

# Caught in the Act of an Insider Attack: Detection and Assessment of Insider Threat

Philip A. Legg, Oliver Buckley, Michael Goldsmith and Sadie Creese  
Cyber Security Centre, University of Oxford, UK.

**Abstract**—The greatest asset that any organisation has are its people, but they may also be the greatest threat. Those who are within the organisation may have authorised access to vast amounts of sensitive company records that are essential for maintaining competitiveness and market position, and knowledge of information services and procedures that are crucial for daily operations. In many cases, those who have such access do indeed require it in order to conduct their expected workload. However, should an individual choose to act against the organisation, then with their privileged access and their extensive knowledge, they are well positioned to cause serious damage. Insider threat is becoming a serious and increasing concern for many organisations, with those who have fallen victim to such attacks suffering significant damages including financial and reputational. It is clear then, that there is a desperate need for more effective tools for detecting the presence of insider threats and analyzing the potential of threats before they escalate. We propose *Corporate Insider Threat Detection (CITD)*, an anomaly detection system that is the result of a multi-disciplinary research project that incorporates technical and behavioural activities to assess the threat posed by individuals. The system identifies user and role-based profiles, and measures how users deviate from their observed behaviours to assess the potential threat that a series of activities may pose. In this paper, we present an overview of the system and describe the concept of operations and practicalities of deploying the system. We show how the system can be utilised for unsupervised detection, and also how the human analyst can engage to provide an active learning feedback loop. By adopting an accept or reject scheme, the analyst is capable of refining the underlying detection model to better support their decision-making process and significant reduce the false positive rate.

## I. INTRODUCTION

The insider-threat problem is one that is constantly evolving and is having devastating impact on organisations worldwide. Those who operate within an organisation are often trusted with highly confidential information such as financial records and customer accounts, and often have detailed knowledge of operational procedures. Furthermore, the set of individuals who operate within the organisation is not always restricted to only employees, since contractors and suppliers may also have some level of access or knowledge of organisational procedures. Any individual who chooses to act maliciously has great potential to cause serious financial and reputational damage to the organisation. Media attention has highlighted numerous cases in recent years of both businesses and governments who have been compromised, where confidential information has

been exfiltrated and exposed. The threat posed by insiders is very real, and is one that requires serious attention by both organisations and individuals alike.

Technological advancements are constantly changing the way that organisations, and the people who act within the organisation, conduct their business. It has become common practice that employees now access documents from organisational file servers, communicate with both internal and external contacts via e-mail, and research information using the Internet. In addition, working practices have changed, so that employees may connect to organisation networks from home, or abroad, to provide flexibility in how we all choose to conduct the work-life balance. From an organisational point-of-view, what these technological advances also introduce is the capabilities of logging user activity, such as what resources have been accessed and at what time. Perhaps unsurprisingly though, is that these log files can become extremely large very quickly due to the sheer amount of activity that a typical organisation would undergo on a daily basis. In order to logging tools to prove effective, there is a need for automated assessment tools that can perform large-scale monitoring of all users, which can then reduce the information load that would be presented to an analyst for investigation.

In this work, we present our systematic approach for insider threat detection and analysis. From activity log data that is collected on users within the organisation, the system constructs a tree-structured profile for each user and for each role. This allows multiple activity types to be represented in a unified approach, that helps to support comparison between different users and associated roles. In addition, the profiles also provide an effective means to derive feature representations of the current behaviour exhibited by each user, which can then be compared against their previous observations and the observations made on their peers. We propose a number of anomaly metrics that are derived by decomposition of multiple groupings of related features. The system can be configured to flag up users who exceed a particular threshold on different anomaly metrics, and can also be configured to identify when users deviation from their peers on each of these metrics. We present a detection front-end that presents the anomaly metrics using a parallel co-ordinates plot, and incorporates an alert list for showing users who are flagged by the system. This allows the analyst to assess not only that a user poses a threat, but also why, as according to the detection scheme. The system also supports an active learning loop, where the analyst can accept or reject results in the alert list, which then feeds back into the detection model.

---

Philip A. Legg is now with the University of the West of England, UK, and Oliver Buckley is now with Cranfield University, UK.

## II. RELATED WORKS

The topic of insider threat, and how this can be detected, has received much attention in the literature. From our *Corporate Insider Threat Detection (CITD)* research project<sup>1</sup>, we have studied the many facets that relate to insider-threat detection, resulting in models to support insider-threat detection [1], a framework for assessing insider attacks based on insider-threat case studies [2], and also the extent of impact that insider-threat has on businesses [3]. Here, we focus primarily on the related works that target the development and implementation of systems that aim to identify insider threats.

Early work on insider-threat detection considered the use of honeypots [4], however as security awareness has increased, those choosing to commit insider attacks are finding more subtle methods to cause harm or defraud their organisations, and so there is a need for more sophisticated prevention and detection. Eldardiry *et al.* [5] propose an insider threat detection system that is based on feature extraction from user activities. However, they do not factor in role-based assessment. In addition, the profiling stage that we perform allows us to extract many more features beyond the activity counts that they suggest. Magklaras and Furnell [6] propose a threat evaluation system that estimates the level of threat that is likely to originate from a particular insider based on certain profiles of user behaviour. Buford *et al.* [7] use situation-aware multi-agent systems as part of a distributed architecture for insider threat detection. Brdiczka *et al.* [8] combine psychological profiling with structural anomaly detection to develop an architecture for insider-threat detection. Eberle *et al.* [9] consider Graph-Based Anomaly Detection as a tool for detecting insiders, based on modifications, insertions and deletions of activities from the graph. Myers *et al.* [10] consider how web server log data can be used to identify malicious insiders who look to exploit internal systems. Nguyen and Reiher [11] propose a detection tool for insider threat that monitors system call activity for unusual or suspicious behaviour. Maloof and Stephens [12] propose a detection tool for when insiders violate need-to-know restrictions that are in place within the organisation. Okolica *et al.* [13] use Probabilistic Latent Semantic Indexing with Users to determine employee interests, which are used to form social graphs that can highlight insiders. Liu *et al.* [14] propose a multilevel framework called SIDD (Sensitive Information Dissemination Detection) that incorporates network-level application identification, content signature generation and detection, and covert communication detection. Parveen *et al.* [15] use stream mining and graph mining to detect insider activity in large volumes of streaming data, based on ensemble-based methods, unsupervised learning and graph-based anomaly detection. Garfinkel *et al.* [16] propose tools for media forensics, as means to detecting insider threat behaviour. Compared to previous works, our focus is on detection tools that account for deviations in both user and role-based behaviour, which we demonstrate through the profiling stages and by how anomaly metrics are extracted. We also focus on the ability to mitigate false positives through active learning, and visualization tools for decision making.

<sup>1</sup><http://www.cs.ox.ac.uk/projects/CITD/>

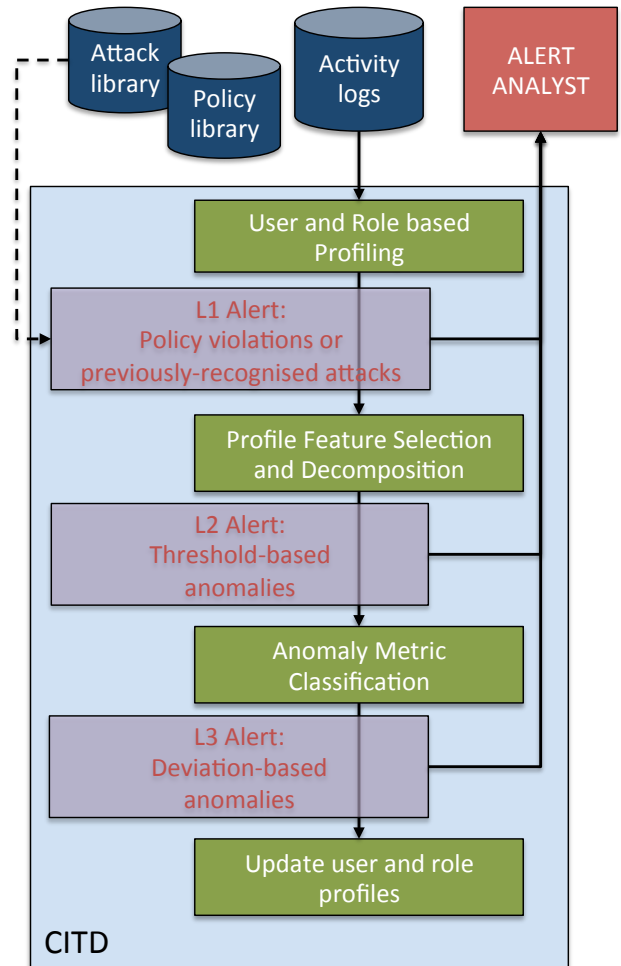


Fig. 1. An overview of the CITD detection tool. The system comprises of three alerting tiers, based on policy violations and previously-known attacks, threshold-based anomalies, and deviation-based anomalies. The system generates user and role based profiles from the observed activity logs. Feature selection is then performed, based on the profile content, to obtain scores for each of the anomaly metrics. Classification is then performed on the anomaly metrics, to determine whether the anomaly deviates significantly from the scope of their normal behaviour. If no alerts are triggered, then the user and role profiles are updated with the daily observations accordingly.

## III. CITD OVERVIEW

The CITD system is designed to detect the presence of insider threats within an organisation. An insider could be anyone who has some degree of access and knowledge regarding organisational resources and procedure, such as financial records, customer accounts, and sensitive documents. An insider threat is an abuse of this access or knowledge that has a detrimental impact on the organisation. Capelli *et al.* from the CERT division at Carnegie Mellon University identify three types of insider threat: IT sabotage, theft of IP, and data fraud [17]. However, given that all organisations will operate in different ways, how insiders choose to conduct these attacks will most likely vary across the different available

opportunities that they have in place. For this reason, there is great need for introducing technological solutions that can help to identify potential threats, based on deviations in their patterns of work compared to previous observations, or the observations of their peers.

Figure 1 provides a high-level overview of the CITD system, and how it can interface with existing infrastructure of the organisation. To facilitate integration of CITD, we assume that the system should be deployed on a standalone machine so as to not conflict or interfere with any existing infrastructure in the organisation. We also assume that the organisation adopts best-practice and maintains log files of user activity on their systems, however we do not make any assumptions about *how* organisations should capture activity logs for their organisation, as there are a number of different vendors and a multitude of approaches that an organisation may choose to adopt. The system can be configured to fetch all available logs, or all logs that correspond to a particular day (e.g., the current day), which are then streamed to the CITD system for processing. In this way, the system can be configured to run either on readily-available data, or can run as a continuous process that fetches new observations as these become available. Additional organisational resources such as policy libraries and attack libraries could also be interfaced with if such data is available, to alert against violations. The requirements of the CITD system are defined as follows:

- The system should be able to receive activity log data for all employees that operate within the organisation, for all systems that exist within the organisation that generate log records.
- The system should be able to construct a profile for each individual who operates within the organisation, and for each job role that a group of individuals act within, that deems what their normal behaviour is based on the observations of the different activities they perform.
- The system should be able to quantify the threat that is posed by each individual who operates within the organisation, based on the amount of deviation between their normal behaviour and their currently-observed behaviour.
- The system should be able to provide detail to the analyst via visual representation, that can inform on who is a potential threat, and what anomalies suggest this.
- The system should be able to take feedback from the analyst to enhance the detection routine, based on the acceptance or rejection of the generated results.

In order to test our approach, we have used synthetic datasets created by CMU/CERT<sup>2</sup>, as well as in-house datasets that have been developed separately from our detection work. The datasets were designed to reflect a realistic organisation, with not only suspicious activity of insider threats, but most importantly, to capture a realistic notion of normal behaviour. This aims to ensure that the detection tool can deal with sufficiently noise data that is a close representation of the activities that real users would conduct. The following sections describe the different stages of the detection process.

<sup>2</sup><https://www.cert.org/insider-threat/tools/>

### A. User and Role-based Profiling

The first component of the CITD system is concerned with user and role-based profiling. We construct the profiles based on a tree structure which provides a standardised, yet flexible, scheme for representing user activity which is later used as the basis for feature selection. As activity logs are streamed, the system will obtain the corresponding user profile based on the observed log, and either update or append the activity as appropriate. The root node defines the user (or role), from which there are three possible branches: daily observations, normal observations, and attack observations. For each of these, the subsequent levels of the tree define the device that captured the log, the activity observed on the device, and the primary attribute associated with the activity. For example, this could be that from PC-012 (device), an email was sent (activity), to john.davis@mycompany.com (attribute). For each node within the profile, a histogram is also maintained that represents the time of day that this observation was made. Our test data consists of five key observation types: logins, usb usage, e-mails, file access, and websites. Should additional logs be available (e.g., VPN, building access) then these could be easily incorporated into the profile also.

### B. Profile Feature Selection

Once we have computed the current daily profile for each user and for each role, we perform feature selection. Since the profile structure is well-defined, it means that comparisons between users, roles, or time steps can be easily made. In particular, the feature sets consists of three main categories: the user's daily observations, comparisons between the user's daily activity and their previous activity, and comparisons between the user's daily activity and the previous activity of their role. Below we describe the different features that we compute from the observed profile:

- 1) *New device for user* — has the user accessed a new device within the organisation.
- 2) *New device for role* — has the role accessed a new device within the organisation.
- 3) *New activity for device for user* — has the user performed a new activity within the organisation.
- 4) *New activity for device for role* — has the role performed a new activity within the organisational network.
- 5) *New attribute for activity for device for user* — has the user performed a new attribute within the organisation.
- 6) *New attribute for activity for device for role* — has the role performed a new attribute within the organisation.
- 7) *New activity for any device for user* — has the user performed a new activity within the organisation.
- 8) *New activity for any device for role* — has the role performed a new activity within the organisation.
- 9) *New attribute for any activity for any device for user* — has the user performed a new attribute within the organisation.
- 10) *New attribute for any activity for any device for role* — has the role performed a new attribute within the organisation.



Fig. 2. Overview of the detection system interface. The systems consists of multiple views, including an alert list (top left), a parallel co-ordinates view that shows how users score against each anomaly metric (top right), a configuration pane (bottom left) and a text output pane (bottom right).

- 11) *Hourly usage count for device* — a 24-bin histogram of the user’s usage.
- 12) *Hourly usage count for activity* — a 24-bin histogram of the user’s usage.
- 13) *Hourly usage count for attribute* — a 24-bin histogram of the user’s usage.
- 14) *Daily usage count for device*
- 15) *Daily usage count for activity*
- 16) *Daily usage count for attribute*

We maintain a matrix of size  $m \times n$ , where  $n$  represents the number of days being observed and  $m$  represents the number of features that are defined from the user profile. Each daily observation of features is then appended to the matrix as an additional row. If it is deemed appropriate to only consider a fixed number of previous days (e.g., 30 days prior to the current observation) then the matrix can be reduced accordingly.

### C. Anomaly Metrics and Classification

Given the feature matrix that we have described previously, the next stage of the system is to perform feature reduction to assess the amount of variance that is exhibited on particular characteristics. To do this, we perform a series of Principal Component Analysis (PCA) decompositions [18], based on the grouping of similar features. The system can support the

use of weighting functions that can be applied to individual features before the decomposition of features to anomaly metrics. Below we describe the different anomaly metrics that are currently supported in the prototype system:

- *Login\_anomaly*
- *Login\_duration\_anomaly*
- *Logoff\_anomaly*
- *USB\_inserstion\_anomaly*
- *USB\_duration\_anomaly*
- *USB\_removal\_anomaly*
- *Email\_anomaly*
- *Web\_anomaly*
- *File\_anomaly*
- *This\_anomaly* — where an anomaly has been observed on ‘this’ device.
- *Any\_anomaly* — where an anomaly has been observed on ‘any’ device.
- *New\_anomaly* — where an anomaly has been triggered by a ‘new’ observation.
- *Hourly\_anomaly* — where an anomaly has been triggered by a time-based observation.
- *Number\_anomaly* — where an anomaly has been triggered by a count-based observation.
- *User\_anomaly* — where an anomaly has been triggered by a ‘user’ comparison.

- *Role\_anomaly* — where an anomaly has been triggered by a ‘role’ comparison.
- *Total\_anomaly* — where an anomaly has been observed over all observed features.

In Figure 2, the detection interface shows the anomaly metric scores for each user by using a parallel co-ordinates plot. Acceptable behaviour would be expected to score low on each anomaly metric, since this would exhibit little deviation from the norm. Users who are deemed to be a threat are likely to score higher than usual on multiple anomaly metrics, suggesting that their currently observed behaviour deviates significantly from their normal. Of course, anomalous activity may not necessarily be threatening, however threatening activities are most likely to be anomalous for at least some of the detection metrics we consider. The parallel co-ordinates plot can also show the classification of user observations, based on the combined scores of the anomaly metrics. In this example, the detection system shows two axes for classification: standard deviation-based classification (*std*), and mahalanobis-based classification (*mahal*). The classification values are given based on the standard deviation and mahalanobis distance, as given by the set of anomaly metric scores for each individual. The greater that the standard deviation, or mahalanobis distance is, the greater the variation in the user’s observed behaviour. User’s who score over the specified deviation threshold are flagged up in the alert list. In this example, it can be seen that the user *lbegum1962* has been flagged up as a high alert for a number of days. The investigation reveals that they scored higher than usual on multiple anomaly metrics (*insert\_anomaly*, *file\_anomaly*, *this\_anomaly*, *new\_anomaly*, *user\_anomaly*, and *total\_anomaly*). As suggested by the alerts, the user was found to be accessing files on the server and using a USB storage device, both of which were new activities for this user to be performing.

By selecting the user from the alert list, the tree-structured profile can also be viewed to provide greater detail on the user behaviour (as shown in Figure 3). In this example, we can compare the histogram of time-observations for file activities (where normal behaviour is shown on the higher branch, and the ‘attack’ behaviour is shown on the lower branch). The histogram shows that there is a peak of activity much later on in the observation that is deemed as an attack, whereas the normal behaviour has a much more regular appearance. Likewise, the accessed files can be compared, that highlights a significant change between the ‘normal’ and ‘attack’ profiles.

#### D. Active learning

We have demonstrated how the system can be deployed to detect anomalous user behaviour that may be deemed threatening to the organisation. However, we also recognise that the limitation of many existing insider-threat detection systems lies in the false positive rate. To mitigate this, we also incorporate a semi-supervised learning approach, also known as *active learning* [19]. Active learning allows the analyst to intuitively incorporate knowledge back into the system that can improve how the underlying detection routine will perform in accordance to the desired outcome of the user.

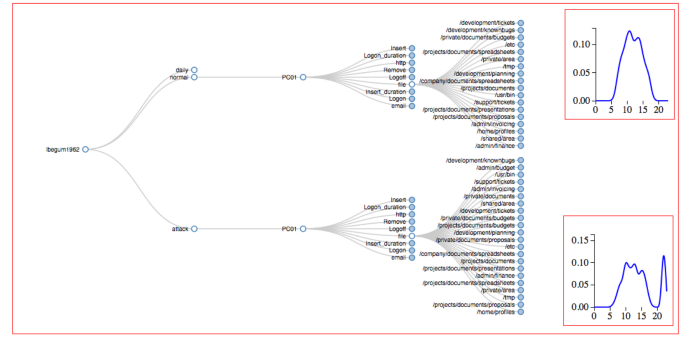


Fig. 3. Tree-structured profiles of user and role behaviours. The profile shows all the devices, activities, and attributes that the user has been observed performing. The probability distribution for normal hourly usage is given in the top-right, and the distribution for the detected attack is given in the bottom-right. Here it can be seen that the user has accessed resources late at night.

In Figure 1, although the insider *lbegum1962* has been correctly identified, it can be seen that a number of false positives are present also. The green and red circles next to each alert list entry correspond to an accept or reject policy, whereby the analyst can choose to accept or reject the result presented by the system based on their findings regarding the true state of the individual. If the analyst chooses to supply this information, then the weights associated with the detection metrics are reconfigured in accordance to the metrics that caused the individual to be flagged, and whether the analyst stated that the result is a reject or an accept. The weights for each anomaly metric are shown by the circular dials by each axis in the parallel co-ordinates plot. Figure 4 shows that the result for *mpowell1969* has been rejected, indicated by the removal of the accept option. This alert corresponds with the *insert\_anomaly*, and so the weighting for this metric has been reduced. On resuming the detection process, the system now only identifies the true insider, *lbegum1962*, dramatically reducing the presence of false positive results.

#### IV. CONCLUSION

In this paper we have presented our approach to insider threat detection, known as CITD. The system has been developed to analyse large-scale data repositories and activity logs to assess the current profile of all individuals who have access to the organisational systems. By incorporating user and role based profiling, the system is capable of obtaining a comprehensive feature set that characterises the user’s recent activity within the organisation. The feature set provides comparative assessment between multiple observations at previous time steps, and between multiple users. We compare a wide range of different metrics, to assess the degree of anomaly that is exhibited across each of these. Notifications are generated for the analyst based on different classification schemes of the anomaly metrics, including both threshold and deviation-based assessments. The system supports user intervention, to accept or reject results that are generated from the current detection model. This serves as a feedback loop to reconfigure the weightings associated with different anomaly metrics, based



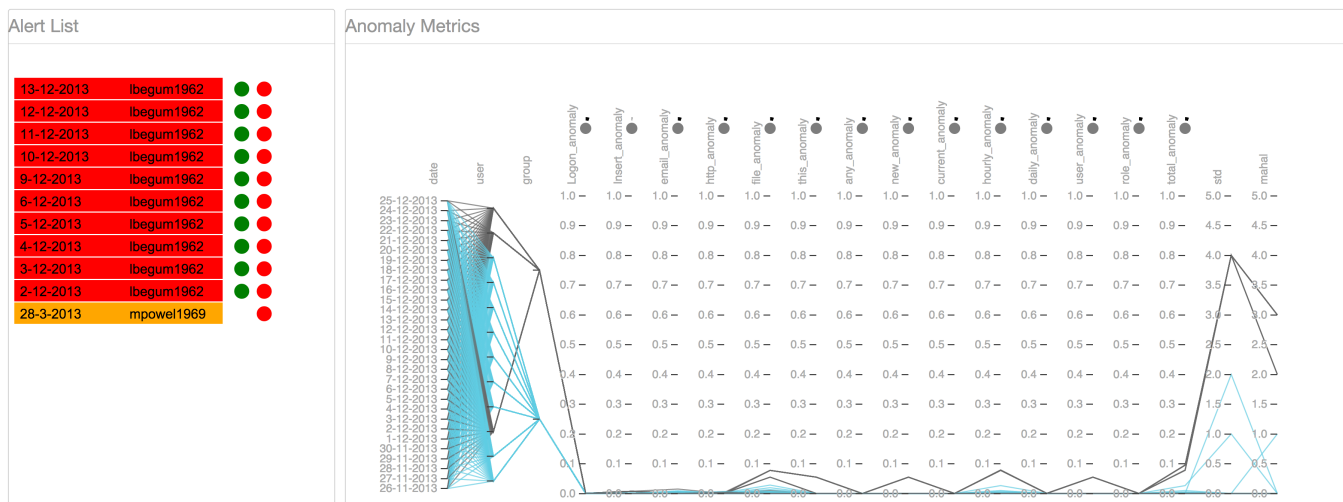


Fig. 4. Detection system as a result of active learning. The analyst has rejected the alert on mpowel1969 (shown by the removal of the accept option). This reconfigures the detection system to downgrade the anomaly associated with this result - in this case *insert\_anomaly* - which can be observed by the circular dials by each anomaly metric. In addition to the alert list, the parallel co-ordinates can be set to present only the 'last 30 days', which provides a clear view of the detected insider *ibegum1962*.

on the desired outcomes of the analyst. We have experimented with a number of synthetic datasets, including third-party examples, where the results obtained through unsupervised detection using the proposed anomaly metrics have been highly encouraging. We are currently in the process of deploying the detection system into a large multi-national corporation, to assess how well the system can perform based on real human activity, and real insider threats.

## REFERENCES

- [1] P. A. Legg, N. Moffat, J. R. C. Nurse, J. Happa, I. Agrafiotis, M. Goldsmith, and S. Creese. Towards a conceptual model and reasoning structure for insider threat detection. *Journal of Wireless Mobile Networks, Ubiquitous Computing and Dependable Applications*, 4(4):20–37, 2013.
- [2] J. R. C. Nurse, O. Buckley, P. A. Legg, M. Goldsmith, S. Creese, G. R. T. Wright, and M. Whitty. Understanding insider threat: A framework for characterising attacks. In *Proc. of the IEEE Symposium on Security and Privacy Workshops (SPW'14)*, pages 214–228, May 2014.
- [3] D. Upton and S. Creese. The danger from within. *Harvard Business Review*, September 2014.
- [4] L. Spitzner. Honeypots: catching the insider threat. In *Proc. of the 19th IEEE Computer Security Applications Conference (ACSAC'03)*, pages 170–179, December 2003.
- [5] H. Eldardiry, E. Bart, J. Liu, J. Hanley, B. Price, and O. Brdiczka. Multi-domain information fusion for insider threat detection. In *Proc. of the IEEE Symposium on Security and Privacy Workshops (SPW'13)*, pages 45–51, May 2013.
- [6] G. B. Magklaras and S. M. Furnell. Insider threat prediction tool: Evaluating the probability of IT misuse. *Computers and Security*, 21(1):62–73, 2002.
- [7] J. F. Buford, L. Lewis, and G. Jakobson. Insider threat detection using situation-aware mas. In *Proc. of the 11th International Conference on Information Fusion*, pages 1–8, 2008.
- [8] O. Brdiczka, J. Liu, B. Price, J. Shen, A. Patil, R. Chow, E. Bart, and N. Ducheneaut. Proactive insider threat detection through graph learning and psychological context. In *Proc. of the IEEE Symposium on Security and Privacy Workshops (SPW'12)*, pages 142–149, May 2012.
- [9] W. Eberle, J. Graves, and L. Holder. Insider threat detection using a graph-based approach. *Journal of Applied Security Research*, 6(1):32–81, 2010.
- [10] J. Myers, M. R. Grimaila, and R. F. Mills. Towards insider threat detection using web server logs. In *Proceedings of the 5th Annual Workshop on Cyber Security and Information Intelligence Research: Cyber Security and Information Intelligence Challenges and Strategies, CSIRW '09*, pages 54:1–54:4, 2009.
- [11] N. Nguyen and P. Reiher. Detecting insider threats by monitoring system call activity. In *Proceedings of the 2003 IEEE Workshop on Information Assurance*, 2003.
- [12] M. A. Maloof and G. D. Stephens. ELICIT: A system for detecting insiders who violate need-to-know. In *Recent Advances in Intrusion Detection*, volume 4637 of *Lecture Notes in Computer Science*, pages 146–166, 2007.
- [13] J. S. Okolica, G. L. Peterson, and R. F. Mills. Using plsi-u to detect insider threats by datamining e-mail. *International Journal of Security and Networks*, 3(2):114–121, 2008.
- [14] Y. Liu, C. Corbett, K. Chiang, R. Archibald, B. Mukherjee, and D. Ghosal. SIDD: A framework for detecting sensitive data exfiltration by an insider attack. In *Proc. of the 42nd Hawaii International Conference on System Sciences (HICSS '09)*, pages 1–10, January 2009.
- [15] P. Parveen, J. Evans, Bhavani Thuraisingham, K.W. Hamlen, and L. Khan. Insider threat detection using stream mining and graph mining. In *Proc. of the 3rd IEEE Conference on Privacy, Security, Risk and Trust (PASSAT)*, pages 1102–1110, October 2011.
- [16] S. L. Garfinkel, N. Beebe, L. Liu, and M. Maasberg. Detecting threatening insiders with lightweight media forensics. In *Technologies for Homeland Security (HST), 2013 IEEE International Conference on*, pages 86–92, Nov 2013.
- [17] D. M. Cappelli, A. P. Moore, and R. F. Trzeciak. *The CERT Guide to Insider Threats: How to Prevent, Detect, and Respond to Information Technology Crimes*. Addison-Wesley Professional, 1st edition, 2012.
- [18] I. Jolliffe. *Principal component analysis*. Wiley Online Library, 2005.
- [19] P. A. Legg, D. H. S. Chung, M. L. Parry, R. Bown, M. W. Jones, I. W. Griffiths, and M. Chen. Transformation of an uncertain video search pipeline to a sketch-based visual analytics loop. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2109–2118, Dec 2013.