: Using Social Networks to Help Detect Malicious Web Content," published by IEEE in 2010.

[8] Appavu, and et al.,"Data mining based intelligent analysis of threatening e-mail," published by Elsevier in knowledge-based systems in 2009.

[9] M. Mahmood Ali, and L. Rajamani, "Phishing Detection in Instant Messengers using Data Mining Approach," proceedings of ObCom 2011, published by Springer-Verlag Berlin Heidelberg 2012, part I, CCIS 269, pp. 490–502, 2012.

[10] Sunitha Ramanujam, and et al., "A Relational Wrapper for RDF Reification," E. Ferrari et al. (Eds.): TM 2009, IFIP AICT 300, pp. 196–214, IFIP International Federation for Information Processing 2009.

[11] http://www.wiki.org