

Intrusion Detection Systems & Multisensor Data Fusion: Creating Cyberspace Situational Awareness

Tim Bass

I. INTRODUCTION

Next generation cyberspace intrusion detection (ID) systems will require the fusion of data from myriad heterogeneous distributed network sensors to effectively create *cyberspace situational awareness*. This article provides a functional overview of how the art and science of multisensor data fusion can enhance the performance and reliability of ID systems. The article also discusses the data fusion inference process and data mining operations, outlines design challenges, and suggests areas for further research and development.

The vast majority of security professionals would agree that real-time intrusion detection systems (IDS) are not technically advanced enough to detect sophisticated cyberattacks by trained professionals. For example, during *The Langley Cyber Attack* the IDS failed to detect substantial volumes of e-mail bombs which crashed critical e-mail servers. Coordinated efforts from various international locations were observed as hackers worked to understand the rules-based filter used in counterinformation operations against massive e-mail bomb attacks [1].

At the other end of the technical spectrum, false alarms from ID systems are problematic, persistent, and preponderant. Numerous systems administrators have been the “subject” of an ID system reporting normal work activities as hostile actions. These types of false alarms result in financial losses to organizations when technical resources are denied access to computer systems or security resources are misdirected to investigate non-intrusion events. In addition, when systems are prone to false alarms, user confidence is marginalized and systems are poorly maintained and under utilized.

ID systems which examine operating system audit trails or network traffic [2] [3] and other similar detection systems have not matured to a level where sophisticated attacks are reliably detected, verified, and assessed. Comprehensive and reliable systems are complex and the technological designs of these advanced systems are only beginning to emerge. There remains much work to be done by IDS engineers in the design, integration, and deployment of efficient, robust, and reliable intrusion detection systems capable of reliably identifying and tracking hostile objects in cyberspace.

Recent industry studies forecast that the consumer market for security assessment tools will grow from approximately \$150M dollars per year in 1999 to over \$600M dollars in the year 2002. In addition, the author recently participated in a Department of Energy workshop which brought together security experts to help the federal government

prioritize a proposed \$500M dollar expenditure for research and development in the area of malicious code, anomalous activity and intrusion detection in the year 2000. Clearly, there are significant technical challenges ahead and a rapidly growing cyberspace intrusion detection marketplace.

The underlying issues and challenges are not unique to intrusion detection systems. Network management is also an expensive infrastructure to operate and more-often-than-not these systems fail to provide network engineers tangible and useful situational information, typically overwhelming operators with system messages and other low-level data. Network management and intrusion detection systems must operate in a uniform and cooperative model, fusing data into information and knowledge, so network operators can make informed decisions about the health and real-time security of their “corner” of cyberspace.

Multisensor data fusion provides an important functional framework for building next generation intrusion detection systems and *cyberspace situational awareness*. There exist significant opportunities and numerous technical challenges for the commercial application of data fusion theory into the art and science of cyberspace intrusion detection. Sections II through IV briefly review intrusion detection concepts and terms, provide an overview of the art and science of multisensor data fusion technology, and introduce the IDS data mining environment as a complementary process to the IDS data fusion model. Future design challenges and areas of further research to develop multisensor data fusion based ID systems are suggested in Section V.

II. INTRUSION DETECTION SYSTEMS OVERVIEW

Defensive information operations and computer intrusion detection systems are primarily designed to protect the *availability*, *confidentiality* and *integrity* of critical information infrastructures. These operations protect information infrastructures against denial of service (DoS) attacks, unauthorized disclosure of information, and the modification or destruction of data. The automated detection and immediate reporting of these events are required to respond to information attacks against networks and computers. In a nutshell, the basic approaches to intrusion detection today may be summarized as *known pattern templates*, *threatening behavior templates*, *traffic analysis*, *statistical-anomaly detection* and *state-based detection*.

Computer intrusion detection systems were introduced in the mid-1980’s to compliment conventional approaches to computer security. Technical writers on IDS often cite Denning’s [2] 1987 seminal intrusion detection model which is built on host-based subject profiles, systems objects, audit logs, anomaly records and activity rules. The underlying

ID model is a rules-based pattern matching system where audits are matched against subject profiles to detect computer misuse based on logins, program executions, and file access.

The subject-anomaly model was applied in the design of many host-based intrusion detection systems, i.e. *Intrusion Detection Expert System (IDES)* [4], *Network Intrusion Detection Expert System (NDIX)* [5] and *Wisdom & Sense (W&S)*, *Haystack*, and *Network Anomaly Detection and Intrusion Reporter (NADIR)* [6]. There are other ID systems based on the Denning model and an excellent survey of these systems may be found in [3]. The basic detection algorithms used in these systems include:

- weighted functions to detect deviations from normal usage,
- covariance-matrix based approaches for normal usage profiling,
- rules-based expert systems approach to detect security events.

The second leading technical approach to present-day intrusion detection is *multi-host network-based*. Heberlein *et al.* extended the Denning model to traffic-analysis on ethernet based networks with the *Network Security Monitor (NSM)* framework [7]. This was further extended with the *Distributed Intrusion Detection System (DIDS)* which combined host-based intrusion detection with network traffic monitoring [3] [8]. Current commercial IDS such as *Real Secure* and *Computer Misuse Detection System (CMDS)* have distributed architectures using either rules-based detection, statistical-anomaly detection, or both.

A significant challenge remains for IDS designers to combine data and information from numerous heterogeneous distributed agents (and managers) into a coherent process which can be used to evaluate the security of cyberspace. First, we review the basic concepts of the art and science of *multisensor data fusion*. This technology is an important avenue on the road toward the development of highly reliable intrusion detection and security-decision systems which identify, track, and assess cyberspace situations with multiple complex threats.

III. IDS DATA FUSION

Multisensor data fusion, or distributed sensing, is a relatively new engineering discipline used to combine data from multiple and diverse sensors and sources in order to make inferences about events, activities, and situations. These systems are often compared to the human cognitive process where the brain fuses sensory information from the various sensory organs, evaluates situations, makes decisions, and directs action.

Data fusion technology has been applied most prominently to military applications such as battlefield surveillance and tactical situation assessment. Data fusion has also emerged in commercial applications such as robotics, manufacturing, medical diagnosis, and remote sensing [10].

The application of data fusion in technical systems requires mathematical and heuristic techniques from fields

such as statistics, artificial intelligence, operations research, digital signal processing, pattern recognition, cognitive psychology, information theory and decision theory [10]. The functional application of multisensor data fusion to the art of intrusion detection is grounded in mathematical theory which is beyond the scope of this article. The interested reader is referred to [9] [10] and [11] for a detailed mathematical discussion.

Input into a data fusion cyberspace ID systems consists of sensor data, commands and *a priori* data from established databases. For example, the system input would be data from numerous distributed packet sniffers, system log-files, SNMP traps and queries, user profile databases, system messages, and operator commands. The output of data fusion cyberspace ID systems would be estimates of the identity (and possibly the location) of an intruder, the intruders activity, the observed threats, the attack rates, and an assessment of the severity of the cyberattack.

In a typical military command and control (C2) system, data fusion sensors are used to observe electromagnetic radiation, acoustic and thermal energy, nuclear particles, infrared radiation, noise and other signals. In cyberspace ID systems the sensors are different because the environmental dimension is different. Instead of a missile launch and supersonic transport through the atmosphere, cyberspace sensors observe information flowing in networks. However, just as C2 commanders are interested in the origin, velocity, threat, and targets of a warhead; network security personnel are interested in the identity, rate of attack, threat, and target of malicious intruders and criminals.

Waltz [9] described the generic sensor characteristics of a multisensor fusion system. These generic characteristics can be applied to next generation cyberspace IDS. We introduce these characteristics based on the generic Waltz model:

Detection Performance is the detection characteristics, i.e. false alarm rate, detection probabilities and ranges, for an intrusion characteristic against a given cyber noise background.

Spatial/Temporal Resolution is the ability to distinguish between two or more cyber intrusions in space or time.

Spatial Coverage is the span of the coverage or field of view for the sensor, (i.e. a the spatial coverage of a network sniffer might be the LAN segment it is monitoring.)

Detection/Tracking Modes is the mode of operation of the sensor, i.e. staring or scanning; single or multiple cyber target tracking, or capable of multimode operation.

Target Revisit Rate is the rate at which a cyber target or intrusion is revisited by the sensor to perform measurements.

Measurement Accuracy is the statistical probability that the cyberspace measurement or observation is accurate.

Measurement Dimensionality is the number or measurement variables between cyber target categories.

Hard vs. Soft Data Reporting is the status of the sensor reports, i.e. can a decision be made without correlation, or does the sensor require confirmation.

Detection/Tracking Reporting is the characteristic of the sensor to report individual cyber events or does the sensor maintain a time-sequence of the events or events.

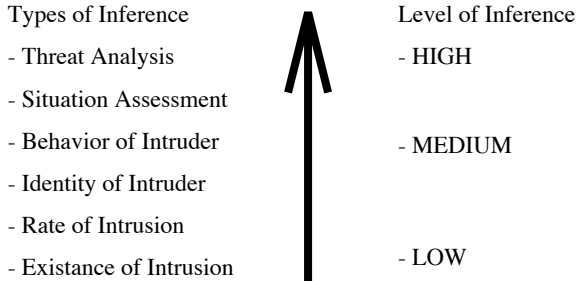


Fig. 1. Hierarchy of IDS Data Fusion Inferences

Real-time human decision making processes are supported by information derived from the fusion process. At the lowest level of inference, a data fusion cyber ID system would indicate the presence of an intruder or an attack. At the highest level the inference could be an analysis of the threat and the vulnerability. Figure 1 illustrates the hierarchy of ID data fusion inferences for a cyber threat.

Decision support systems for situational awareness are tightly coupled with data fusion systems. The basic decision system *observe - orient - decide and act, OODA*, is the classic decision support mechanism used in military information operations. OODA provides a cognitive mapping of the lowest level of cyber inference to knowledge based personnel actions. This cyber-fusion process requires the utilization of techniques ranging from processing algorithms and statistical estimations, to heuristic methods such as template correlation, or expert systems to assess situations and threats in cyberspace.

The IDS *observe* functions include the technical and human collection of data, including intrusion detection sensors, network sniffers, and computer system log files. The *orient* function includes data mining concepts to discover or learn previous unknown characteristics in the recorded data and computer files. The *orient* function also encompasses the application of templates for intrusion detection and association in data fusion processes. In the *decision* function, cyber information is further refined into threat knowledge which is used in the determination of an appropriate action or countermeasures. *Act* functions include both automated and human responses. Simple responses to cyberattacks may be automated; however, more complex decisions will always require human intervention.

The OODA decision-support process may be mapped into the three levels of abstractions. *Data* is the measurements and observations. *Information* is the data placed in context, indexed, and organized. *Knowledge* or intelligence is information explained and understood. These abstrac-

tions make up the intrusion detection data fusion model, illustrated in Figure 2, introduced by Waltz [12] for physical targets.

Cyberspace situational data is collected from sniffers and other intrusion detection sensors with primitive observation identifiers, times of observation, and descriptions. This raw data will require calibration or filtering and is commonly referred to as *Level 0 Refinement* in fusion models. All of these measurements must be aligned to a common frame of reference. This alignment is referred to as *Level 1 Object Refinement*.

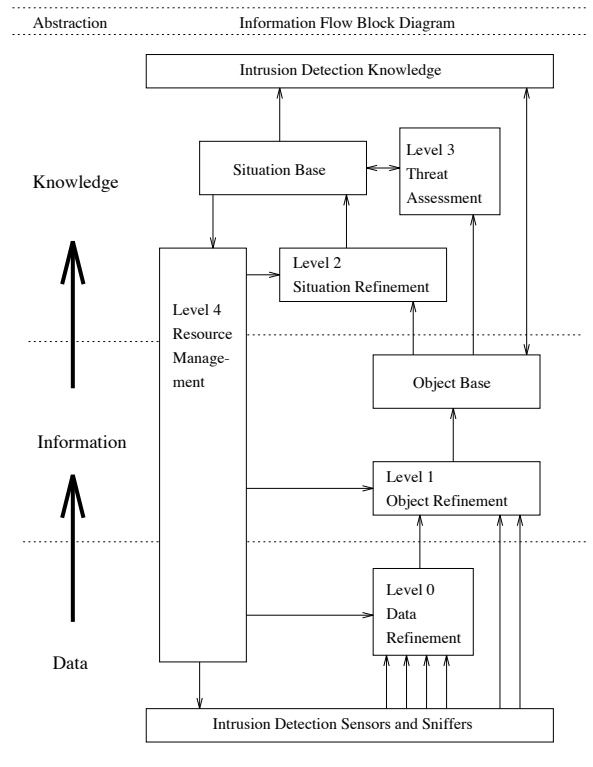


Fig. 2. Intrusion Detection Data Fusion

In Level 1 Object Refinement, data is correlated in time (and space if required) and data is assigned weighted metrics based on the relative importance. Observations may be associated and paired in this step of the process and classified according to intrusion detection primitives.

After objects have been aligned, correlated and placed in context in an information base, aggregated sets of objects are then detected by their coordinated behavior, dependencies, common points of origin, common protocols, common targets, correlated attack rates, or other high-level attribute. This step, *Situation Refinement*, provides situational knowledge and awareness.

Situation knowledge of cyberspace is used to analyze objects and aggregated groups against existing intrusion detection templates to provide an assessment of the current situation and suggest or identify future threatening attacks or cyberspace activity. Correlation between the *Level 3 Threat Assessment* and the security policy and objectives

determine the implications of the current *situation base*. The entire process is refined via *Level 4 Resource Management* based on the current situational awareness (and additional data as required) to further refine detection. For example, certain objects and subjects of interest may receive a higher processing priority, forming an intrusion detection - data fusion feedback loop.

The ID model described above is a deductive process used to detect previously known patterns in many sources of data by searching for specific intrusion signatures and templates in data streams to understand the state of the network security. As networks continue to evolve in complexity, the number of objects, situations, threats, sensors and data streams dramatically increase, presenting a very complex challenge for advanced IDS designers. Some of the potential applications and challenges are suggested later in this article.

IV. IDS DATA MINING

Intrusion detection cyberspace *data mining* is an off-line knowledge creating process where large sets of previously collected data is filtered, transformed, and organized into information sets. This information is used to discover hidden but previously undetected intrusion patterns.

Data mining is often called *knowledge discovery* and is distinguished from the data fusion process by two important characteristics, *inference method* and *temporal perspective* [12]. Data fusion uses known intrusion detection templates and pattern recognition. Data mining processes search for hidden patterns based on previously undetected intrusions to help develop new detection templates. In addition, data fusion focuses on the current state of the network based on past data; data mining focuses on new or hidden patterns in old data to create previously unknown knowledge, illustrated in Figure 3.

Raw data from relevant network management and intrusion detection systems are collected and indexed in the data warehouse. A major technical issue is how to reconcile the raw data from many different formats and inconsistent data definitions. This process is a part of the *data cleansing* operation. Data cleaning performs checks to insure that collected data is in correct ranges and limits, evaluates the overall consistency of the data, and insures that all indexed and referenced data and hierarchical relationships exist.

The initial data sets that will be used in data mining operation are selected in the *data selection and transformation* process. Data mining is normally performed on a small set and then extended to larger sets as patterns emerge and are validated. The *data mining operation* is performed on the selected data sets in either manual or automated modes. Waltz summarizes these operations in [12] for the physical realm:

Clustering is when data is segmented into subsets that share common properties.

Association is the analysis of both the cause-and-effect and the structure of relationships between data sets.

Statistical Analysis is performed to determine the like-

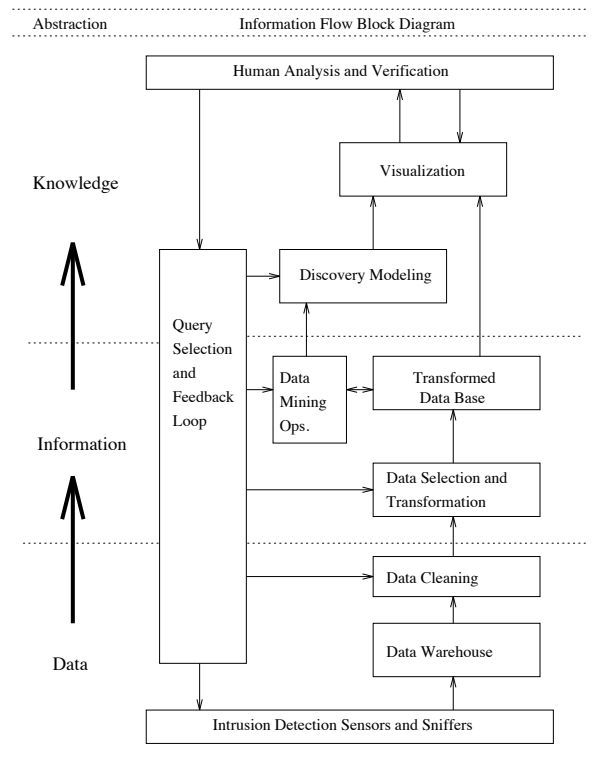


Fig. 3. Intrusion Detection Data Mining

lihood of characteristics and associations in selected data sets.

Rule Abduction is the development of IF-THEN-ELSE rules that describe associations, structures and the test rules.

Link or Tree Abduction is performed to discover relationships between data sets and interesting connecting pattern properties.

Deviation Analysis locates and analyzes deviations from normal statistical behavior.

Neural Abduction is the process of training artificial neural networks to match data, extract node weights and structure (similar to abducted rule sets).

As cybersensor information is mined into new ID knowledge, refined models are developed which seek to predict future events based on historical data. This process is known as *discovery modeling*. In addition, analysts require *visualization* tools to support the very well developed human process of pattern recognition. The entire data mining process is refined by adjusting parameters, sets, and associations in lower level processes.

Both the data mining and fusion process are in the very early stages of technical development. However, as network continue to grow and the expanding realms of cyberspace evolve, the marketplace will drive ID systems toward next generation capabilities. Integrated reasoning and decision support tools are emerging requirements for robust and reliable intrusion detection in complex internetworks.

V. CHALLENGES IN IDS FUSION

The preceding discussion starts to illustrate the complexity of designing reliable intrusion detection systems. These systems are required to fuse data and information from heterogeneous distributed cybersensors, where *cybersensors* are broadly defined as all hardware-software devices collecting cyberspace situational information (e.g. processor and network events that may be evidence of intrusion). One of the first challenges is to extend the groundwork introduced by Denning in [2] to develop a structured meta-language for generic intrusion detection - network management objects. A standard meta-language is required for Level 0 and Level 1 Object Refinement, data storage, cleansing, and primitive correlation.

Data refinement is simplified when a common meta-language for both intrusion detection and network management exist. The temporal calibration of numerous streams of raw data from heterogeneous sources are also required. Internetworking protocols are evolving and may be used to synchronize object and events in a distributed Internet environment. However, the security of TCP/IP information flows remain a critical issue.

Correlation in physical space compares observations to a physical coordinate system (e.g. the Euclidean distance between two measurements) to determine if there is a common source. Correlation in cyberspace requires the comparison of observations based on a different set of parameters such as source (IP address), network path, session flow, behavior or timing.

The automated identification and tracking of dynamic intrusion subjects (suspected intrusion events) in cyberspace are also formidable technical challenges. Imagine, if you will, intruders executing TCP-based attacks from numerous geographically dispersed networks, or initiates attacks with one network connection and continues with another, sequentially changing IP addresses. Tracking and assessing the threat of these classifications of cyberattacks require new technical solutions. These topics have not been adequately addressed, however the threats to critical infrastructures are emerging.

Hall [10] discusses mathematical techniques for multisensor data fusion. The application of these techniques to cyberspace ID systems is also quite complex. At the lowest level of inferences is the process of *data association*. These are example fusion concepts related to data association which are also requirements for cyberspace ID systems:

Gating: methods used to eliminate unlikely associations to reduce the number of associated pairs of network events to evaluate.

Association: the selection of metrics used to quantify the closeness or similarity between observed events.

Assignment: selection of the events to declare to be associated with the intrusion hypothesis, and hypothesis processing.

Parametric data is used to estimate basic parametrics

of network events. *Estimation* theory is required to infer intrusion attack rates, attack targets, origins and other cyberspace situational parametrics. The estimation and detection process is highly mathematical and processor intensive, drawing from sub-disciplines such as *optimization*, *least squares estimation*, and *sequential estimation*. Also required for cyber ID systems are complex error analysis algorithms and stochastic models for noise and cyber false alarm estimation [10].

The *identity declaration and pattern recognition* phase of the fusion model is a difficult technical problem because the level of inference very high. This is often done by extracting *features* which are abstractions of raw data. The basic parametric for pattern recognition is templating. Elementary forms of templating are used in current state-of-the-art ID systems. Future systems tracking coordinated multifaceted cyberspace attacks require *cluster analysis* techniques, *adaptive neural networks* and rules-based *knowledge systems*.

Classical Inference, *Bayesian Inference*, *Dempster-Shafer Method*, *Generalized EPT*, and *Heuristic Methods* are a few of the mathematical methods that are required in the decision-level identity fusion process. The reader interested in these techniques is kindly referred to both [10] and [11]. The application of these technologies to intrusion detection and network monitoring is required to realize the *cyberspace situational awareness* required for advanced ID systems.

The highest level of inferences, *knowledge fusion*, is also a very complex and challenging area. Visualize next generation ID systems which identify and track multiple hostile information flows for targets, attack rate, and severity in cyberspace. Determining the origin of highly sophisticated attacks in cyberspace will continue to grow in complexity as attackers become more cyberspace astute. The time allowed for network operators to trace (multiple) attack origins is a function of the attack rate and the potential damage (situation assessment). These are just a few of the exciting requirements of cyberspace ID systems. Dreaming, brainstorming, developing and articulating the engineering requirements for these next generation systems is the first step.

VI. CONCLUSION

The current state-of-the-art of ID systems is relatively primitive with respect to the recent explosion in computer communications, cyberspace, and electronic commerce. Organizations fully realize that cyberspace is a complex realm of vital information flows with both enabling and inhibiting technical factors. Identifying, tracking, classifying, and assessing hostile and inhibiting activities in this ever growing complex dimension is an enormous and fascinating technical challenge.

Multisensor data fusion is a multifaceted engineering approach requiring the integration of numerous diverse disciplines such as statistics, artificial intelligence, signal processing, pattern recognition, cognitive theory, detection theory, and decision theory. The art and science of data fusion is directly applicable in cyberspace for intrusion and attack

detection.

Dynamic cyber data mining operations are required to develop new intrusion detection models based on historical data in data warehouses. Hence, a significant research and development effort is required to bring next generation intrusion detection systems into the commercial marketplace. I hope this article, in some small way, stimulates the neurons of engineers and scientists interested in Internet security, and in particular, the research and development of advanced ID systems and cyberspace situational awareness.

REFERENCES

- [1] Bass, T., Freyre, A., Gruber, D. and Watt., G., *E-Mail Bombs and Countermeasures: Cyber Attacks on Availability and Brand Integrity*, IEEE Network, pp. 10-17, Vol. 12, No. 2., March/April 1998.
- [2] Denning, D., *An Intrusion-Detection Model*, IEEE Transactions on Software Engineering, Vol. SE-13, No. 2, pp. 222-232, February 1987.
- [3] Mukherjee, ., Heberlein, L., and Levitt, K., *Network Intrusion Detection*, IEEE Network Magazine, Vol. 8. No. 3, pp. 26-41, May/June 1994.
- [4] Denning, D. et al., *A Prototype IDIES: A Real Time Intrusion Detection Expert System*, Computer Science Laboratory, SRI International, August 1987.
- [5] Bauer, D. and Koblentz, M., *NDIX - An Expert System for Real-Time Network Intrusion Detection*, Proceedings of the IEEE Computer Networking Symposium, pp. 98-106, April 1988.
- [6] Hochberg, et al., *NADIR: An Automated System for Detecting Network Intrusion and Misuse*, Computers & Security, Elsevier Science Publishers, pp. 235-248, 1993.
- [7] Heberlein, L. et al., *A Network Security Monitor*, Proceedings of the IEEE Computer Society Symposium, Research in Security and Privacy, pp. 296-303, May 1990.
- [8] Snapp, S. et al., *A System for Distributed Intrusion Detection*, Proceedings of IEEE COMPCON, pp. 170-176, March 1991.
- [9] Waltz, E. and Llinas, J., *Multisensor Data Fusion*, Artech House, Boston, MA, 1990.
- [10] Hall, D., *Mathematical Techniques in Multisensor Data Fusion*, Artech House, Boston, MA, 1992.
- [11] Varshney, P., *Distributed Detection and Data Fusion*, Springer-Verlag, New York, NY, 1996.
- [12] Waltz, E., *Information Warfare Principles and Operations*, Artech House, Boston, MA, 1998.

Tim Bass (bass@silkroad.com) is an independent writer and consulting engineer specializing in network management and network security.