



Contents lists available at SciVerse ScienceDirect

Journal of Network and Computer Applications

journal homepage: www.elsevier.com/locate/jnca

Review

A survey of network flow applications

Bingdong Li^{a,b,*}, Jeff Springer^a, George Bebis^b, Mehmet Hadi Gunes^b^a Department of Information Technology, University of Nevada-Reno, NV 89557, United States^b Department of Computer Science & Engineering, University of Nevada-Reno, NV 89557, United States

ARTICLE INFO

Article history:

Received 28 June 2012

Received in revised form

17 September 2012

Accepted 12 December 2012

Available online 4 January 2013

Keywords:

Machine learning

NetFlow

Network traffic analysis

Network security

sFlow

ABSTRACT

It has been over 16 years since Cisco's NetFlow was patented in 1996. Extensive research has been conducted since then and many applications have been developed. In this survey, we have reviewed an extensive number of studies with emphasis on network flow applications. First, we provide a brief introduction to sFlow, NetFlow and network traffic analysis. Then, we review the state of the art in the field by presenting the main perspectives and methodologies. Our analysis has revealed that network security has been an important research topic since the beginning. Advanced methodologies, such as machine learning, have been very promising. We provide a critique of the studies surveyed about datasets, perspectives, methodologies, challenges, future directions and ideas for potential integration with other Information Technology infrastructure and methods. Finally, we concluded this survey.

© 2013 Elsevier Ltd. All rights reserved.

Contents

1. Introduction	568
2. Background	569
2.1. Net Flow	569
2.2. sFlow	569
2.3. IPFIX	569
2.4. Network flow analysis	570
3. Perspectives	570
3.1. Network monitoring, measurement and analysis	570
3.1.1. Network monitoring	570
3.1.2. Application monitoring	570
3.1.3. Host monitoring	571
3.2. Network application classification	571
3.2.1. Peer to peer network	571
3.3. User identity inferring	571
3.4. Security awareness and intrusion detection	571
3.4.1. Top N	571
3.4.2. Port scanning	572
3.4.3. Denial of service	572
3.4.4. Worms	572
3.4.5. Botnet	572
3.4.6. Policy validation	572
3.5. Issues of data error	573
4. Methodologies	573
4.1. Statistics	573

* Corresponding author at: Department of Information Technology, University of Nevada-Reno, NV 89557, United States. Tel.: +1 775 682 6805; fax: +1 775 784 4529.
E-mail addresses: bingdongli@unr.edu, libingdong@gmail.com (B. Li).

- 4.2. Machine learning 573
 - 4.2.1. Application classification 574
 - 4.2.2. Security awareness and anomaly detection 574
- 4.3. Profiling 574
- 4.4. Behavior-based approaches 575
- 4.5. Visualization 575
- 4.6. Anonymization 576
- 4.7. Analysis systems 576
 - 4.7.1. Optimization 576
 - 4.7.2. Sampling 576
 - 4.7.3. Distributed analysis system 577
- 5. Discussion 577
 - 5.1. Datasets 577
 - 5.2. Research perspectives 577
 - 5.3. Methodologies 577
 - 5.4. Challenges 577
 - 5.4.1. Feature representation and selection 577
 - 5.4.2. Real time analysis and data storage 577
 - 5.5. Future directions 578
 - 5.5.1. Distributed data collection and analysis 578
 - 5.5.2. Advanced analysis methodologies 578
 - 5.5.3. Integration 578
- 6. Conclusion 578
- Acknowledgments 578
- References 578

1. Introduction

Computer networks are playing a very important role in our daily life. Our dependency on computer networks is growing tremendously. Understanding what information flows in a computer network is important not just for network administrators but also for accounting, network planning, network security, forensics and counter-terrorism. Many governments require Internet Service Providers (ISP) to have capabilities of ‘lawful interception’ (LI) network traffic (Baker et al., 2012). Moreover, network flow can provide information for business relationship (Kind et al., 2006).

Network flow records high-level descriptions of Internet connections but not the actual data transferred. Collection and analysis of network flow information is more efficient than deep packet inspection and protects the privacy of users. This information helps to uncover both external activities and internal activities such as network misconfiguration and policy violation. Network flow information is supported by a wide range of products including Cisco NetFlow (Kerr and Bruins, 2001), Juniper’ cflowd, NetStream, and sFlow. These systems are all similar to NetFlow systems, and will be referred to as NetFlow-like in this survey.

It has been over 16 years since Cisco’s NetFlow was developed by Darren and Barry Bruins in 1996 (Kerr and Bruins, 2001). Research in network flow analysis has become very active in the recent years as observed in Fig. 1. It is necessary to look back what perspectives have been achieved, and what methods have been used and are more effective in order to move forward. This paper presents a survey of NetFlow-like applications that studies were published between 1998 and early 2012. Figure 1 shows the published paper distribution with respect to publication year. Our objective is to provide a better understanding of major achievements in the field by reviewing state of the art of approaches, perspectives, important issues and challenges, and suggesting directions for future research. Note that, we have focused mostly on studies using NetFlow-like data as input, emphasizing some of the latest approaches rather than attempting to provide a complete historical review of network flow applications.

Related reviews discussing similar aspects to this survey but not specific to NetFlow-like applications can be found in Introduction to Cisco IOS NetFlow (2012) for IP-Flow based intrusion detection, Zhu et al. (2008) for botnet detection, Nguyen and Armitage (2008) for

Internet traffic classification using machine learning, and Sommer and Paxson (2010) for discussion of using machine learning for network intrusion detection.

Traditionally, NetFlow-like analysis systems have been used for network monitoring, planning and billing. Recent research approaches, however, have focused more on network security analysis with the objective of detecting anomalous activities that traditional security infrastructures, such as intrusion detection systems (IDS), firewalls and anti-virus tools, cannot handle. These approaches employ advanced techniques such as machine learning. Moreover, new NetFlow-like analysis system design is moving toward distributed systems to provide more scalability, robustness and computational power for real-time in depth analysis.

Despite the popularity of sFlow and its wide deployment, few studies have focused on using sFlow as their data source.

The rest of this paper is organized as follows. Section 2 provides a brief introduction to NetFlow, sFlow, IPFIX, and network flow analysis. Section 3 reviews the key perspectives which have been addressed in the literature including networking monitoring, analysis and management, application classification, inferring user identity, and network security awareness. Section 4 explains the

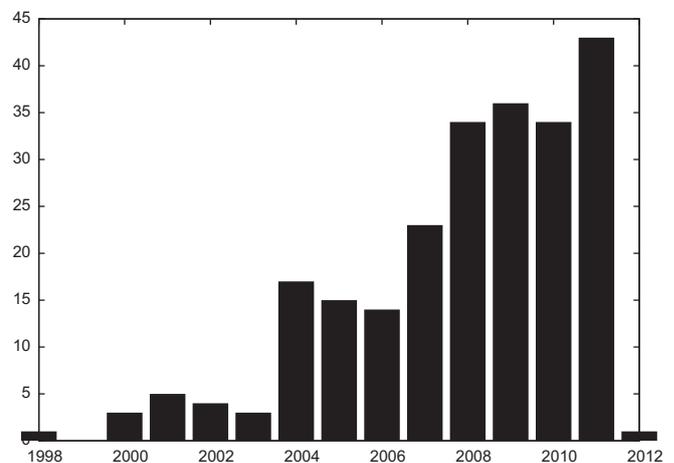


Fig. 1. Publications by year.

key methodologies and tools which have been employed for data analysis including statistic, machine learning, profiling, behavior-based approaches, sampling, visualization, and computational infrastructures. Section 5 discusses the limitations of existing approaches, challenges and future directions. Finally, Section 6 presents our conclusions.

2. Background

Network flow is defined as an unidirectional or bidirectional sequence of packets between two endpoints (from server to client or from client to server) with some common attributes. The most important key fields include: first and last time of the flow received, source/destination IP address, source/destination port number, layer 3 protocol type, type of services, bytes transferred, and input logical interface. Additional fields may be included that depend on the NetFlow version or configuration for export. They provide a rich set of traffic statistics including user, protocol, port, and type of services which can be used for a wide variety of purposes such as network security, network monitoring, traffic analysis, capacity planning, traffic classification, accounting, and billing. The general process of working with NetFlow includes capturing, sampling, generating, exporting, collecting, analyzing and visualizing.

There are various systems that capture NetFlow IP operational data from network links or devices: Ntop (Deri, 2012), NG-MON (Han et al., 2002), NetFlow, NetFlow-lite, sFlow, cflowd, Net-Stream, etc. In addition, IPFIX (2012) defines the standard IP flow format for exportation. NetFlow and sFlow are widely used systems while IPFIX is a new standard. In this section, we will give a brief introduction about NetFlow, sFlow, IPFIX, and traffic analysis.

2.1. Net Flow

NetFlow is a traffic monitoring technology developed by Darren and Barry Bruins in 1996 at Cisco (Kerr and Bruins, 2001). It defines how a router exports information and statistics of routed sockets. As a de facto industry standard, it is a built-in feature of most routers and switches from Cisco, Juniper, and other vendors. Network devices look at the packets arriving on the interfaces, and capture traffic statistics per flow based on configuration for sampling or filtering, then they create a flow cache, aggregate and export the data through UDP or Stream Control Transport Protocol (SCTP). NetFlow cache entry is created by the first packet of a flow, maintained for similar flow characteristics, and exported to collectors periodically based on flow timers or flow cache management. The export formats of version 1 to 8 are fixed. After version 9, extensibility and flexibility are added to integrate with MPLS, IPv6 and BGP, and user defined records. NetFlow versions 5 and 9 are the most popular versions. Sampled NetFlow is a variant originally introduced by Cisco to reduce computational burden by reducing number of NetFlow. It can be configured as predetermined n th packet or randomly selected interval. Figure 2 presents the basic process of NetFlow formation, exportation, storage and analysis. Due to the great value of network traffic and limited computational resources (memory, CPU and bandwidth), technologies of caching, sampling and UDP exportation were used. These can cause quality issues for the collected NetFlow data: (1) some new flows will not be counted when cache is full; (2) sampling reduces the accuracy of flows, especially when sampling rate is adjusted by the traffic rate; (3) exported flow records do not necessarily correspond to the order in which the flow traffic arrived at the router. There are varieties of NetFlow collectors and analysis tools

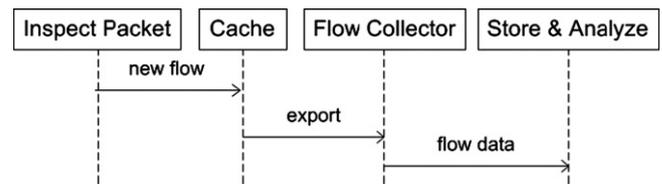


Fig. 2. NetFlow process.

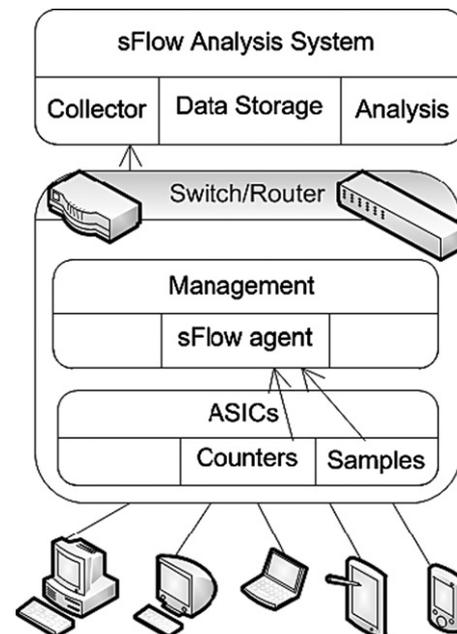


Fig. 3. sFlow process.

from commercial vendors such as Cisco, freeware or developed in-house for special purposes (Introduction to Cisco IOS NetFlow, 2012; NetFlow applications, 2012).

2.2. sFlow

Packet sampling of traffic flow (Duffield, 2004) has a long history before NetFlow was developed. sFlow was developed by InMon Inc. and has become an industry standard defined in RFC 3176. It is a technology using simple random sampling and supported by Alcatel, Extreme, Force10, HP, Hitachi by embedding the sFlow agent within switches and routers. The sFlow agent is a software process that combines interface counters and flow samples into sFlow datagrams and immediately sends them to sFlow collectors via UDP. Immediate forwarding of data minimizes memory and CPU usage. Packets are typically sampled by Application-Specific Integrated Circuits (ASICs) to provide wire-speed performance. sFlow data contains complete packet header and switching/routing information, and provides up to the minute view of the network traffic. sFlow is able to run at layer 2 and capture non-IP traffic as well. The sFlow collectors are servers that collect the sFlow datagrams. The official sFlow web site (sFlow Collectors, 2012) provides a list of available sFlow collectors. Figure 3 presents the basic components and processes of sFlow analysis.

2.3. IPFIX

IP Flow Information Export protocol (IPFIX) is an IETF standard for exporting network flow based on NetFlow version 9, and is defined in RFC 5101 for information transmitting protocols, RFC 5102 for information model, and RFC 5103 for exporting bidirectional flow.

IPFIX was designed to meet the fast growing requirements to observe network traffic, provide an extensible and flexible data model that can be customized, and support reliable and secure data transfer through SCTP, TCP and UDP. IPFIX flow definition is less restrictive than traditional flow definition. As standardization is underway, more vendors are going to support IPFIX.

2.4. Network flow analysis

Network flow analysis is the process of discovering useful information by using statistics or other sophisticated approaches. The basic process includes capturing, collecting and storing data, aggregating the data for query and analysis, and analyzing the data and results for useful information. This information is mostly related to network management, measurement, and network security. There are different ways to collect network flow data: SNMP, NetFlow, raw packet, or auditing data from network infrastructure such as IDSeS, Firewalls, and VPN gateways. Typically, there are two strategies: depth-first when there is known information and clear purpose, or breadth-first when looking for a general view of the network without a particular purpose.

Deep packet inspection needs packet level information and consumes more computational resources. Flow level analysis, such as NetFlow and sFlow, consumes less computational resources. There are many products and tools developed by industry or open source community. Analysis of network flow information has become crucial as the Internet has become the living blood in our society and is expanding at a fast pace around the world. There are many challenges in analyzing network flow data, such as huge amount of data due to networks becoming larger and faster, limited high-level information, and complex statistical properties. As discussed in Sections 3 and 4, various perspectives have been analyzed and many algorithms have been developed.

3. Perspectives

In this section, we survey the main research perspectives of network flow applications. In particular, we cover network monitoring, measurement and analysis, application classification, user identity inferring, security awareness and intrusion detection, and issues related to error and bias in NetFlow collection and analysis. Figure 4 shows the distribution of published papers with respect to five main perspectives: monitoring, classification, user identifying, security, and issues related to errors. As it can be observed, network security has been the main research topic using NetFlow data.

3.1. Network monitoring, measurement and analysis

Network monitoring and measurement provide valuable information to network administrators, ISPs and content providers. Compared to other technologies, such as SNMP or Windows Management Instrumentation(WMI), network flow data contain additional information for further analysis. For example, they can provide bandwidth analysis, specific protocol monitoring, and system performance, etc. Monitoring based on NetFlow can be categorized as:

- *Network monitoring*: provides information about routers and switches as well as network-wide basis view, and is used for problem detection along with efficient troubleshooting.
- *Application monitoring*: provides information about application usage over the network, and is used for planning and allocation of resources.
- *Host monitoring*: provides information about user utilization of network and applications, and is used for planning, network access control, violating security policy.

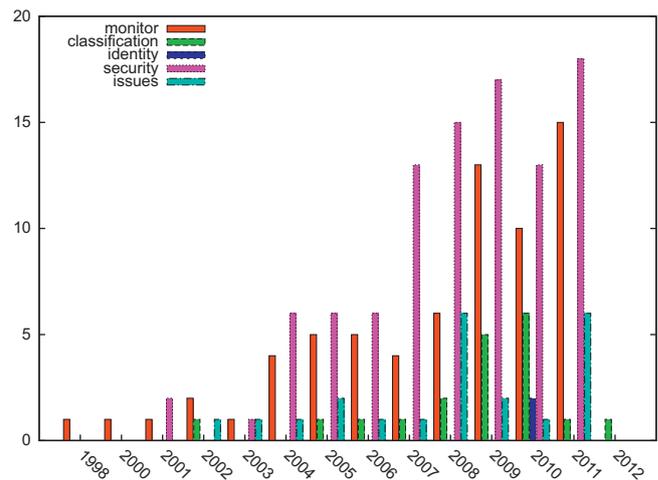


Fig. 4. Analyzed perspectives by years.

- *Security monitoring*: provides information about network behavior changes, and is used to identify DoS attacks, viruses and worms, and network anomalies.
- *Accounting and Billing*: provides network metering, and is used for billing.

In this section, we focus on network, application, user and resource monitoring while Section 3.4 focuses on security related monitoring. In the following, we discuss specific research perspectives.

3.1.1. Network monitoring

Many aspects of networks have been studied using NetFlow data, including network performance based on round-trip time (Strohmeier et al., 2011), delay measurement (Köandgel, 2011; Lee et al., 2011), connectivity (Schatzmann et al., 2011), misuse of bandwidth (Mansmann et al., 2009), traffic characterization (Kundu et al., 2009), finding heavy hitters (Truong and Guillemin, 2009), monitoring for special purpose QoS (Li, 2010), and diagnostic of troubleshooting (Sukhov et al., 2009).

3.1.2. Application monitoring

Liu and Huebner (2002) investigated the stochastic characteristics of distributions of flow length, packet size, throughput, etc. for the popular and bandwidth consuming applications. Kalafut et al. (2009) proposed a heuristic method to differentiate wanted and unwanted traffic based on the sampled NetFlow data.

VoIP. Voice over IP (VoIP) service is widely used, however, it introduces security threats that include Session Initiation Protocol (SIP) scanning, SIP flooding, and Real-time Transport Protocol (RTP) flooding. Lee et al. (2009) developed a system that can monitor VoIP service and detect VoIP network threats based on NetFlow statistics and behaviors. Kobayashi and Toyama (2007) presented a method for measuring VoIP traffic fluctuation by using NetFlow and sFlow based on the variance of the interval of the target RTP packets. Lucas (Deri, 2012) provided an open source VoIP monitoring system based on protocol characteristics.

Mobile network. Sinha et al. (2003) analyzed the flow-level upstream traffic behavior from Broadband Fixed Wireless (BFW) and Digital Subscriber Link (DSL) to provide traffic characteristics of these networks. Moghaddam and Helmy (2011) used wireless NetFlow data to measure and simulate user behavior and provide information for future mobile network design.

IPv6. With the transition from IPv4, there is a need to understand IPv6 usage including user behavior, traffic volume, transitional

technologies, assignment of IPv6 address, IPv6 percentage of network traffic. NetFlow can provide this information including application types, usage of transitional technologies of IPv6 to IPv4, interface identifier assignment schemes, etc. (Shen et al., 2009; Zhang and Meng, 2011).

3.1.3. Host monitoring

Host profile and relationships in the network can be used for resource planning as well as for network security analysis. Caracas et al. (2008) proposed an algorithm based on NetFlow data to describe the dependencies among computer systems, software components, and services. Kind et al. (2006) presented a method to uncover the relationships between IT infrastructures using NetFlow data. Chen et al. (2011) developed novel heuristics to analyze characteristics and correlations between inter-data centers and client traffic, provided insights into data center design and operation. Several methods have been proposed for profiling behaviors on the end hosts (Wei et al., 2006; Xu et al., 2011, 2005). Behavior-based approaches will be discussed in detail in Section 4.4.

3.2. Network application classification

Network application classification classifies network traffic into certain application categories which can be coarse- or fine-grained. Network application classification is a challenging task because of obfuscation techniques such as content encryption, dynamic ports, and proprietary communication protocols. Classification approaches can be divided into four categories: port based, payload-based, heuristic based using transport layer statistics, and machine learning based. Port based approaches are no longer reliable because of certain applications that randomly assign ports. Payload based approaches do not work on encrypted traffic, and are resource-intensive and scale poorly with high bandwidth. Approaches based on heuristic and machine learning approaches provide alternative methods.

There are many reasons for network traffic classification. Network administrators need information for applications running at the network (e.g., file sharing), whether they are legitimate users or worms. ISPs and content providers need the information for quality of service assurances. Research in this area has conducted for over 10 years, but it is still growing. There is a list of 68 published papers and 86 datasets collected in CAIDA web pages (Internet Traffic Classification, 2012). There are several surveys on network application classification using traffic classification approaches (Kim et al., 2008) and machine learning approaches (Nguyen and Armitage, 2008). We discuss machine learning approaches in detail in Section 4.2.1. It is worth mentioning that Perelman et al. proposed a method that investigates the application signatures of web browsers, mail client, or media-players in network flow (Perelman et al., 2011). Peer to peer networks have become a major security concern and the focus of most network classification studies. We discuss peer to peer classification in detail below.

3.2.1. Peer to peer network

Peer to Peer (P2P) networks have been widely utilized for file-sharing, video distribution and voice communications. They consume more Internet traffic than traditional applications, and have been a concern for network administrators and a challenge for network security. There is interest from ISPs and network administrators to identify and control the P2P network traffic (Gossett et al., 2010; Zha and He, 2011). NetFlow provides an alternative approach that is more efficient in terms of storage and processing than deep packet inspection (DPI). Recently, there has been

considerable effort on NetFlow P2P analysis. These include methods based on: (a) default P2P port for heavy-hitters (Wagner et al., 2006), (b) port usage pattern of specific P2P network such as BitTorrent (Bo et al., 2009; Gossett et al., 2010), (c) flow statistic characteristics such as packet length and time-interval (Bo et al., 2009; Qun et al., 2011; Xu et al., 2009), (d) TCP flags that a host, as both client and server, send/receive a packet with both SYN and ACK at the same time (Jinsong et al., 2009), (e) machine learning that using features such as IP address and port, packet size, bytes exchanged. Among machine learning approaches we discuss in Section 4.2.1, six out of eight classify P2P traffic.

3.3. User identity inferring

Identifying a person based on extrinsic biometric is not new; well-known examples include signatures and keystrokes. Inferring user identity based on network flow patterns, however, is a new field. Melnikov and Schönwälder (2010) discussed the potential of inferring user identity using NetFlow feature distribution and cross-correlation of various trace parameters and relationships among packets. Even though the reported results were preliminary, additional research may yield more promising results.

3.4. Security awareness and intrusion detection

In this section, we focus on security related awareness, detection and monitoring. Table 1 provides a list of studies that provide perspectives on security awareness and intrusion detection. Table 3 also lists approaches that use machine learning approaches. IDSes can be categorized based on how they identify intrusions: anomaly-based, misuse-based (knowledge-based or signature-based), or combination of both anomaly and misuse-based (Sperotto et al., 2010). Alternatively, IDSes can be categorized based on what they target: host-based, network-based or both.

Network anomaly detection refers to finding patterns that are not expected users behaviors, also known as anomaly-based IDS. Compared with misuse-based IDS, these patterns are previously unknown. Most content-oriented systems belong to knowledge-based detection, which looks for known signatures of malware by inspecting traffic packets. Most behavior-oriented systems belong to anomaly-based detection, which differentiate anomalous behavior from normal behavior. NetFlow based IDSes use existing NetFlow data and limited information, and avoid privacy issues compared to content-oriented approaches. However, NetFlow based IDSes are more difficult because of limited information in the NetFlow data. Consequently, recent research has shown that machine learning approaches are better than statistical and streaming methods. Sperotto et al. (2010) conducted an overview of IP flow-based intrusion detection that focused on flow-based IDS, concept of flows, classification of attacks, and defense techniques. In the following, we discuss perspectives of security awareness and intrusion detection that can be achieved using NetFlow data.

3.4.1. Top N

Top N refers a set of statistic and models of NetFlow data. They reflect the basic network status. It is relatively simple with NetFlow analysis. It can be used to find the big talkers or heavy-hitters. It also can be used for abnormal traffic detection (Zhang, 2009).

Table 1
Summary of security awareness and intrusion detection.

Year	Methodology	Perspective
2001 (Erbacher, 2001)	Histogram and chart	IDS
2001 (Kotsokalis et al., 2001)	Statistic	DoS and DDos
2004 (Yin et al., 2004)	Links between machines or domains	IDS
2004 (Kim et al., 2004)	Statistic patterns	DoS and DDos
2005 (Dubendorfer and Plattner, 2005)	Host behavior based	Worm outbreaks
2005 (Dubendorfer et al., 2005)	IP aggregation	Detection and monitoring
2006 (Ren et al., 2006)	Flow aggregation	IDS
2007 (Rehak et al., 2007)	Trust and reputation model	IDS
2007 (Dressler et al., 2007)	Flow signature and honeypot logs	Worm detection
2008 (Chan et al., 2008)	Heuristics	Worm detection
2008 (Zhenqi and Xinyu, 2008)	Statistic	Anomaly detection
2009 (Krmicek et al., 2009)	Heuristics	NAT detection
2009 (Vykopal et al., 2009)	Decision tree	Dictionary attack
2009 (Frias-Martinez et al., 2009)	K-means	Behavior-based NAC
2009 (Yin and Nianqing, 2009)	Information theory	Risk detection
2009 (Zhang, 2009)	Statistic	Top N detection
2009 (Vliek, 2009)	Statistic	Spam machines
2010 (Čeleda et al., 2010)	NBA	Malware
2010 (Hsiao et al., 2010)	Spatial-temporal aggregating	Malicious website detection
2011 (Sawaya et al., 2011)	Statistic of host behavior	Attack detection
2011 (Ke-xin and Jian-qi, 2011)	Dynamic entropy	DoS
2011 (Galtsev and Sukhov, 2011)	Statistic	DDoS and port scan
2011 (François et al., 2011)	Host behavior and PageRank	Botnets detection
2011 (Sperotto and Pras, 2011)	Time series	IDS
2011 (Francois et al., 2011)	PageRank	Botnets detection

3.4.2. Port scanning

Port scanning is the act of systematically scanning a computer's ports, and is usually done by using small packets that probe the target machines. In most network attacks, port scanning is the first reconnaissance step. Port scanning can be classified into three categories: scanning many ports on a single host, scanning a single port on many hosts, and combination of both. Detection of port scanning is addressed in most studies cited in Table 1. Approaches include host incoming/outgoing connections, probability of entropy, Bayesian logistic regression, distances from baseline models, and machine learning.

3.4.3. Denial of service

A denial-of-service (DoS) or distributed denial-of-service (DDoS) attack is an attempt to make the target host or network resource unable to respond to its requests. Detection of DOS or DDoS is addressed in most flow based IDSes. Gao et al. (2006) proposed a resilient DoS detection based on sketch-based schemes that use a hash table for storing aggregated flow measurement. Kim et al. (2004) described different DoS attacks based on traffic patterns and presented a network anomaly detection method that can detect flooding attacks. New developments include using novel dynamic entropy to measure the anomaly (Ke-xin and Jian-qi, 2011), an attack detection method based on statistic aggregation that can detect DDoS and port scanning (Galtsev and Sukhov, 2011). Table 1 provides a list of related studies.

3.4.4. Worms

A worm is a standalone malicious program that replicates across the networks by exploiting software vulnerabilities or tricking users to execute it by social engineering. Worms can cause mildly annoying effects, damaging data or software, DoS, stealing data, etc. Detection of worms can be categorized as trap-oriented, packet-oriented and connection-oriented (Chan et al., 2008). Detection of port scanning is one of the important steps for worm detection, and hence many similar approaches are used in

both types of detection. NetFlow-like approaches are connection-oriented and include: analysis of host behavior on the basis of incoming and outgoing connections (Dubendorfer and Plattner, 2005), correlation between NetFlow data and honeypot logs (Dressler et al., 2007), and detecting hit-list worms using protocol graphs (Collins and Reiter, 2007). Chan et al. (2008) proposed FloWorM system that includes tracker, analyzer and reporter based on NetFlow data. Abdulla et al. (2011) presented a support vector machine (SVM) method based on the fact that a scanning activity or email worm initiates a significant amount of traffic without DNS.

3.4.5. Botnet

Botnets are malware at the infected target and controlled by a remote entity known as bot-master. They have become one of the major security threats credited for DDoS, spamming, phishing, identity theft, and other cyber crimes. Many botnets rely on communication channels varying from centralized IRC and HTTP to decentralized P2P networks. Detection of a botnet is relatively more difficult than detection of port scanning and worms. Chan et al. (2008) conducted a survey on understanding, detection and tracking, and defending against botnets. Recent approaches use advanced methodologies and combine host and network level information. Zeng et al. (2010) proposed a method that combined host and network-level information with protocol-independent detection. BotCloud's detection is based on MapReduce and combining host and network approaches (François et al., 2011). BotTrack's tracking is based on PageRank of NetFlow data and host behavioral model (François et al., 2011). Finally, Barsamian used a network statistical behavioral model for botnet detection (Barsamian, 2009), and Weststrate used heuristic methods to find botnet servers (Weststrate, 2009).

3.4.6. Policy validation

Peer-to-peer networks can be used legitimately, or misused by botnet, or violate network usage policy. Section 3.2.1 details peer-to-peer classification using NetFlow data. Krmicek et al.

(2009) proposed an approach to detect the use of unauthorized Network Address Translation (NAT) via a heuristic method based on NetFlow data. NetFlow data can also provide information about legitimate flows denied by the security policy, and help network administrator with troubleshooting. Frias-Martinez (Frias-Martinez et al., 2009) proposed a behavior-based network access control mechanism with a true rejection rate of 95%.

3.5. Issues of data error

NetFlow data is exported using UDP. Data can be lost due to overloaded segments between routers and collectors, an overload of collectors with benign traffic increases, burst nature of NetFlow traffic, or attacks in progress. Similarly, errors may happen in the process of sampling, transporting and collecting. In order to address these problems, several methods have been proposed. Cohen et al. (2008) proposed a framework for calculating confidence intervals to address the estimation errors in a multistage combination of sampling and aggregation. Trammell et al. (2011) characterized, quantified, and corrected timing errors, which are consequence of Cisco NetFlow version 9 protocol design that estimates the true base time from derived base time information. Rohmad et al. (2008) proposed an enhanced NetFlow version 9 using nProbe GPL. Fioreze et al. (2009) investigated the trustfulness of NetFlow measurements and found that octets and packets are reliably reported, but the flow duration of samples are shorter than the actual duration. Zhu et al. (2011) studied the errors of utilized bandwidth measurement of NetFlow, and provided guidance for correctly estimating the utilized bandwidth. Finally, Ricciato et al. (2011) described a methodology to estimate one-way packet loss from IPFIX or NetFlow flow records.

4. Methodologies

In this section, we review various methodologies used to analyze NetFlow data. Figure 5 provides a chronological summary of the methodologies discussed in this section. As it can be observed, a considerable number of studies have focused on using machine learning algorithms and real time analysis.

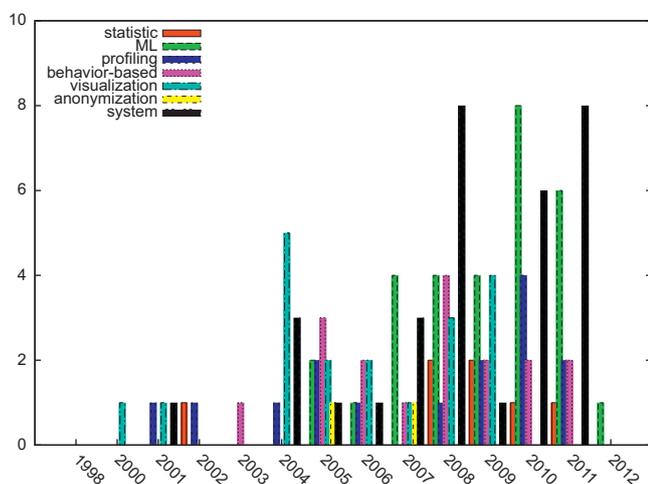


Fig. 5. Methods by years.

4.1. Statistics

Statistic approaches are the most common methods in NetFlow analysis. In general, it is the basic step before applying heuristic-based approaches, machine learning and visualization. NetFlow data contains statistics of network flow information generated and exported from routers. Duffield et al. (2002) investigated the resource usage of NetFlow formation and exportation as well as statistical properties of original traffic from sampled traffic data. Proto et al. (2010) proposed a statistical model for network intrusion detection system. Sawaya et al. (2011) proposed an approach of attack detection based on traffic flow statistics of hosts. Barsamian (2009) proposed a botnet detection method using statistical signatures. Bin et al. (2008) proposed an analysis and monitoring system using NetFlow statistic, and an IDS based on variance similarity.

Compared to other approaches, statistical approaches are usually easier to implement, provide accurate results and consume less resources. However, statistical approaches are good only for known cases and lack the ability to adapt to new cases.

4.2. Machine learning

Machine learning represents a collection of methods for discovering knowledge by searching for patterns. Machine learning refines and improves its knowledge base by learning from experience. The basic learning types are listed below:

- *Classification*: classify inputs to labeled outputs.
- *Clustering*: group inputs into clusters.
- *Association*: discover interesting relations between features.
- *Prediction*: predict outcome in terms of a numeric quantity.

Machine learning schemes include information theory, neural networks, support vector machines, genetic algorithms, and many more (Sommer and Paxson, 2010). Machine learning applications require the collection of training and test datasets and depend on algorithms for feature extraction, feature selection, and learning. Initially, the system is trained using example data to learn specific data associations; then, the system is deployed in a similar environment where test data is used for classification. In this section, we provide a survey of machine learning approaches in NetFlow applications, which include traffic classification, anomaly detection, and security awareness.

Selecting an appropriate set of features for a specific problem is critical. Example of features are shown in Tables 2 and 3, and are categorized as: (1) *basic features* such as NetFlow data fields, source and destination IP address and port, network interface, transport protocol, type of service, start and finish timestamps, cumulative TCP flags, number of bytes and packets transmitted, and MPLS labels; (2) *derived features* such as flow length (finish time–start time), average packet size (bytes/number of packets), average flow rate (bytes/length), average packet rate (number of packet/length), aggregation of IP subnet and traffic load bytes, percentage of traffic load on a node, percentage of traffic load at the current sub-tree with time period and aggregation threshold (Wagner et al., 2011, 2010); (3) *application specific heuristics* such as webmail traffic (Schatzmann et al., 2010) that has properties as close service proximity, daily and weekly pattern, and duration of client session; and (4) *advanced features* such as abacus signature, degree distribution, self-similarity of flow interval, entropy, kernel function, mutual information and Hellinger distance (Valenti and Rossi, 2011), or data fusion with other log files such as Snort, DNS related requests (Abdulla et al., 2011) (number of DNS requests, response, normals, and anomalies for each host over a certain period of time).

Table 2
Summary of machine learning approaches of network application classification.

Year	Algorithm	Accu. (%)	Feature	Application
2007 (Jiang et al., 2007)	NBKE	91	Basic ^a and derived ^b	P2P, email, Multi-media
2009 (Carela-Espanol et al., 2009)	DT	90	Basic	P2P, VoIP, DNS, email, FTP
2010 (Chaudhary et al., 2010)	Clustering	90	Application ^c	SNMP, email, DNS, IRC
2010 (Rossi and Valenti, 2010)	SVM	90	Advanced ^d	P2P
2010 (Schatzmann et al., 2010)	SVM	94	Application	Webmail
2010 (Barlet-ros and Cabellos-aporicio, 2010)	DT	90	Basic and derived	P2P, HTTP, VoIP, DNS, FTP, email, games
2011 (Valenti and Rossi, 2011)	SVM	70	Advanced	P2P
2012 (Liang and Jian, 2012)	BN	95	Derived	BULK, email, P2P

^a Basic NetFlow data fields.

^b Calculation and aggregation of basic features.

^c Application specific properties from basic and derived features.

^d Abstract information from basic and derived features.

Table 3
Summary of machine learning approaches of anomaly detection.

Year	Algorithm	Feature	Dataset	Perspective
2005 (Lakhina et al., 2005)	Cluster	Advanced ^a	Internet	Anomaly
2007 (Liu et al., 2007)	Multiclass SVM	Advanced	Internet	NSSA
2008 (Wang and Guo, 2008)	GA-based	Derived ^b	Non-NetFlow ^c	DDoS
2010 (Wagner et al., 2010)	Kernel	Derived	Internet	Monitoring
2010 (Strasburg et al., 2010)	SVM	Derived	Intranet	Masquerade
2011 (Abdulla et al., 2011)	SVM	Application ^d	Non-NetFlow	Worm
2011 (Wagner et al., 2011)	SVM	Derived	Internet	Attacks
2011 (Wagner et al., 2011)	SVM	Advanced	Internet	Attacks
2011 (Winter et al., 2011)	SVM	Basic ^e and derived	Non-NetFlow	IDS

^a Abstract information from basic and derived features.

^b Calculation and aggregation of basic features.

^c Simulation or log data.

^d Application specific properties from basic and derived features.

^e Basic NetFlow data fields.

Methods for feature selection include symmetric uncertainty (Jiang et al., 2007), information gain (Strasburg et al., 2010), subgroup, keyword selection, gradually reduction based on efficiency (Liu et al., 2007), and rough sets. The type of datasets and features being employed are very important for a successful machine learning approach. Typically, a large dataset is necessary to cover various relations in the data, including temporal and spatial relations. Training data has to be attack-free or attack-specific, both of which are difficult to obtain. The datasets in Table 3 can be categorized as (a) Internet backbone of more than one week period, (b) Internet backbone of less than one week period, (c) Intranet of more than two weeks, and (d) simulated data or honeypot log.

4.2.1. Application classification

Nguyen and Armitage (2008) surveyed the application of machine learning techniques for traffic classification from 2004 to 2007; even though NetFlow was not specified as analysis dataset, but the basic methodologies are applicable to NetFlow data. Kim et al. (2008) conducted an evaluation of traffic classification using traces with collected payloads. Their evaluation included seven machine learning algorithms: *Naive Bayes* (NB), *Naive Bayes Kernel Estimation* (NBKE), *Bayesian Network* (BN), *C4.5 Decision Tree* (DT), *k-Nearest Neighbors*, *Neural Networks*, and *Support Vector Machines* (SVM). They concluded that SVM consistently achieved higher accuracy. Soysal and Schmidt (2010) conducted more specific evaluations and comparisons of BN, DT and *Multilayer Perceptrons* on flow-based network traffic classification using flow trace data. They concluded that BN and DT are suitable for Internet traffic flow classification.

Nor and Mohd (2009) evaluated a large number of machine learning algorithms in terms of their performance on NetFlow data with the objective of classifying HTTP, gmail, and video streaming. The highest accuracy machine learning algorithms had an accuracy more than 99.33%. Unfortunately, they did not provide information about the features used. Table 2 summarizes the algorithms, accuracy, features and data types for traffic classification using NetFlow data. Since accuracy varies considerably, there is a need to evaluate these algorithms and features on the same dataset.

4.2.2. Security awareness and anomaly detection

Table 3 provides a summary of machine learning algorithms for anomaly detection in terms of algorithms, features, and research perspectives. The highest reported detection rate is 98% (Winter et al., 2011). Sommer and Paxson (2010) found that applying machine learning for network anomaly detection is harder than in other domains. This is mainly due to the great variety of traffic and the fundamental nature of machine learning approaches that are better suited at finding similarities than identifying relationships that are not present in the training data.

4.3. Profiling

Network profiling is an important step for further analysis. Various profiling levels have been discussed in the literature including user, application, host, and network profiling.

- *User profiling*: There is limited work on the user profiling based on NetFlow data. Melnikov and Schönwälder (2010) proposed

a set of correlation and distribution of user flow data related to time and packet to identify a users. Different user behavior based approaches have employed various features that are discussed in Section 4.4. User profiling research, however, may provide helpful information for network security in the future.

- **Application profiling:** Liu and Huebner (2002) discussed the stochastic characteristics of some of the most popular applications (i.e., FTP, HTTP, SNMP, NNTP, DNS, and Napster): flow length and time by probability density function and tail distribution, average packet size distribution, and average throughput distribution. Karagiannis et al. (2005) proposed traffic patterns of social behavior, function (provider or consumer), and application ports that were used to classify traffic based on heuristic rules.
- **Host profiling:** Wei et al. (2006) proposed an approach for Internet host profiling using a data structure that can be expressed in XML-like format at Listing 1, where the communication similarity is the average of Dice similarity values for the host. Xu et al. (2011) proposed an approach based on bipartite graphs to represent host communication and one-mode projection of bipartite graphs to capture the social-behavior similarity of end hosts as Fig. 6. For networks with few hosts, we need more detailed information for further analysis. Minarik et al. (2009) proposed a host behavior profiling based on the bi-directional NetFlow that use communicating peers (number of servers contacted, clients answered, and single flows), amount of traffic (amount of requests, replies, and single flows), and traffic structure (number of client, server and single flows). Frias-Martinez et al. (2009) defined a host behavior profile that contains seven features: the total number of flows, average flow size, average flow duration, total number of packets contained in all flows, average number of packets per flow, total number of unique IP addresses contained in all flows, and average packet size.
- **Network profiling:** Cho et al. (2001) proposed *Aguri* tree, an aggregation-based traffic profile that aggregates small volume flows with a fixed number of nodes in an IP tree for spatial measurement. Jiang et al. (2010) characterized network prefix-level traffic profiling as daily traffic volume, distributions (over time, direction, applications, and flow size), and ratio of upload-download. Lakhina et al. (2004) described *Origin-Destination* flows using a routing metric, and further analyzed using *Aguri* tree to include time, features (i.e., source and destination address, source and destination port) and volume to represent both time and space attributes (Lakhina et al., 2005).

4.4. Behavior-based approaches

Recently, behavior-based approaches to network security have received attention (Geer, 2006). Compared to signature-based approaches, behavior-based approaches first learn normal

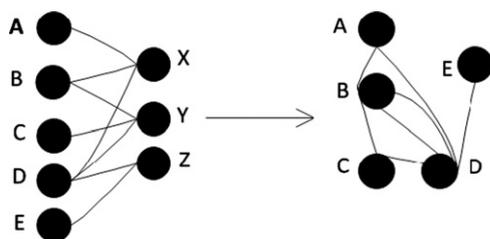


Fig. 6. Bipartite graph (left) and one-mode projection (right).

behaviors, and then detect anomalies. This approach has been applied with many research perspectives: application classification, anomaly detection, zero day attack detection, network access control (Frias-Martinez et al., 2009), and network design (Sinha et al., 2003). Types of behavior-based approaches include threshold, statistic and learning-based. Levels of behavior-based approaches include ISP-based Internet backbone behavior (Dubendorfer and Plattner, 2005; Xu et al., 2011, 2005), network behavior (Dubendorfer and Plattner, 2005; Rehak et al., 2008; Vlieg, 2009), user behavior (Melnikov and Schönwälder, 2010), host behavior (Dubendorfer and Plattner, 2005; Karagiannis et al., 2005; Mansman et al., 2007; Taylor et al., 2008; Celeda et al., 2010; Xu et al., 2011) and application (or protocol) behavior (Karagiannis et al., 2005; Lee et al., 2009; Soysal and Schmidt, 2010).

Box 1–Internet Host Profile (Courtesy of Wei et al., 2006).

```

A survey of network flow applications
< host >
  ip_address
  daily_destination_number
  daily_byte_number
  average_TTL
  < tcp_service > port1port2... < /tcp_service >
  < udp_service > port1port2... < /udp_service >
  < communication >
    < tcp_communication >
      destination_address
      daily_byte_num
      daily_connection_num
      average_duration_time
      < port > port1port2... < /port >
    < /tcp_communication >
    < udp_communication >
      destination_address
      daily_byte_num
      daily_packet_number
      < port > port1port2... < /port >
    < /udp_communication >
  < /communication >
  communication_similarity
< /host >

```

4.5. Visualization

Network visualization provides interactive visual displays for exploration of network traffic. It is a challenging task to visualize a large amount of information and provide sufficient level of detail to be meaningful and useful. Visualization can be at different levels of network abstraction (i.e., whole network, individual machine, and between whole network and individual machine), and described by different mechanisms (histogram, chart or glyph-based and 3D graph). Table 4 presents a chronological summary of related studies with their abstract level, mechanism of data processing and visualization, and research issues. Most applications use statistics and chart methods; whereas few applications use advanced methodologies such as machine learning, graph theory and quad-tree. In terms of research perspectives, most of them focus on security detection while others provide network monitoring.

Besides the approaches summarized in Table 4, several other projects are worth mentioning. *NFSen* (*NFSen*—*Netflow* Sensor, 2012) is an open source, graphical web based front-end tool.

Table 4
Summary of NetFlow visualization applications.

Year	Abstract	Mechanism	Perspective
2000 (Plonka, 2000)	Network	Aggregate traffic volume of protocols, chart	Network protocol and traffic amount
2001 (Erbacher, 2001)	Multiple	Histogram and chart	IDS
2004 (Yin et al., 2004)	Multiple	Links between machines or domains, graph	IDS
2004 (Lakkaraju et al., 2004)	Multiple	Activities of IP, histogram, glyph-based graph	Security situational awareness
2004 (Ball et al., 2004)	Multiple	Map between internal and external traffic, graph	Security
2004 (McPherson et al., 2004)	Network	Aggregation based on port, chart and graph	Security event detection
2005 (Dubendorfer et al., 2005)	Network	IP aggregation of traffic bursts, chart & graph	Worm detection and backbone monitoring
2005 (Patwari et al., 2005)	Network	Manifold learning, chart	Monitoring, detection
2006 (Oberheide et al., 2006)	Individual	Extended the quad-tree, 3D navigation and playback	Internet traffic
2006 (Oslebo, 2006)	Network	Network statistics of protocols, chart	Network statistics
2006 (Ren et al., 2006)	Multiple	Statistic, flow aggregation, chart & graph	IDS
2007 (Mansman et al., 2007)	Multiple	Host behavior, force-directed graph	Host behavior
2008 (Minarik and Dymacek, 2008)	Individual	Graph theory, graph	Network traffic
2008 (Fischer et al., 2008)	Network	TreeMap with splines, chart and graph	Network security monitor
2008 (Taylor et al., 2008)	Multiple	Aggregate data per port, 3D graph	Intrusive behavior
2009 (Singh et al., 2009)	Network	Based on simple K-means clustering, chart	Detect anomalies
2009 (Choi et al., 2009)	Network	Pattern of shape, graph	Network attacks
2009 (Taylor et al., 2009)	Multiple	Aggregate and map, graph and chart	Network monitoring
2009 (Goodall and Sowul, 2009)	Multiple	Aggregate, tree view, Geo-location, chart and graph	Network security
2012 (Shelley and Gunes, 2012)	Multiple	Sphere	Network traffic

It aggregates network traffic by protocols, direction or hosts using charts, and is used for network investigation. *AURORA* (2012) is an IBM research project for traffic analysis and visualization. It was designed for large networks, supports multiple levels of abstraction, and uses chart and graph to visualize traffic, anomaly detection or real time traffic flow. Finally, the Spinning Cube of Potential Doom (Lau, 2004) is a 3D display of network links for anomaly detection and visualized as a cube.

4.6. Anonymization

There is a need to anonymize NetFlow data to protect the privacy when the data is shared among parties. There are several approaches for NetFlow specific anonymization. Slagell and Luo (2005) proposed an anonymization tool for sharing network logs for computer forensics. Their tool can anonymize the common fields in multiple ways. Similarly, Foukarakis et al. (2007) proposed an anonymization tool with flexible features and high-performance.

4.7. Analysis systems

As the traffic volume is very large, methodologies to improve performance of capturing, collection, and analysis are needed. There are three commonly used methods to reduce data size: aggregation, filtering, and sampling (Duffield, 2004). In the following, we will survey optimization, sampling, and distributed analysis systems.

4.7.1. Optimization

Optimization can be applied in many stages of the NetFlow analysis process: capturing, collecting and analyzing. Bouhtou and Klopfenstein (2007) proposed mathematical models to select the NetFlow interfaces based on robust optimizations to deal with probabilistic constraints. Schatzmann et al. (2011) proposed a method of *Successive c-Optimal Design* to select NetFlow interfaces and find the optimal sampling rates. Hu et al. (2009) proposed an entropy based adaptive flow aggregation algorithm to improve efficiency of storage and export, and improve the accuracy of legitimate flows. Zadnik et al. (2005) proposed an architecture of network flow monitoring adapter based on hardware platform *COMBO6*, which is able to monitor one million simultaneous flows

on a 2 Gbps link. Nagaraj et al. (2008) proposed an efficient aggregation techniques to speed up querying based on attributes and filter condition of queries.

4.7.2. Sampling

Sampling network flow reduces the burden of handling massive volumes of flow data in collection, storage and analysis. Duffield (2004) conducted a review of Internet measurement sampling in 2004, focusing on classical sampling methods, new applications and sampling methods, and applications areas. In 2007, Haddadi et al. (2008) revisited the issues of NetFlow sampling which focuses on data distortion and techniques for the compensation of data distortion.

Sampling methods, impact of sampling, integration of system-wide sampling, and recovering sampled data from distortion are mentioned in below studies. Duffield (2004) and Duffield et al. (2001) developed a size-dependent sampling scheme suitable for billing purposes. Estan et al. (2004) proposed an Adaptive NetFlow which dynamically adapts the sampling rate to achieve robustness without sacrificing accuracy. Brauckhoff et al. (2006) evaluated the impact of sampling on anomaly detection metrics using flows with the *Blaster* worm, and found that entropy-based features are less affected. Bartlet-ros and Cabellos-aparicio (2010) analyzed the impact of sampling on the accuracy of traffic classification using machine learning methods, and proposed a solution to reduce the impact. Cheng and Gong (2007) proposed a resource-efficient sampling system that combines three models: a pre-sampling model that records the estimated value rather than the measured value, a sampling and holding model that process the sampled packets to update the cache, and a non-uniform sampling model and keep the long flows in cache. Hao et al. (2007) developed a sampling scheme based on sampling two-runs to improve time and memory efficiency. Han et al. (2008) proposed a *pFlours* tool that fetches a packet and performs sampling to eliminate the synchronization problem during network traffic sampling. Duffield and Grossglauser (2008) discussed trajectory sampling, methods to eliminate duplications, and methods to join incomplete trajectories. Sekar et al. (2008) presented a system-wide approach that samples as a router primitive. To identify high-rate flow, Zhang et al. (2010) developed two methods: fixed sample size test which uses user specified accuracy, and truncated sequential probability test

through sequential sampling. Lee et al. (2011) proposed a method for related sampling where flows from the same application session are given higher probability. Bartos et al. (2011) proposed adaptive, feature-aware statistical sampling techniques to reduce the impact of sampling on anomaly detection.

4.7.3. Distributed analysis system

More applications demand real time analysis, advanced detection and classification. Centralized analysis systems face the difficulties of performance, scalability, and robustness. Although sampling provides an approach to reduce those burdens, there are tasks that cannot be based on sampling data. Distributed systems provide new mechanisms for capturing, accounting and monitoring (Morariu and Stiller, 2010). Several distributed analysis systems have been mentioned below. Kitatsuji and Yamazaki (2004) proposed a real-time system with a bit-pattern based flow definition and round-robin mechanism to balance packet streams. Sekar et al. (2008) proposed cSamp, a monitoring tool based on a coordinating mechanism for flow sampling, hash-based packet selection, and workload distributed. Morariu and Stiller (2011) proposed a distributed IP traffic analysis system. DiCAP (Morariu and Stiller, 2008), a flow capturing system, uses round-robin and a Distributed Hash Table (DHT) to distribute the workload and uses off-the-shelf hardware at network links. DIPStorage (Morariu et al., 2008) is a distributed flow storage platform for IP flow records based on DHT. SCRIPT (Morariu et al., 2010) is a distributed flow analysis framework that distributed flow records equally to multiple nodes. Others used peer-to-peer communication infrastructure (François et al., 2011; Gao et al., 2011) and map-reduce for efficient computation. Recent studies use the existing Hadoop based clustering platform and the map-reduce framework. Lee et al. (2011) proposed using Hadoop based map-reduce to process packet trace files. François et al. (2011) proposed botnet detecting system based on Hadoop based clustering and PageRank. Morken (2010) compared two map-reduce frameworks of Apache Hadoop and Nokia Disco, and concluded that Nokia Disco provides fast response time while Hadoop provides rich features, and map-reduce model is a very good approach for flow filtering and aggregation.

5. Discussion

Even though a large body of research has focused on traffic flows, many issues remain open. In particular, NetFlow data analysis is challenging because of the difficulty in collecting real data, huge datasets with limited information, and lack of systematic methodologies. In this section, we discuss our view about datasets, research perspectives, methodologies, challenges, and possible future research directions.

5.1. Datasets

Because of privacy and other concerns, researchers lack effective traffic flow datasets. Simulated data and other log data have been used as alternatives. Even though there is some real data, this data is either old or does not cover a large enough time period. Acquiring training datasets is another challenge for supervised machine learning. Moreover, there is no publicly available data for comparing different methodologies. Accurate analysis depends on real-time data collection. In surveyed papers, very few discuss a real time data collection solution (Gao et al., 2011).

Despite the popularity of sFlow and its wide deployment, few papers have focused on sFlow as their data source.

5.2. Research perspectives

Current studies have covered most perspectives of network monitoring, measurement, and network security. Application of network flow data in network monitoring is more successful than in network security, while real time network security is in high demand for network management. Basic top N data is not enough to understand the current complex network security situations. More specific perspectives such as referring user identity will provide clear information for security and forensic purpose. New perspectives will probably from network security because network security is becoming more important and challenge.

5.3. Methodologies

Heuristic approaches are easier to implement and seem more effective than machine learning approaches; however, practical experiences and findings are difficult to gain. Statistical approaches with heuristic methods give accurate results for known situations. For situations involving anomaly, more research is needed to develop advanced approaches that leverage information theory, machine learning and data mining. Much of the work has been limited to specific problems such as port-scan, DoS, or worms. A system that covers wide network security situations is needed for network security administrators. Moreover, visualization needs to focus more on IT operations and provides easy to understand and helpful information. For machine learning approaches, feature selection is a very important step that needs to be specific to the problem. Currently, there is no study available for understanding and comparing the effect of feature selection in the context of NetFlow data. Integrating data from other IT infrastructures will provide more information. As there is no publicly available dataset for comparing different approaches, researchers use their own private datasets in their experiments.

5.4. Challenges

With the constantly changing nature of networks, new applications and protocols being added to the Internet, network analysis will have to keep up with the speed of changes. For example, IPv6 addresses can be randomly generated and may not be identified as a unique host or user. Since IPv6 over IPv4 packets can bypass firewalls (Gregg et al., 2011), new approaches for IPv6 measurements are needed. New applications and protocols, faster Internet speeds with increased backbone bandwidth, and more complex content will make the analysis more difficult. In particular, the cloud computing that is based on moving contents to cloud services will make the analysis more complicated. In the following, we discuss specific challenges.

5.4.1. Feature representation and selection

Because NetFlow data only provides the header information, representing and selecting a set of appropriate features is challenging. For a specific task, the key question is how to effectively represent and extract features, and how to select the right features for a specific problem. With NetFlow version 9, it would be important to effectively leverage these new information.

5.4.2. Real time analysis and data storage

Analysis results need to be available in real time or within some fairly short period of time as the traffic is flowing. Furthermore, data needs to be continuously stored for certain amount of time for future need. Real time data collection is a challenging task because of the data size and the nature of the network traffic.

Real time analysis requires understanding the dynamic nature of network traffic. As Weinberger (2011) pointed out “that is the face of knowledge in the age of the Net: never fully settled, never fully written, never entirely done”.

5.5. Future directions

Despite significant work in the field, future research is needed to address the above mentioned challenges.

5.5.1. Distributed data collection and analysis

Real time analysis is in high demand in network security. Centralized analysis systems have difficulty dealing with huge data and real time analysis. Scalability and robustness are required to analyze data from multiple collectors. New technologies, such as *Apache Hadoop* related distributed data collection and analysis systems, open up more opportunities for re-thinking the network traffic analysis. Distributed applications and *map-reduce* model will provide more power and bring more insight and understanding.

5.5.2. Advanced analysis methodologies

Advanced methodologies using behavior-based features have the potential to mine helpful information. As Sommer and Paxson (2010) pointed out, machine learning algorithms excel at finding similarity rather than at identifying anomalous behaviors. To make machine learning approaches more accurate and efficient, there is a need for better understanding of different types of features and heuristics for specific goals. In practice, selecting and understanding an effective set of features is challenging and labor-intensive.

5.5.3. Integration

Integrating with existing network infrastructures (e.g., IDS, firewall and VPN gateway), integrating with log file event activities as well as integrating with host IDS (e.g., meta-events) all show a trend. NetFlow analysis can fill in the gap that IDS, firewall and host-based anti-virus tools cannot provide. It can provide monitoring, reporting, security altering, validating policy and configuration, assisting for forensic investigation, and serving as complimentary approaches for other network applications. Correlating with existing network infrastructures (e.g., NIDS may alert for an attack then NetFlow data will validate the alert) can give a high probability factor to remove false positives. Liu et al. (2007) proposed a method using Snort logs and NetFlow data fusion with SVMs to create network security awareness. Integrating together with other approaches (such as deep packet inspection), NetFlow-like approaches can provide a breadth-first approach for early investigation, and cover more hierarchies of network layers.

6. Conclusion

In this paper, we performed a comprehensive survey of network flow applications. First, we provided a brief background information on sFlow, NetFlow, and network traffic analysis. We covered the state of the art in network monitoring, analysis and management, application classification, user identity inferring, and network security awareness. We found that network security has been an important research topic, and has covered various aspects of network security issues. We then surveyed the state of the art of methodologies related to statistics, machine learning, profiling, behavior-based approaches, visualization, anonymization, and analysis systems. We found that advanced methodologies such as machine learning has been an important approach,

and applied mostly on application classification and network security awareness. Then, we critiqued the surveyed research with emphasis on datasets, research perspectives, methodologies, challenges, and pointed out possible directions for future research.

Acknowledgments

We are grateful to Dong Yu of Microsoft Research for valuable feedback. George Bebis is a Visiting Professor in the Department of Computer Science at King Saud University (KSU), Riyadh, Saudi Arabia.

References

- AURORA: Traffic analysis and visualization. < <http://www.zurich.ibm.com/aurora/> >. Retrieved September 13, 2012.
- Abdulla SA, Ramadass S, Altaher A, Nassiri AA. Setting a worm attack warning by using machine learning to classify netflow data. *International Journal of Computer Applications* 2011;36(December (2)):49–56.
- Ball R, Fink GA, North C. Home-centric visualization of network traffic for security administration. In: *Proceedings of the 2004 ACM workshop on visualization and data mining for computer security, VizSEC/DMSEC '04*. New York, NY, USA: ACM; 2004. p. 55–64.
- Barlet-ros P, Cabellos-aporicio A. Analysis of the impact of sampling on NetFlow traffic classification. *Methodology* 2010;55(5):1083–99.
- Barsamian AV. Network characterization for botnet detection using statistical-behavioral methods. Master's thesis, Thayer School of Engineering, Dartmouth College, USA; June 2009.
- Bartos K, Rehak M, Krmicek V. Optimizing flow sampling for network anomaly detection. In: *Wireless communications and mobile computing conference (IWCMC), 2011 7th international, July 2011*. p. 1304–9.
- Baker F, Foster B, Sharp C. Cisco architecture for lawful intercept in IP networks. < <http://rfc-ref.org/RFC-TEXTS/3924/index.html> >. Retrieved June 3, 2012.
- Bin L, Chuang L, Jian Q, Jianping H, Ungsunan P. A NetFlow based flow analysis and monitoring system in enterprise networks. *Computer Networks* 2008;52(5):1074–92.
- Bo X, Ming C, Fei L, Na W. P2P flows identification method based on listening port. In: *2nd IEEE international conference on broadband network multimedia technology, 2009. IC-BNMT '09, 2009*. p. 296–300.
- Bouhtou M, Klopfenstein O. Robust optimization for selecting netflow points of measurement in an IP network. In: *IEEE GLOBECOM 2007, IEEE global telecommunications conference, 2007*. p. 2581–5.
- Brauckhoff D, Tellenbach B, Wagner A, May M, Lakhina A. Impact of packet sampling on anomaly detection metrics. In: *Proceedings of the 6th ACM SIGCOMM on Internet measurement IMC 06, 2006*. p. 159.
- Caracas A, Kind A, Gantenbein D, Fussenegger S, Dechouniotis D. Mining semantic relations using NetFlow. In: *3rd IEEE/IFIP international workshop on business-driven IT management, 2008. BDIM 2008, April 2008*. p. 110–1.
- Carela-Espanol V, Barlet-Ros P, Solé-Pareta J. Traffic classification with sampled netflow. Technical Report 2, Technical report, Universitat Politècnica de Catalunya, 2009.
- Čeleda P, Vykojal J, Plesník T, Trunečka M, Krmíček V. Malware detection from the network perspective using netflow data. In: *3rd NMRG workshop on NetFlow/IPFIX usage in network management, 2010*.
- Chan Y-T, Shoniregun CA, Akmayeva GA. A NetFlow based Internet-worm detecting system in large network. In: *Third international conference on digital information management, 2008. ICDIM 2008, 2008*. p. 581–6.
- Chaudhary UK, Papapanagiotou I, Devetsikiotis M. Flow classification using clustering and association rule mining. In: *15th IEEE international workshop on computer aided modeling, analysis and design of communication links and networks (CAMAD), 2010*. p. 76–80.
- Chen Y, Jain S, Adhikari V, Zhang Z-L, Xu K. A first look at inter-data center traffic characteristics via yahoo! datasets. In: *INFOCOM, 2011 Proceedings IEEE, April 2011*. p. 1620–8.
- Cheng G, Gong J. A resource-efficient flow monitoring system. *Communications Letters, IEEE* 2007;11(June (6)):558–60.
- Cho K, Kaizaki R, Kato A. Aguri: an aggregation-based traffic profiler. In: *Proceedings of the second international workshop on quality of future internet services, COST 263*. London, UK: Springer-Verlag; 2001. p. 222–42.
- Choi H, Lee H, Kim H. Fast detection and visualization of network attacks on parallel coordinates. *Computers Security* 2009;28(5):276–88.
- Cohen E, Duffield N, Lund C, Thorup M. Confident estimation for multistage measurement sampling and aggregation. In: *Proceedings of the 2008 ACM SIGMETRICS international conference on measurement and modeling of computer systems SIGMETRICS 08, (i), 2008*. p. 109.
- Collins MP, Reiter MK. Hit-list worm detection and bot identification in large networks using protocol graphs. In: *Proceedings of the 10th international conference on recent advances in intrusion detection, RAID'07*. Berlin, Heidelberg: Springer-Verlag; 2007. p. 276–95.

- Deri L. ntop. <<http://www.ntop.org>>. Retrieved June 3, 2012.
- Deri L. Open source VoIP traffic monitoring. <<http://131.114.21.22/VoIP.pdf>>. Retrieved June 3, 2012.
- Dressler F, Jaegers W, German R. Flow-based worm detection using correlated honeypot logs. In: Proceedings of 15th GI/ITG fachtagung kommunikation in verteilten systemen (KIVS 2007), 2007. p. 181–6.
- Dubendorfer T, Plattner B. Host behaviour based early detection of worm outbreaks in internet backbones. In: 14th IEEE international workshops on enabling technologies infrastructure for collaborative enterprise WETICE05, (c), 2005. p. 166–71.
- Dubendorfer T, Wagner A, Plattner B. A framework for real-time worm attack detection and backbone monitoring. In: First IEEE international workshop on critical infrastructure protection IWCIP05, 2005. p. 3–12.
- Duffield N. Sampling for passive internet measurement: a review. *Statistical Science* 2004;19:472–98.
- Duffield N, Grossglauser M. Trajectory sampling with unreliable reporting. *IEEE/ACM Transactions on Networking* 2008;16(February (1)):37–50.
- Duffield N, Lund C, Thorup M. Charging from sampled network usage. In: Proceedings of the 1st ACM SIGCOMM workshop on internet measurement, IMW '01. New York, NY, USA: ACM; 2001. p. 245–56.
- Duffield N, Lund C, Thorup M. Properties and prediction of flow statistics from sampled packet streams. In: Proceedings of the 2nd ACM SIGCOMM workshop on Internet measurement, IMW '02. New York, NY, USA: ACM; 2002. p. 159–71.
- Erbacher RF. Visual behavior characterization for intrusion detection in large scale systems. In: Proceedings of the IASTED international conference on visualization, imaging, and image processing, 2001. p. 54–9.
- Estan C, Keys K, Moore D, Varghese G. Building a better NetFlow. *ACM SIGCOMM Computer Communication Review* 2004;34(4):245.
- Fioreze T, Granville LZ, Pras A, Sperotto A, Sadre R. Self-management of hybrid networks: Can we trust netflow data? In: IM09 IFIP/IEEE international symposium on integrated network management 2009, 2009. p. 577–84.
- Fischer F, Mansmann F, Keim DA, Pietzko S. Large-scale network monitoring for visual analysis of attacks. In: Visualization for computer security 5th international workshop VizSec 2008, 2008 proceedings, vol. 72(1–3), Cambridge, MA, USA, September 15, 2008. p. 1–8.
- Foukarakis M, Antoniadis D, Antonatos S, Markatos EP. Flexible and high-performance anonymization of NetFlow records using anontool. In: 2007 third international conference on security and privacy in communications networks and the workshops SecureComm, 2007. p. 33–8.
- François J, Wang S, State R, Engel T. BotTrack: tracking botnets using NetFlow and PageRank. *NETWORKING 2011* 2011;6640:1–14.
- François J, Wang S, Bronzi W, State R, Engel T. BotCloud: detecting botnets using MapReduce. In: 2011 IEEE international workshop on information forensics and security (WIFS), 2011. p. 1–6.
- Frias-Martinez V, Sherrick J, Stolfo SJ, Keromytis AD. A network access control mechanism based on behavior profiles. In: Computer security applications conference, 2009. ACSAC '09. Annual, 2009. p. 3–12.
- Galtsev AA, Sukhov AM. Network attack detection at flow level. *Aerospace*, 2011.
- Gao L, Yang J, Zhang H, Zhang B, Qin D. FlowInfra: a fault-resilient scalable infrastructure for network-wide flow level measurement. In: Network operations and management symposium (APNOMS), 2011 13th Asia-Pacific, 2011. p. 1–8.
- Gao Y, Li Z, Chen Y. A DoS resilient flow-level intrusion detection approach for high-speed networks. In: 26th IEEE international conference on distributed computing systems, 2006, ICDCS 2006, 2006. p. 39.
- Geer D. Behavior-based network security goes mainstream. *Computer* 2006; 39(March (3)):14–7.
- Goodall JR, Sowul M. VIAssist: visual analytics for cyber defense. In: IEEE conference on technologies for homeland security, 2009. HST '09, May 2009. p. 143–50.
- Gossett AM, Papapanagiotou I, Devetsikiotis M. An apparatus for P2P classification in Netflow traces. In: GLOBECOM workshops (GC Wkshps), 2010 IEEE, 2010. p. 1361–6.
- Gregg M, Matousek P, Sveda M, Podermanski T. Practical IPv6 monitoring-challenges and techniques. In: 2011 IFIP/IEEE international symposium on integrated network management (IM), May 2011. p. 650–3.
- Haddadi H, Landa R, Moore AW, Bhatti S, Rio M, Che X. Revisiting the issues on netflow sample and export performance. In: Third international conference on communications and networking in China, 2008. ChinaCom 2008. 2008. p. 442–6.
- Han B-J, Lee J-H, Sohn S-G, Ryu J-H, Chung T-M. pFlours: a new packet and flow gathering tool. In: 10th international conference on advanced communication technology, 2008. ICACT 2008, vol. 1, 2008. p. 731–6.
- Han S-h, Kim M-s, Ju H-t, Hong Jw-k. The architecture of NG-MON: a passive network monitoring system. In: Proceeding of 13th IFIP/IEEE international workshop on distributed systems: operations and management, 2002. p. 16–27.
- Hao F, Kodialam M, Lakshman TV, Mohanty S. Fast memory efficient flow rate estimation using runs. *IEEE/ACM Transactions on Networking* 2007;15(6): 1467–77.
- Hsiao H-W, Chen D-N, Wu TJ. Detecting hiding malicious website using network traffic mining approach. In: 2nd international conference on education technology and computer (ICETC), 2010, vol. 5, June 2010. p. V5-276–80.
- Hu Y, Chiu D-m, Lui JCS, Member S. Entropy based adaptive flow aggregation. *IEEE/ACM Transactions on Networking* 2009;17(3):698–711.
- Internet Traffic Classification. <<http://www.caida.org/research/traffic-analysis/classification-overview/>>. Retrieved June 3, 2012.
- Introduction to Cisco IOS ® NetFlow—a technical overview. <http://www.cisco.com/en/US/prod/collateral/iOSSwrel/ps6537/ps6555/ps6601/prod_white_paper0900aecd80406232.html>. Retrieved June 3, 2012.
- IPFIX. <<http://datacenter.ieff.org/wg/ipfix/>>. Retrieved September 13, 2012.
- Jiang H, Ge Z, Jin S, Wang J. Network prefix-level traffic profiling: characterizing, modeling, and evaluation. *Computer Networks* 2010;54(December (18)): 3327–40.
- Jiang H, Moore AW, Ge Z, Jin S, Wang J. Lightweight application classification for network management. In: Proceedings of the 2007 SIGCOMM workshop on Internet network management INM 07, 2007. p. 299.
- Jinsong W, Weiwei L, Yan Z, Tao L, Zilong W. P2P traffic identification based on NetFlow TCP flag. In: International conference on future computer and communication, 2009. ICFC 2009. April 2009. p. 700–3.
- Kalafut AJ, Van Der Merwe J, Gupta M. Communities of interest for internet traffic prioritization. In: IEEE INFOCOM workshops 2009, 2009. p. 1–6.
- Karagiannis T, Papagiannaki K, Faloutsos M. BLINC: multilevel traffic classification in the dark. *SIGCOMM Computer Communication Review* 2005;35(August (4)):229–40.
- Ke-xin Y, Jian-qi Z. A novel DoS detection mechanism. In: 2011 international conference on mechatronic science, electric engineering and computer (MEC), 2011. p. 296–8.
- Kerr DR, Bruins BL. Network flow switching and flow data export. US Patent Number US 6243667, 2001.
- Kim H, Claffy KC, Fomenkov M, Barman D, Faloutsos M, Lee K. Internet traffic classification demystified: myths, caveats, and the best practices. In: Proceedings of the 2008 ACM CoNEXT conference, CoNEXT '08, New York, NY, USA: ACM; 2008. p. 11:1–12.
- Kim M-S, Kong H-J, Hong S-C, Chung S-H, Hong JW. A flow-based method for abnormal network traffic detection. In: Network operations and management symposium, 2004. NOMS 2004. IEEE/IFIP, vol. 1, April 2004. p. 599–612.
- Kind A, Gantenbein D, Etoh H. 2006 Relationship discovery with netflow to enable business-driven IT management. In: The first IEEE/IFIP international workshop on business-driven IT management, 2006, BDIM '06, April 2006. p. 63–70.
- Kitatsujii Y, Yamazaki K. A distributed real-time tool for IP-flow measurement. In: 2004 international symposium on applications and the Internet, 2004. Proceedings, 2004. p. 91–8.
- Köandgel J. One-way delay measurement based on flow data: quantification and compensation of errors by exporter profiling. In: 2011 International conference on information networking (ICOIN), 2011. p. 25–30.
- Kobayashi A, Toyama K. Method of measuring VoIP traffic fluctuation with selective sFlow. In: 2007 International symposium on applications and the internet workshops. IEEE; 2007. p. 89.
- Kotsokalis C, Kalogeras D, Maglaris B. Router-based detection of DoS and DDoS attacks. In: Proceedings of HP OpenView university association HPOVUA 8th annual workshop, 2001.
- Krmicek V, Vykopal J, Krejci R. Netflow based system for NAT detection. In: Proceedings of the 5th international student workshop on emerging networking experiments and technologies. ACM; 2009. p. 23–4.
- Kundu SR, Pal S, Basu K, Das SK. An architectural framework for accurate characterization of network traffic. *IEEE Transactions on Parallel and Distributed Systems* 2009;20(1):111–23.
- Lakhina A, Crovella M, Diot C. Mining anomalies using traffic feature distributions. *SIGCOMM Computer Communication Review* 2005;35(August (4)):217–28.
- Lakhina A, Papagiannaki K, Crovella M, Diot C, Kolaczky ED, Taft N. Structural analysis of network traffic flows. *SIGMETRICS Performance Evaluation Review* 2004;32(June (1)):61–72.
- Lakkaraju K, Yurcik W, Lee AJ. NvisionIP: netflow visualizations of system state for security situational awareness. In: Proceedings of the 2004 ACM workshop on visualization and data mining for computer security, vol. 29, 2004. p. 65–72.
- Lau S. The spinning cube of potential doom. *Communications of the ACM* 2004;47(June (6)):25–6.
- Lee C-y, Kim H-k, Ko K-h, Kim J-w. A VoIP traffic monitoring system based on NetFlow v9. *International Journal of Advanced Science and Technology* 2009;4:1–8.
- Lee M, Hajjat M, Kompella RR, Rao S. RelSamp: preserving application structure in sampled flow measurements. In: INFOCOM, 2011 Proceedings IEEE, April 2011. p. 2354–62.
- Lee M, Member S, Duffield N. Opportunistic flow-level latency estimation using consistent NetFlow. *IEEE/ACM Transaction on Networking* 2011:1–14.
- Lee Y, Kang W, Lee Y. A hadoop-based packet trace processing tool. In: Proceedings of the third international conference on traffic monitoring and analysis, TMA'11. Berlin, Heidelberg: Springer-Verlag; 2011. p. 51–63.
- Li Y. Study of the monitoring model for securities trading network Quality of Service. In: 2nd international conference on information science and engineering (ICISE), 2010, 2010. p. 1–4.
- Liang C, Jian G. Fast application-level traffic classification using Net Flow records. *Journal on Communications* 2012;33(1):145–52.
- Liu D, Huebner F. Application profiling of IP traffic. In: 27th annual IEEE conference on local computer networks, 2002. Proceedings. LCN 2002, 2002. p. 220–9.
- Liu X-W, Wang H-Q, Liang Y, Lai J-B. Heterogeneous multi-sensor data fusion with multi-class support vector machines: creating network security situation awareness. In: 2007 international conference on machine learning and cybernetics, vol. 5, 2007. p. 2689–94.
- Mansman F, Meier L, Keim DA. Visualization of host behavior for network security. *Network Security* 2007:187–202.

- Mansmann F, Fischer F, Keim DA, Pietzko S, Waldvogel M. Interactive analysis of netflows for misuse detection in large IP networks. In: Müller P, Neumair B, Rodosek GD, editors. DFN-Forum Kommunikationstechnologien, LNI, vol. 149, GI, 2009. p. 115–24.
- McPherson J, Ma K-L, Krystosk P, Bartoletti T, Christensen M. PortVis: a tool for port-based detection of security events. In: Proceedings of the 2004 ACM workshop on visualization and data mining for computer security, VizSEC/DMSEC '04, ACM: New York, NY, USA; 2004. p. 73–81.
- Melnikov N, Schönwälder J. Cybermetrics: user identification through network flow analysis. In: Stiller B, De Turck F, editors. Mechanisms for autonomous management of networks and services. Lecture notes in computer science, vol. 6155. Berlin, Heidelberg: Springer; 2010. p. 167–70.
- Minarik P, Dymacek T. NetFlow data visualization based on graphs. Visualization for computer security, 2008. p. 144–51.
- Minarik P, Vykopal J, Krmicek V. Improving host profiling with bidirectional flows. In: International conference on computational science and engineering, 2009. CSE '09, vol. 3, 2009. p. 231–7.
- Moghaddam S, Helmy A. SPIRIT: a simulation paradigm for realistic design of mature mobile societies. In: Wireless communications and mobile computing conference (IWCMC), 2011 7th international, July 2011. p. 232–7.
- Morariu C, Kramis T, Stiller B. DIPStorage: distributed storage of IP flow records. In: 16th IEEE workshop on local and metropolitan area networks, 2008. LANMAN 2008. 2008. p. 108–13.
- Morariu C, Racz P, Stiller B. SCRIPT: a framework for scalable real-time ip flow record analysis. In: Network operations and management symposium (NOMS), 2010 IEEE, April 2010. p. 278–85.
- Morariu C, Stiller B. DiCAP: distributed packet capturing architecture for high-speed network links. In: 33rd IEEE conference on local computer networks, 2008. LCN 2008, 2008. p. 168–75.
- Morariu C, Stiller B. Distributed architecture for real-time traffic analysis. In: Proceedings of the mechanisms for autonomous management of networks and services, and 4th international conference on autonomous infrastructure, management and security, AIMS'10. Springer-Verlag: Berlin, Heidelberg; 2010. p. 171–4.
- Morariu C, Stiller B. An open architecture for distributed IP traffic analysis (DITA). In: 2011 IFIP/IEEE international Symposium on integrated network management (IM), May 2011. p. 952–7.
- Morken JT. Distributed NetFlow processing using the map-reduce model. Master's thesis, Norwegian University of Science and Technology; 2010.
- Nagaraj K, Naidu KVM, Rastogi R, Satkin S. Efficient aggregate computation over data streams. In: IEEE 24th international conference on data engineering, 2008. ICDE 2008. April 2008. p. 1382–4.
- NetFlow applications <<http://netflow.caligare.com/applications.htm>>. Retrieved September 13, 2012.
- NFSen—Netflow Sensor. <<http://nfsen.sourceforge.net/>>. Retrieved September 13, 2012.
- Nguyen TTT, Armitage G. A survey of techniques for Internet traffic classification using machine learning. Communications Surveys Tutorials, IEEE 2008;10(4): 56–76.
- Nor SM, Mohd AB. Towards a flow-based internet traffic classification for bandwidth optimization. International Journal of Computer Science and Security 2009;3(3):146.
- Oberheide J, Goff M, Karir M. Flamingo: visualizing internet traffic. In: Network operations and management symposium, 2006. NOMS 2006. 10th IEEE/IFIP, April 2006. p. 150–61.
- Oslebo A. Stager a web based application for presenting network statistics. In: Network operations and management symposium, 2006. NOMS 2006. 10th IEEE/IFIP, April 2006. p. 1–15.
- Patwari N, Hero AO, Pacholski A. Manifold learning visualization of network traffic data. In: Proceeding of the 2005 ACM SIGCOMM workshop on mining network data MineNet '05, 2005. p. 191.
- Perelman V, Melnikov N, Schönwälder J. Flow signatures of popular applications. In: Agoulmine N, Bartolini C, Pfeifer T, O'Sullivan D, editors. Integrated network management, IEEE; 2011. p. 9–16.
- Plonka D. FlowScan: a network traffic flow reporting and visualization tool. In: Proceedings of the 14th USENIX conference on system administration. Berkeley, CA, USA: USENIX Association; 2000. p. 305–18.
- Proto A, Alexandre LA, Batista ML, Oliveira IL, Cansian AM. Statistical model applied to netflow for network intrusion detection. Transactions on Computational Science 2010;11:179–91.
- Qun W, Xiyue D, Lu H. Novelty P2P flow analysis system. In: 7th International conference on wireless communications, networking and mobile computing (WiCOM), 2011, 2011. p. 1–4.
- Rehak M, Pechoucek M, Celeda P, Krmicek V, Grill M, Bartos K. Multi-agent approach to network intrusion detection. In: Proceedings of the 7th international joint conference on autonomous agents and multiagent systems demo papers. International foundation for autonomous agents and multiagent systems, 2008. p. 1695–6.
- Rehak M, Pechoucek M, Minarik P. Collaborative attack detection in high-speed networks. Analysis. Lecture notes in artificial intelligence, vol. 4696. 2007. p. 73–82.
- Ren P, Gao Y, Li Z, Chen Y, Watson B. IDGraphs: intrusion detection and analysis using stream compositing. Computer Graphics and Applications IEEE 2006;26(2):28–39.
- Ricciato F, Strohmeier F, Dorfinger P, Coluccia A. One-way loss measurements from IPFIX records. In: 2011 IEEE international workshop on measurements and networking proceedings (M N), 2011. p. 158–63.
- Rohmad MS, Azmat F, Manaf M, Manan J-I. Enhanced Netflow version 9 (e-Netflow v9) for network mediation: structure, experiment and analysis. In: International symposium on information technology, 2008. ITSIM 2008, vol. 3, 2008. p. 1–6.
- Rossi D, Valenti S. Fine-grained traffic classification with netflow data. In: Proceedings of the 6th international wireless communications and mobile computing conference on ZZZ IWCMC 10, 2010. p. 479.
- sFlow Collectors. <<http://www.sflow.org/products/collectors.php>>. Retrieved September 13, 2012.
- Sawaya Y, Kubota A, Miyake Y. Detection of attackers in services using anomalous host behavior based on traffic flow statistics. In: 2011 IEEE/IPSJ 11th international symposium on applications and the internet (SAINT), July 2011. p. 353–9.
- Schatzmann D, Leinen S, Kögel J, Mühlbauer W. FACT: flow-based approach for connectivity tracking. In: Spring N, Riley GF, editors. PAM, Lecture notes in computer science, vol. 6579. Springer; 2011. p. 214–23.
- Schatzmann D, Mühlbauer W, Spyropoulos T, Dimitropoulos X. Digging into HTTPS: flow-based classification of webmail traffic. In: Proceedings of the 10th annual conference on internet measurement, IMC '10. New York, NY, USA: ACM; 2010. p. 322–7.
- Sekar V, Reiter MK, Willinger W, Zhang H, Kompella RR, Andersen DG. cSamp: a system for network-wide flow monitoring. In: Proceedings of the 5th {USENIX} {NSDI}, San Francisco, {CA}, April 2008.
- Shelley DS, Gunes MH. GerbilSphere: inner sphere network visualization. Computer Networks 2012;56(3):1016–28.
- Shen W, Chen Y, Zhang Q, Chen Y, Deng B, Li X, et al. Observations of IPv6 traffic. In: ISECS international colloquium on computing, communication, control, and management, 2009. CCCM 2009, vol. 2, 2009. p. 278–82.
- Singh MP, Subramanian N, Rajamenakshi R. Visualization of flow data based on clustering technique for identifying network anomalies. In: IEEE symposium on industrial electronics applications, 2009. ISIEA 2009, vol. 2, 2009. p. 973–8.
- Sinha A, Mitchell K, Medhi D. Flow-level upstream traffic behavior in broadband access networks: DSL versus broadband fixed wireless. In: 3rd IEEE workshop on IP operations and management, 2003. (IPOM 2003), 2003. p. 135–41.
- Slagell AJ, Luo K. Sharing network logs for computer forensics: a new tool for the anonymization of netflow records. In: Workshop of the 1st international conference on security and privacy for emerging areas in communication networks 2005, 2005. p. 37–42.
- Sommer R, Paxson V. Outside the closed world: on using machine learning for network intrusion detection. In: 2010 IEEE symposium on security and privacy, 2010. p. 305–16.
- Soysal M, Schmidt EG. Machine learning algorithms for accurate flow-based network traffic classification: evaluation and comparison. Performance Evaluation 2010;67(June (6)):451–67.
- Sperotto A, Pras A. Flow-based intrusion detection. In: Agoulmine N, Bartolini C, Pfeifer T, O'Sullivan D, editors. Integrated network management. IEEE; 2011. p. 958–63.
- Sperotto A, Schaffrath G, Sadre R, Morariu C, Pras A, Stiller B. An overview of ip flow-based intrusion detection. Communications Surveys Tutorials, IEEE 2010;12(3):343–56.
- Strasburg C, Krishnan S, Dorman K, Basu S, Wong JS. Masquerade detection in network environments. In: 10th IEEE/IPSJ international symposium on applications and the Internet (SAINT), 2010, July 2010. p. 38–44.
- Strohmeier F, Dorfinger P, Trammell B. Network performance evaluation based on flow data. In: Wireless communications and mobile computing conference (IWCMC), 2011 7th International, July 2011. p. 1585–9.
- Sukhov AM, Sidelnikov DI, Galtsev AA, Platonov AP, Strizhov MV. Active flows in diagnostic of troubleshooting on backbone links. CoRR, abs/0911.2, 2009.
- Taylor T, Brooks S, McHugh J. NetBytes viewer: an entity-based NetFlow visualization utility for identifying intrusive behavior. VizSEC 2007, 2008. p. 101–14.
- Taylor T, Paterson D, Glanfield J, Gates C, Brooks S, McHugh J. FloVis: flow visualization system. In: Conference for homeland security, 2009. CATCH '09. Cybersecurity applications technology, March 2009. p. 186–98.
- Trammell B, Tellenbach B, Schatzmann D, Burkhart M. Peeling away timing error in netflow data. In: Spring N, Riley GF, editors. PAM, Lecture notes in computer science, vol. 6579. Springer; 2011. p. 194–203.
- Truong P, Guillemin F. A heuristic method of finding heavy hitter prefix pairs in IP traffic. Communications Letters, IEEE 2009;13(October (10)):803–5.
- Valenti S, Rossi D. Identifying key features for P2P traffic classification. In: 2011 IEEE international conference on communications, ICC. IEEE; 2011. p. 1–6.
- Vliek G. Detecting spam machines, a netflow-data based approach. February 2009.
- Vykopal J, Plesnik T, Minarik P. Network-based dictionary attack detection. In: 2009 international conference on future networks, March 2009. p. 23–7.
- Wagner A, Dubendorfer T, Hammerle L, Plattner B. Flow-based identification of P2P heavy-hitters. International conference on internet surveillance and protection, 00(c), 2006. p. 15.
- Wagner C, Francois J, State R, Engel T. Machine learning approach for IP-flow record anomaly detection. In: Proceedings of the 10th international IFIP TC 6 conference on Networking—volume part I, NETWORKING'11. Berlin, Heidelberg: Springer-Verlag; 2011. p. 28–39.
- Wagner C, Francois J, State R, Engel T. DANAK: finding the odd! In: 5th International conference on network and system security (NSS), 2011, 2011. p. 161–8.
- Wagner C, Wagener G, State R, Engel T, Dulaunoy A. Game theory driven monitoring of spatial-aggregated IP-Flow records. In: 2010 International conference on network and service management (CNSM), 2010. p. 463–8.

- Wang S, Guo R. GA-based filtering algorithm to defend against DDoS attack in high speed network. In: Fourth international conference on natural computation, 2008. ICNC '08. vol. 1, 2008. p. 601–7.
- Wei S, Mirkovic J, Kissel E. Profiling and clustering internet hosts. In: DMN'06, 2006. p. 269–75.
- Weinberger D. The machine that would predict the future. *Scientific American* 2011;305(6):52–7.
- Weststrate H. Botnet detection using netflow information finding new botnets based on client connections. In: *Structure*, 2009.
- Winter P, Hermann E, Zeilinger M. Inductive intrusion detection in flow-based network data using one-class support vector machines. In: 4th IFIP international conference on new technologies, mobility and security (NTMS), 2011, 2011. p. 1–5.
- Xu B, Chen M, Hu C. DEAPFI: a distributed extensible architecture for P2P flows identification. In: IEEE international conference on network infrastructure and digital content, 2009. IC-NIDC 2009. 2009. p. 59–64.
- Xu K, Wang F, Gu L. Network-aware behavior clustering of Internet end hosts. In: INFOCOM, 2011 Proceedings IEEE, April 2011. p. 2078–86.
- Xu K, Zhang Z-L, Bhattacharyya S. Profiling Internet backbone traffic: behavior models and applications. *SIGCOMM Computer Communication Review* 2005; 35(August (4)):169–80.
- Yin K, Nianqing T. Study on the risk detection about network security based on grey theory. In: International forum on information technology and applications, 2009. IFITA '09, vol. 1, May 2009. p. 411–3.
- Yin X, Yurcik W, Treaster M, Li Y, Lakkaraju K. VisFlowConnect: netflow visualizations of link relationships for security situational awareness. In: Proceedings of the 2004 ACM workshop on visualization and data mining for computer security, VizSEC/DMSEC '04. ACM; 2004. p. 26–34.
- Zadnik M, Pecenka T, Korenek J. Netflow probe intended for high-speed networks. In: International conference on field programmable logic and applications, 2005. 2005. p. 695–8.
- Zeng Y, Hu X, Shin KG. Detection of botnets using combined host- and network-level information. In: *Symposium a quarterly journal in modern foreign literatures*, 2010. p. 291–300.
- Zha W, He J. On campus network P2P and its link control. In: 2011 international conference on consumer electronics, communications and networks (CECNet), April 2011. p. 5086–9.
- Zhang H. Study on the TOPN abnormal detection based on the netflow data set. *Computer and Information Science* 2009;2(3):103–8.
- Zhang J, Meng S. A design of NetFlow traffic statistic and analysis system for process of the transition of commercialization of IPV6. In: 2011 International conference on computer science and service system (CSSS), June 2011. p. 963–5.
- Zhang YB, Fang BB, Luo H. Identifying high-rate flows based on sequential sampling. *leice Transactions on Information and Systems* 2010;E93-D(5):1162–74.
- Zhenqi W, Xinyu W. NetFlow based intrusion detection system. In: 2008 International conference on multimedia and information technology, 2008. p. 825–8.
- Zhu H, Zhang X, Ding W. Research on errors of utilized bandwidth measured by netflow. In: 2011 Second international conference on networking and distributed computing (ICNDC), 2011. p. 45–9.
- Zhu Z, Lu G, Chen Y, Fu ZJ, Roberts P, Han K. Botnet research survey. In: *Computer software and applications*, 2008. COMPSAC '08. 32nd annual IEEE international, 2008. p. 967–72.