

# NTCS: A Real Time Flow-based Network Traffic Classification System

Silas Santiago Lopes Pereira, Jorge Luiz de Castro e Silva, José Everardo Bessa Maia

Department of Statistics and Computing

UECE - State University of Ceará

Fortaleza - Ceará - Brazil

Email: {silas, jlcs}@larces.uece.br, jose.maia@uece.br

**Abstract**—The continual growth of high speed networks is a challenge for real-time network analysis systems. The real time traffic classification is an issue for corporations and ISPs (Internet Service Providers). This work presents the design and implementation of a real time flow-based network traffic classification system. The classifier monitor acts as a pipeline consisting of three modules: packet capture and preprocessing, flow reassembly, and classification with Machine Learning (ML). The modules are built as concurrent processes with well defined data interfaces between them so that any module can be improved and updated independently. In this pipeline, the flow reassembly function becomes the bottleneck of the performance. In this implementation, was used a efficient method of reassembly which results in a average delivery delay of 0.49 seconds, approximately. For the classification module, the performances of the K-Nearest Neighbor (KNN), C4.5 Decision Tree, Naive Bayes (NB), Flexible Naive Bayes (FNB) and AdaBoost Ensemble Learning Algorithm are compared in order to validate our approach.

## I. INTRODUCTION

The Internet traffic is changing continuously and this contribute to difficult the characterization of network behavior and structure. Massive games [1] and cloud and grid services increase every day their percentage participation in total network traffic. Traffic monitoring Systems generally make use of flow information. Examples are NetFlow [2] or IETF IPFIX [3], which defines a standard to exporting flow information by routers and switches. Such systems are widely used in network service providers and corporations to gain knowledge about critical business applications, analyze communication patterns prevalent in traffic, collect data for account, or detect anomalous traffic patterns [4]. A vital issue for corporations and ISPs (Internet Service Providers) is to identify traffic application types which are transmitted on their networks [5].

Pattern recognition and machine learning models have given significant attention to semi-supervised learning [6]. In network traffic areas, encryption and processing restrictions, protocol obfuscation and use of ephemeral ports make the task of construct classification models difficult. The large amount of Internet traffic flowing through networks makes the use of approaches that combine labeled and unlabeled data to construct accurate classifiers suitable.

There are a large number of papers in the traffic monitoring and traffic classification area. Most papers usually focus on either traffic flow reassembly or traffic classification and identification, but not on their combination. This paper describes

the architecture of a real time Internet traffic classifier monitor for use in corporate networks. It also evaluates different machine learning methods for network traffic classification. The classifier monitor is based on concept of bidirectional flow. This means that the fundamental object to be classified in a determined pattern is the traffic flow, either complete or as subflow. A flow is defined by one or more packets between a host pair with the same quintuple: source and destination IP address, source and destination ports and protocol type (ICMP, TCP, UDP) [9].

The remainder of this paper is organized as follows. Section II overviews the related work about flow reassembly and traffic classification. In section III, we describe the design and implementation of the classifier monitor. Section IV details the data collection used for evaluate the NTCS and describes how the experiments were performed. Section V presents and discusses the performance tests results. Section VI ends with some conclusions and future work.

## II. RELATED WORK

Statistical classification is based on collecting statistical information on properties of the traffic flow, and relies on the assumption that each category has a particular distribution of properties which represents it and can be used to identify it [10]. The statistical traffic classification using machine learning techniques have been widely explored in the recent years.

There are a limited number of tools available in the literature for traffic classification [7]. The NetAI tool is able to perform online and offline feature extraction, although not directly perform traffic classification. The FullStats is able to extract an extensive set of characteristics, but from a offline trace. The GTVS, which is a DPI (*Deep Packet Inspection*) based software, allows the the labeling of traces in a semi-automated manner. The only two traffic classification tools which implement machine learning method are Tstat 2.0 and TIE. The Tstat makes use of the *packet size* and *interpacket time* features in a Bayesian framework for identification of Skype and obfuscated P2P file sharing. Although the tool has a limited number of applications, it can extract a large number of characteristics. The TIE software platform is available to the research community, and which allows the development of classification methods. The framework provides traffic capture

and processing, feature extraction, and online classification. At the moment, a few number of features is available on TIE. Systems like Bro, which is able to collect flow statistics and perform payload-based classification at high speed rates, are limited when the traffic is encrypted [8]. Here, we briefly review some important approaches to stream reassembly (subsection II-A) and traffic classification (subsection II-B).

#### A. Stream Reassembly

In [11], the authors present an efficient TCP stream reassembly mechanism for real time network traffic processing at high speeds. The mechanism uses the recently-accessed-first principle to reduce the search cost of a connection for each packet arrival. Moreover, to improve the search process, the system keeps established and not established TCP connections in different structures. Experimental results based on network traffic captured in a typical gigabit gateway showed the proposed policy, in comparison of traditional one (RFC 793 [12]), was more efficient and could attend the real time property requisite of traffic analysis systems in gigabit networks.

In [13], the author presents a TCP stream reassembly mechanism designed and implemented to an network-based intrusion detection system. The system receives individual packets from network and performs signature detection from the payload. The approach is described as follows: First, the system associates each received packet to its corresponding TCP connection, based on the quadruple composed of source IP address and port, and destination IP address and port. Then, the system checks the packet sequence number and determines if this packet is the next expected packet for the respective connection. If true, the packet is sent for signature detection.

In [14], to improve the forensic analysis, the authors present a TCP stream evaluation methodology which consists of estimate the TCP reassembly accuracy by the precise identification of potential errors at packet and stream levels hidden in this process. This approach can be used for derivation and computation of reassembly errors. In the proposed TCP reassembly model, a session counting algorithm is presented, which defines a flow as a set of TCP packets with same values to source IP and port, and destination IP and port, and a flow can have multiples sessions delimited by well defined phases of connection establishment and termination. From two traffic captures obtained with *Tcpdump* tool, the authors used a libpcap-based program [15] to read the packet traces and evaluate the reassembly and error verification approaches, which was experimentally validated with known analysis tools as *Tcptrace* and *Tcpflow*.

#### B. Traffic Classification

In [16], it is evaluated the effectiveness of machine learning techniques for the real time traffic classification problem using statistical attributes derived from first packets of each flow. The authors utilized port-base method for labeling flows into categories. Since the traffic traces are anonymized by privacy reasons, it precludes the inference of applications that generated the flows. Although this approach can consequently

introduce incorrectly labeled flows, the authors argue that, for the studied ports, the non-labeled flow instances percentage is low and most of the traffic belongs to standard applications. The acquired results showed that the classification with decision trees had the highest accuracy and performance in comparison of other classifiers. In addition, subflow-based classifiers can reach high accuracy values while the computational complexity is reduced.

In [17], the machine learning approach applied to traffic classification using only transport-layer statistics is explored. This approach seeks to circumvent the problem that many network applications, such as P2P protocols, make use of dynamic port numbers and content encryption to avoid detection. This becomes inefficient the traditional approaches of port mapping and content analysis. The author evaluate the impact on the performance of data dimensionality, selected attributes, and machine learning used algorithms, which were, respectively, TAN, C4.5, NBTree, RandomForest and Distance Weighted KNN. The application of classification techniques and discriminant selection based on genetic algorithms together can dramatically reduce the learning and modeling time with little variation in the classification process accuracy.

Bar-Yanai et al. (2010) [10] proposes a hybrid combination of k-means and k-nearest neighbor geometrical classifiers. The proposed statistical classifier was integrated into a real-time embedded environment implemented on a SCE2020 Cisco platform and the algorithm was tested at a full line rate. The average accuracy rate of K-Means and k-nearest neighbor algorithm is of 83% and 99.1%, respectively, and the accuracy of the hybrid approach is very similar to the k-NN technique, with marginal accuracy loss. The time complexities of the hybrid solution are much lower than those of the k-NN algorithm, which has the better accuracy but higher complexity. The authors conclude that their proposed algorithm works well on platforms with limited memory, CPU resources and real-time requests. The technique works as follows: Initially, the training dataset is divided in N smaller datasets, each one corresponding to one of the known application categories presented in the training data. After this, each partition is grouped with the K-Means clustering method for a certain number of K clusters, so that each cluster is labeled naturally by the application protocol in which it was constructed. The second training stage consists in redistribution of the entire sample set of cluster centers defined in the previous stage. The online classification is performed by the execution of the KNN ( $k = 1$ ) classification method over the members of the closest cluster to an unseen flow.

Mallapragada et al. (2009) [6] propose a boosting framework for binary classification which combines advantages of graph-based approaches and ensemble methods. The strategy is to improve the performance of a given supervised learning algorithm using unlabeled data. The proposed algorithm is named *SemiBoost*, and is a general and efficient approach that allows the choice of a base classifier which is well-suited to a specific task. Like other boosting algorithms, the classification accuracy is improved iteratively, but *SemiBoost* selects

unlabeled data along the iterations. The proposed algorithm combines similarity information with classifier predictions to obtain the most confident pseudo-labels. The authors used WEKA [18] software in the implementation of benchmark semi-supervised approaches. From the evaluation and comparison of the method with three state-of-art semi-supervised algorithms (TSVM, LDS and LavSVM) and in 16 different datasets indicate significant gain of the proposed approach, when evaluated with the base classifiers Decision Stump, J48 and SVM.

Szabó et al. (2011) [19] proposes a framework for traffic classification which uses machine learning techniques and only information from the packet headers. One component of the proposed framework is the combination of classification and clustering algorithms to make the identification system robust under different network conditions. The training and evaluation of classification system were performed with the traffic flow obtained from measurements in networks with different access technologies and different locations, in order to make the traffic characteristics varied. The authors found that these clustering and classification methods resulted with different performance results when used to identify traffic from unknown networks. They also verified that clustering algorithms have proved to be more robust with network parameter changes while classification methods can learn about a specific network more accurately. The authors present and evaluate two different combinations of classification and clustering approaches that result in accuracy increase when comparing to standalone cases. The first combination, named classification with clustering information, means that each training flow has its respective cluster number as a new attribute for supervised classification. The cluster numbers of training flows are obtained from the previous clustering of training data. Since the clustering information attribute may be neglected or considered with low-importance by a supervised technique, this approach cannot always improve the global accuracy. The second approach, named model refinement with per cluster based classification, initially applies unsupervised learning to generate clusters. A separate classification model is then built for the set of flows of each cluster. In the evaluation phase at an unseen flow, the unsupervised method results with the number for the most similar cluster for which the associated model is used to evaluate the flow. This approach always considers the clustering results with high importance and the supervised techniques can build simple models since each group contains a limited number of flow types. This implies that the impact of over-fitting with the classification model is reduced. The per cluster based classification scheme outperforms the standalone supervised and unsupervised methods. The proposed methods obtained 93% and 75% of TP ratio for the evaluation on the same network and the cross-checks on other networks, respectively.

Erman et al. (2007) [8] proposes and evaluates a semi-supervised methodology for statistical traffic classification that is able to accommodate known and unknown applications. The authors mention three main advantages to their proposed semi-

supervised method: Initially, fast and accurate classifiers can be designed since a training dataset that consists with a small number of labeled flows and a large number of unlabeled flows. Second, the approach is robust and can handle both previously unknown applications and new patterns of existing applications. Moreover, network operators can insert unlabeled flows to improve the classifier performance, allowing the iterative development. Finally, the proposed approach can be integrated with solutions that collect flow statistics. The semi-supervised model combines supervised and unsupervised methods in two steps: The approach first uses the K-Means clustering technique to partition a training dataset composed of a few labeled flows and abundant unlabeled flows. The second step utilizes the available labeled flows to build a cluster with application mapping for which clusters without labeled flows remain unmapped, which corresponds to flows that possibly do not belong to any known application. The authors found that the proposed model was able to identify a variety of different applications with a high rate of accuracy, such as Web, P2P, FTP and E-mail. The flow and byte accuracy was above 98% and 93%, respectively. Furthermore, the authors noted that datasets with large number of flows consistently achieve a high classification of accuracy. The authors verify that despite labeling tools, labeling a large dataset can be expensive and difficult. In practice, labeling only a fraction of the training flows is sufficient to obtain high levels of accuracy.

### III. THE CLASSIFIER MONITOR

This section details the design and operation of our classifier monitor, and concludes with the presentation of the modules which compose the system.

#### A. Architecture

The monitor works as a three-stage pipeline, with a collect and preprocessing module, a flow reassembly module, and an attribute extraction and classification module. For the purpose of pipeline, the time is divided in intervals of 30s. This value was chosen arbitrarily. Once the monitor starts, three parallel processes are in execution on each interval: the packet capture, the flow reassembly of the previous interval packet capture, and the flow classification for the collection occurred in two delay intervals. Another parallel process is responsible for closing old connections periodically, in order to reduce the use of memory and processing during reassembly process. This approach allows the classifier monitor reach a response time of  $(30 + \alpha)$  seconds, where  $\alpha$  is the necessary time to performs the reassembly of the captured data in a given interval. In summary, the monitor works with a *quantum* of 30s of traffic and with an average delay of  $\alpha$  seconds in the flow reassembly, feature extraction, and classification. The average found value achieved in the current implementation was  $\alpha = 0.49$  seconds.

The Figure 1 exhibits the capturing and processing environment of the monitoring and classification system. Basically, we assume that the traffic is mirrored by a network border router to a network interface monitored by the system. The system periodically performs the processing and categorization

of captured data, and presents the obtained information from monitoring and classification process.

The Figure 2 exhibits the layered structure of implemented classifier monitor, whose tasks modules are online traffic collection from a network point, preprocessing for flow re-assembly, extraction and selection of statistical attributes, flow labeling since payload analysis or port-based method (only during training step), training with a supervised machine learning technique and classification, using the ML model built from training data. The classifier monitor performs the packet traffic capture continuously. In the training phase, the captured packets are sent to a reassembly process, which associate each packet to its respective flow. A parallel process extracts statistical information from the packet headers, selects the most relevant attributes using an attribute selection algorithm, and labels the flows with well-known ports method [20]. The traffic flows, which are disposed in a spatial representation (each flow is an instance with a set of characteristics), are used to training a selected supervised classification method. In the evaluation phase, the unlabeled flows obtained with the collection, reassembly, and attribute extraction are finally evaluated by the classifier.

Since the modules are implemented as concurrent processes, the data interfaces between each module allows that each one can be improved and updated independently. To the packet capture and preprocessing, the input and output are respectively the TCP packet and a data structure with only the necessary packet information to the monitor (e.g. bytes, timestamp, flags, no payload). For the flow reassembly module, the input and output are respectively a set of pre-processed packets and the data structure which represents the reconstructed flow. For the classification module, the features vector is the input, and the output is the class.

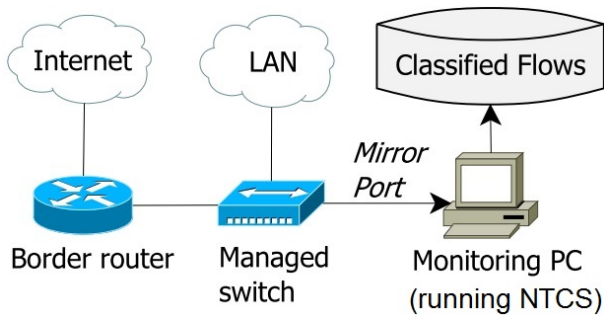


Fig. 1. Traffic Capture Environment .

The monitor was developed in C#.Net programming language using Visual Studio Integrated Development Environment. The online packet capture is performed based on sequential reading and processing of each packet contained from a network interface. The monitor also adopts a given timeout for packet capture and presentation of results. The reason for implementation of a flow reassembly algorithm, despite given the existence of several tools and libraries to reach this goal, as libNIDS [21], TcpTrace [22], and WireShark [23],

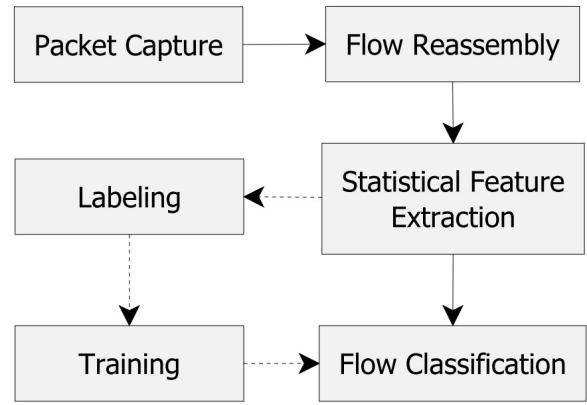


Fig. 2. Block Diagram of Classifier Monitor.

is the possibility to evaluate different approaches for subflow classification, as mentioned in [24]. Moreover, the evaluation of approaches for real time TCP stream reassembly becomes possible, as in [25] and [26], which are fundamental in the development of a high-speed traffic classification system.

### B. Packet capture and Preprocessing

The traffic capture at packet level, following by data processing and visualization, is a common demand on traffic volume monitoring or traffic supervision tasks. Therefore, it is essential that packet capture process must capture completed packets so that the stream is reassembled correctly

For the task of packet capture, we use "TCP Session Reconstruct Tool" [27], which is a C# utility for packet capture and reconstruction of complete and incomplete TCP sessions. This tool is available with CPOL license, and is based on libnids library [21] and Wireshark. It uses a reconstruction TCP algorithm named TcpRecon. TcpRecon reconstructs a bidirectional flow, holds each flow in a dictionary structure and recovers its payload. We reuse this software by replacing the TcpRecon policy by our proposed reassembly scheme. The TCP flow Reassembly is explained in the next subsection.

### C. Flow Reassembly

A reassembly function associates a TCP packet with its respective stream. The purpose of this function is to recover the initial state, emitted by sender, from the captured TCP packets [14]. It is essential that reassembly, which is applicable to a diversity of network traffic analysis systems, such as intrusion detection and prevention, content inspection and network forensics, can be performed as fast as possible to handle high traffic rates, especially in high-speed networks [11].

Although the RFC 973 [14] presents a standard specification of TCP protocol (the TCP protocol is specified in a large number of RFCs), there are different implementations, which make TCP reassembly a difficult task. Different reassembly tools hold their own stream concept specifications. For example, the Tcpflow tool binds a tuple with a stream, while the Tcptrace

tool associates a session with a given stream. The *Tcptrace* and *Tcpflow* tools group the data sent in each direction in different streams. In a stream generated by the *Wireshark* tool, data from sender and receiver are grouped in the same stream.

*Recently-accessed-first Principle:* In a high-speed network, there may be hundreds of thousands of simultaneous connections, so that the reassembly system, which maintains a structure for storing connection records, searches in this structure the corresponding record for each collected packet. The search becomes expensive and needs to be optimized as the number of connections increases [11]. The idea of recently-accessed-first principle is to bring the most recently accessed connection records to top of connection record list.

Since in a TCP connection the data transfer obeys the TCP/IP specification, the locality principle applied to the arrival of packets on network is based on assumptions that, given a packet and its subsequent, it is presumable both belong to the same connection, and considering the packet arrival of a certain connection, the next packet of this connection will come soon [11]. From the described principle about TCP packet arrival, the recently-access first principle is adopted in search of connection record set. For every successful search, the accessed record is moved to top of the set, so that after a sequence of access, the most frequently accessed nodes will be in the beginning of set. Thus, they can be accessed more quickly and the search efficiency is improved [11].

This principle, although efficient for clients, present bad performance on server or links with a lot of traffic flows. This is because the locality principle is lost in these scenarios [28]. However, a concurrent process in our monitor finishes old flows periodically in order to circumvent this issue.

This work uses the same concept of TCP stream presented in [14], and the recently-access-first principle presented in [11] to optimize the reassembly step of our software-based solution. Differently of [11], which uses two hash tables for connection management, we use a simple list to store connections. The adopted reassembly policy was based on mechanism proposed in [13] for TCP session reassembly. The applied reassembly approach was validated experimentally with the *Tcptrace* and *Tcpflow* tools. The reassembly approach is described in algorithm 1 and works as follows: for each TCP packet received, the system searches the corresponding connection in the connection record list. If the record is valid, the package is inserted into this one. If is an invalid record and the packet contains a SYN flag, a new connection is created for the packet. If the record is invalid and the packet does not contain the SYN flag, the packet is dropped. If the packet containing the FIN or RST flags, the connection is terminated.

#### D. Flow Labeling

Labeling is a necessary step for training and subsequent evaluation of classifiers. Although the utilization of port-based method [20] to traffic flow labeling can introduce errors due its increasing ineffectiveness since flows can be incorrectly labeled, the existence of some inaccurate values in data sets is a common machine learning problem, and a good ML scheme

```

Packet ← GetPacket()
quintuple ← Packet.Quintuple
if Packet.SYN then
    C ← GetNotEstablishedConn(quintuple)
    if C is invalid then
        conn ← NewConnection(Packet)
        InsertNotEstablishedConn(conn)
    else
        InsertPacket(C, Packet)
    end if
else
    C ← GetEstablishedConn(quintuple)
    if C is invalid then
        C ← GetNotEstablishedConn(quintuple)
        if C is invalid then
            DropPacket(Packet)
        else
            InsertInEstablishedConn(C)
            RemoveFromNotEstablishedConn(C)
            InsertPacket(C, Packet)
            if Packet.FIN or Packet.RST then
                ConcludeConn(C, Packet)
            end if
        end if
    else
        InsertPacket(C, Packet)
        if Packet.FIN or Packet.RST then
            ConcludeConn(C, Packet)
        end if
    end if
end if

```

**Algorithm 1:** TCP Flow Reassembly Algorithm

must have the capability to deal with this situation [16]. Although this labeling method was implemented in the first prototyping of NTCS, other sophisticated labeling techniques can be subsequently incorporated.

#### E. Traffic Classification

Real time internet traffic classification enables the solution of hard network management problems by Internet service providers and their equipment suppliers. Network operators, especially in high-speed networks, need to know about the current traffic to quickly respond to support diverse business goals [29].

The presented approach uses real traces in the evaluation of Naive Bayes, Kernel Naive Bayes, C4.5 Decision Tree and K-Nearest Neighbor methods for Internet traffic classification using statistical information derived from packet headers. The Weka [18] (Waikato Environment for Knowledge Analysis) is an Open Source tool implemented in JAVA, and contains a collection of machine learning algorithms for data mining problems. The NTCS utilizes the Weka libraries for the training and evaluation of the machine learning methods. We utilize the IKVM [30] software to enable Java and .NET

interoperability. This tool allows to generate the Weka *dlls*, which were imported to C# code. Thus, it is possible to use the weka classifiers in the classification module of our classifier monitor.

1) *K-Nearest Neighbor*: Among the various supervised statistical methods for pattern recognition, the Nearest Neighbor (NN) technique is the one that achieves better results, without priori assumptions requirement about the training samples distributions [31]. The algorithm assumes all instances correspond to points on an  $n$ -dimensional space  $\mathbb{R}^n$ . A new instance  $X = x_1, x_2, \dots, x_n$ , in which  $x_1, x_2, \dots, x_n$  are the corresponding attributes, is classified by computation of its Euclidean distance to the training instances, and then categorized with the label of the nearest training instance [17].

The KNN (*K-Nearest Neighbor*) classifier extends this idea through the selection of the  $k$  nearest neighbors and the classification of each new instance with the most common class among them [31]. The Euclidean distance between two instances  $X$  and  $Y$  is defined by the following expression, where  $x_k$  and  $y_k$  denote, respectively, the values for the  $k$ -th attribute of  $X$  and  $Y$  instances:

$$d(X, Y) = \sqrt{\sum_k^n (x_k - y_k)^2} \quad (1)$$

2) *Naive Bayes*: The NB classifier is a simple technique that can be applied to the Internet traffic classification problem [32]. Consider a random variable  $C$  that denotes the class on an instance, and a random variable vector  $X$  which represents the observed values of attributes. Furthermore, assume  $c$  a label of a determined class and  $x$  an attribute values vector. Consider a test instance  $x$  to be classified. The most probable class will be the one of most high value to  $P(C = c|X = x)$ , which is the probability of occurrence of class  $c$  class given the instance  $x$ . The following expression presents the Bayes rule, applied to calculate this probability, where  $X = x$  corresponds to  $X_1 = x_1 \wedge X_2 = x_2 \wedge \dots \wedge X_k = x_k$  event, and  $P(C = c)$  represents the priori probability of  $c$ , which is the probability of obtain the class  $c$  without consider the training data:

$$p(C = c|X = x) = \frac{p(C = c)p(X = x|C = c)}{p(X = x)} \quad (2)$$

A common assumption, which is not inherent in Naive Bayesian approach but still often used, is that for each class, the numerical attributes values are normally distributed. According to [33], despite this assumption does not reflects the Internet traffic reality, such approach outperforms some more complex models.

3) *Kernel Naive Bayes*: The KNB learning technique is a generalization of the NB and algorithmically similar to this one in every aspects, except in computation of distribution function  $P(X = x|C = c)$  for continuous attributes, which can be replaced by a variety of non-parametric estimation methods, among them, the kernel density estimation, which, as the name suggests, uses kernel estimation methods instead of simple Gaussian approach [34]. The distribution function is shown

by the following expression, where  $n$  represents the number of training instances that belongs to class  $c$ , and  $\mu_i = x_i$ :

$$p(X = x|C = c) = \frac{1}{n} \sum_i^n g(x; \mu_i, \sigma_c) \quad (3)$$

The intention of KNB method is that the kernel estimation allow the technique has a good performance in domain that violate the normality assumption. According to [9], the simple assumption about the normality of the discriminators is inaccurate and eminent problems arise when the actual distribution is multimodal, and this situation may indicate that the considered class is too large or other distribution must be used for the data analysis.

4) *C4.5 Decision Tree*: The C4.5 method, which is used to generate Univariate decision tree [35], is a statistical classifier because its decision trees can be used in classification [36]. The classifier was created by Ross Quinlan [37] and is an extension of ID3 (Iterative Dicotomizer 3) algorithm, which builds simple decision trees. C4.5 makes use of information entropy to build decision tree in the same way as the ID3.

The C4.5 makes various decisions based on data takes into account all available input features [38]. The algorithm recursively chooses the feature with highest normalized information gain at each node of the tree. The choosed feature splits the data into subsets inriched at one class or the other [36]. The information gain indicates how well a decision will separate the output class from most [38]. We utilized the J48 learning method, which is an open source Java implementation of the C4.5 algorithm.

5) *AdaBoost Ensemble Learning Algorithm*: AdaBoost is an ensemble learning algorithm which alters weak learning classifiers by assigning weights with the overral goal to correctly classify the instances. Generic Boosting manipulate de training data in order to generate distinct classifiers and improves the classification accuracy iteratively. At each iteration, the weights are ajusted so that the weights of the misclassified examples are increased and those of correctly classified examples are decreased. All instances are used in each iteration. Each classifier generated at each round of AdaBoost has an associated weight, which means that a classifier with high degree of correctly identified instances will have more weight in the total ensemble. The final decision of the ensemble is the weighted-majority voting of the hypothesis generated of each generated classifier [38]. See [39] for more details.

## IV. METHODOLOGY

This section starts with a presentation of the network traffic data used to evaluate the NTCS, and finishes with the description of the statistical features.

### A. Data Collection and Experiments

The performance of packet capture and reassembly modules was evaluated for capacity verification, under variable load conditions. The classifier monitor was performed with a Core i5 computer with CPU 2.30 GHz and 4GB of memory. The

TABLE I  
CHARACTERISTICS OF THE TRACE  $T_1$

Characteristic	Description
Packets Number	614282
Capture Size	565.62MB
Capture Duration	3516.79s
Average packet Size	920.78 Bytes
Average capture Rate	1.28 Mbps

TABLE II  
CHARACTERISTICS OF THE TRACE  $T_2$

Characteristic	Description
Packets Number	1579921
Capture Size	1.88GB
Capture Duration	1355.83s
Average packet Size	1195.16 Bytes
Average capture Rate	11.14 Mbps

trace-driven simulation allows flexibility in the evaluation of distinct classifiers and reassembly approaches. This is because the different executions of our system for the same packets trace always generate the same flows set. Without this determinism, it would be extremely difficult to reproduce the same results of an online packet capture, given the possibility of delay and packet loss, for example. For confidently evaluate the online monitor, we used traffic traces collected in a host connected to the broadband Ethernet 100Mbps. Each flow in a reassembly process is configured with a 60 seconds timeout, in order to avoid the storage of old or idle connections, which consume memory and processing resources unnecessarily. This means that TCP flows whose duration is greater than this value are finished by the collector process periodically. We compare the time complexities and the number of reconstructed flows of our flow reassembly module and the external tools Tcptrace, TcpFlow, TcpRecon and Wireshark. Using Weka resources, we use 10-fold cross validation for accuracy evaluation of the aforementioned classification models.

The used traffic traces characteristics, referred as  $T_1$  and  $T_2$ , are presented in Tables I and II. They reflect the communication between a host at State University of Ceara (UECE) and the Internet.

Table III presents the summary of the identified applications flows in the current traces, which were: *Www* (World Wide Web), *Https* (Http protocol over TLS/SSL), *Ftp* (File Transfer Protocol), *Xvttp* (Xvttp Protocol) and *Isakmp* (Isakmp Protocol). The most representative categories in  $T_1$  traffic traces are *Www* and *Ftp* applications. In  $T_2$  traffic trace, the *Https* and *Isakmp* application have a greater number of flow instances. In our study, the classification and training steps are performed at the end of packet capture simulation and reassembly. This methodology is necessary to evaluate the modules of our classifier monitor.

TABLE III  
CATEGORY OF APPLICATIONS

Category	Description	$T_1$	$T_2$
<b>Www-http</b>	World Wide Web HTTP	1022	337
<b>HTTPs</b>	Http Protocol over TLS/SSL	139	27
<b>FTP</b>	File Transfer Protocol	1808	-
<b>Domain</b>	Domain Name Server	-	1
<b>Total</b>		2969	365

## B. Statistical Features

In order to evaluate the classification process, we considerate the following features: *elapsed time between the first and last packets, number of packets, number of bytes, the number of all packets with at least a byte of TCP data payload, the number of all packets seen with the PUSH bit set in the TCP header, and the median and the variance of the number of bytes in IP packet*. Since each attribute is computed for both directions of flow (uplink and downlink), each flow instance has 14 statistical discriminators plus the class label. There was not a proper selection of attributes in this study. We chose some of the most often found attributes in previous published work [9] which could be calculated from the data contained in the header of packets without examining their payload.

## V. RESULTS AND DISCUSSION

Table IV presents some performance metrics of our classifier monitor for the used traces. We observed that the highest achieved throughput for the packet capture and reassembly modules was 3.95 flows per second (fps) for  $T_1$  traffic trace. This means that this is the number of traffic flows are delivered by reassembly process at every second. Although this is a low value, since the mean packet capture throughput of  $T_1$  is only 1.28 Mbps, the reassembly process achieves a throughput of 24997.25 fps at one of packet capture intervals. The average packet capture and reassembly rate, expressed by  $Mbits/(T_{CO} + T_{RE})$ , was 1014.73 Mbps, where  $T_{CO}$  and  $T_{RE}$  are the total duration times of packet capture and reassembly, respectively. Similarly, the same performance metrics are presented for  $T_2$  traffic trace. We can observe no bottlenecks in the reassembly process, which could support the considerate traffic. Our software-based monitor is effective to work in real time for a corporate network, for example. The average delivery delay  $\alpha$  for  $T_1$  and  $T_2$  traffic traces is 0.49s and 7.50s, respectively. This means that this is the average reassembly duration for ever quantum of 30s. We can observe that delivery delay varies greatly between the two compared traces. Although they are from the same packet capture point, the two trace files are essentially different. The throughput collection of the  $T_2$  is much larger than  $T_1$ . Furthermore, the traffic load in  $T_2$  is larger than  $T_1$ . This means that there are less traffic to process in  $T_1$ , and consequently, its delivery delay is lower than the other trace.

In Table V, our system is compared with the Tcptrace, TcpFlow, TcpRecon and Wireshark tools. The TcpRecon was

TABLE IV  
PERFORMANCE OF MONITORING AND REASSEMBLY PROCESSES

Metric	$T_1$	$T_2$
TCP Connections Number	2969	365
Max Capture & Reassembly Throughput	3.95 fps	2.89 fps
Max Reassembly Throughput	24997.25 fps	128.21 fps
Average Capture & Reassembly Rate	1014.73 Mbps	635.34 Mbps
Average Delivery Delay	0.49s	7.50s
Total Monitor Time	75.67s	310.84s

modified to use a flow timeout of 60 seconds. We can observe that the number of flows is not the same between the tools, because the divergence of the used traffic flow concept, as explained previously. We can observe that our adopted reassembly approach execution time is lower than the other tools. Our reassembly scheme was implemented in TCP session Reconstruction Tool, replacing the TcpRecon default policy. In summary, the difference between these two policies is the adopted recently-accessed-first principle and the use of different data structures to hold established and not established TCP connections.

TABLE V  
COMPARISON WITH EXTERNAL TOOLS

Approach/Tool	Flows Number	Reassembly Time
Proposed Approach	2969	75.67s
TcpFlow	5894	118.87s
Tcptrace	3044	612.95s
Wireshark	3036	182.69s
TcpRecon	3036	96.97s

Since TcpRecon and our proposed scheme are written in same language and uses the same packet capture libraries, we also compare the performance of these two policies one of each other. The confidence interval estimation of an event population will have greater reliability if the event is executed at least 30 times [40]. We executed and measured the elapsed times of the aforementioned TCP reassembly policies. The policies were evaluated over the already presented datasets. We computed the average execution time and confidence level for each TCP policy. We consider a high confidence level of 95%. The resulting confidence levels for TcpRecon and our adopted reassembly scheme are presented in Table VI. For the  $T_1$  traffic trace, our scheme obtain a time complexity advantage of 20.31 seconds. For the  $T_2$  traffic trace, there is also a reduction of 9.39 seconds with our approach.

TABLE VI  
PERFORMANCE COMPARISON OF REASSEMBLY POLICIES

Traffic Trace	TcpRecon	Proposed Policy
$T_1$	91.85 ± 5.61 seconds	71.54 ± 3.57 seconds
$T_2$	380.36 ± 14.82 seconds	370.97 ± 7.72 seconds

The Table VII presents the main results about the clas-

sification process. We can observe that C4.5 Decision Tree was able to categorize on average 87.40% and 89.86% of the traffic correctly for the two traffic traces. The AdaBoost ensemble algorithm, using the DecisionStump classifier, was able to categorize on average 78.17% and 89.58% of the traffic correctly. The KNN technique, with  $k = 10$ , was able to categorize on average 86.25% and 91.50% of the traffic correctly, against 72.48% and 80.00% for NB classifier. The duration of classification phase was a few seconds, and the results aim to validate the previous phases of the classifier monitor.

TABLE VII  
GLOBAL ACCURACY PER TRACE

Classifier	$T_1$	$T_2$
<b>C4.5 Decision Tree</b>	87.40%	89.86%
<b>AdaBoost(DecisionStump)</b>	78.17%	89.58%
<b>K-Nearest Neighbor</b>	86.25%	91.50%
<b>Naive Bayes</b>	72.48%	80.00%
<b>Flexible Naive Bayes</b>	64.09%	88.76%

## VI. CONCLUSION

This paper presented the architecture, implementation, and performance of an Internet traffic classifier monitor. The monitor is composed of three modules which were implemented as concurrent processes: capture and preprocessing, flow reassembly, and classification. For the  $T_1$  traffic trace, the throughput reassembly module of the current implementation is 24997.25 flows per second. The average delivery delay is 0.49 seconds. For the classification module, the C4.5 algorithm outperforms KNN and AdaBoost classifiers with average accuracy of 87.40% and 89.86% against 72.48% and 80% for the KNN and AdaBoost methods, respectively.

Future directions for this research includes to incorporate subflow based classification in NTCS to reduce response time. Second, we aim to verify the performance impact of our classifier monitor at gigabit links, which are becoming increasingly common at computer networks. Finally, we could also prototype NTCS with NetFPGA hardware, since the implementation of network systems in hardware is essential for any real time application, particularly in gigabit networks [41].

## REFERENCES

- [1] Y.-T. Han and H.-S. Park, "Game traffic classification using statistical characteristics at the transport layer." *ETRI journal*, vol. 32, no. 1, 2010.
- [2] B. Claise, G. Sadasivan, V. Valluri, and M. Djernaes, "Cisco systems netflow services export version 9," RFC 3954, October, Tech. Rep., 2004.
- [3] B. Claise, M. Fullmer, P. Calato, and R. Penno, "Ipfix protocol specifications," *draftietf-ipfix-protocol-03.txt*, 2004.
- [4] P. Siska, M. Stoecklin, A. Kind, and T. Braun, "A flow trace generator using graph-based traffic classification techniques," in *Proceedings of the 6th International Wireless Communications and Mobile Computing Conference on ZZZ*. ACM, 2010, pp. 457–462.
- [5] H. Kim, K. Claffy, M. Fomenkov, D. Barman, M. Faloutsos, and K. Lee, "Internet traffic classification demystified: myths, caveats, and the best practices," in *Proceedings of the 2008 ACM CoNEXT conference*. ACM, 2008, p. 11.



- [6] P. Mallapragada, R. Jin, A. Jain, and Y. Liu, "Semiboost: Boosting for semi-supervised learning," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 31, no. 11, pp. 2000–2014, 2009.
- [7] A. Dainotti, A. Pescapè, and K. Claffy, "Issues and future directions in traffic classification," *Network, IEEE*, vol. 26, no. 1, pp. 35–40, 2012.
- [8] J. Erman, A. Mahanti, M. Arlitt, I. Cohen, and C. Williamson, "Offline/realtime traffic classification using semi-supervised learning," *Performance Evaluation*, vol. 64, no. 9–12, pp. 1194–1213, 2007.
- [9] A. Moore and D. Zuev, "Internet traffic classification using bayesian analysis techniques," in *Proceedings of the 2005 ACM SIGMETRICS international conference on Measurement and modeling of computer systems*. ACM, 2005, p. 60.
- [10] R. Bar-Yanai, M. Langberg, D. Peleg, and L. Roditty, "Realtime classification for encrypted traffic," *Experimental Algorithms*, pp. 373–385, 2010.
- [11] B. XIONG, C. Xiao-su, and C. Ning, "A Real-Time TCP Stream Reassembly Mechanism in High-Speed Network," *JOURNAL OF SOUTHWEST JIAOTONG UNIVERSITY*, vol. 17, no. 3, 2009.
- [12] J. Postel, "Rfc 793: Transmission control protocol," *DARPA Internet Program Protocol Specification*, 1981.
- [13] P. Agarwal, "TCP Stream Reassembly and Web based GUI for Sachtet IDS," Master's thesis, Indian Institute of Technology Kanpur, Kanpur, India, 2007.
- [14] G. Wagener, A. Dulaunoy, and T. Engel, "Towards an estimation of the accuracy of tcp reassembly in network forensics," in *Future Generation Communication and Networking, 2008. FGCN'08. Second International Conference on*, vol. 2. IEEE, 2008, pp. 273–278.
- [15] V. Jacobson and S. McCanne, "libpcap: Packet capture library," *Lawrence Berkeley Laboratory, Berkeley, CA*, 2009.
- [16] Y. Wang and S. Yu, "Machine Learned Real-Time Traffic Classifiers," in *Intelligent Information Technology Application, 2008. IITA'08. Second International Symposium on*, vol. 3. IEEE, 2009, pp. 449–454.
- [17] L. Jun, Z. Shunyi, L. Yanqing, and Z. Zailong, "Internet traffic classification using machine learning," in *Second International Conference on Communications and Networking in China, 2007. CHINACOM'07, 2007*, pp. 239–243.
- [18] E. Frank, M. Hall, and L. Trigg, "Weka 3-Data Mining with Open Source Machine Learning Software in Java," *The University of Waikato*, 2000.
- [19] G. Szabó, J. Szüle, Z. Turányi, and G. Pongrácz, "Multi-level machine learning traffic classification system," in *ICN 2012, The Eleventh International Conference on Networks*, 2012, pp. 69–77.
- [20] IANA. (2014, May) Internet assigned numbers authority. [Online]. Available: <http://www.iana.org>
- [21] R. Wojtczuk, "Libnids homepage. available at <http://libnids.sourceforge.net/>," 2012.
- [22] S. Ostermann, "Tcptrace homepage. available at <http://www.tcptrace.org/>," 2012.
- [23] A. Orebaugh, G. Ramirez, and J. Burke, *Wireshark and Ethereal network protocol analyzer toolkit*. Syngress Media Inc, 2007.
- [24] G. Maiolini, A. Baiocchi, A. Rizzi, and C. Di Iollo, "Statistical classification of services tunneled into ssh connections by a k-means based learning algorithm," in *Proceedings of the 6th International Wireless Communications and Mobile Computing Conference*. ACM, 2010, pp. 742–746.
- [25] S. Nor, "Near Real Time Online Flow-Based Internet Traffic Classification Using Machine Learning (C4. 5)," *International Journal of Engineering (IJE)*, vol. 3, no. 4, p. 370, 2009.
- [26] L. Bernaille, R. Teixeira, I. Akodkenou, A. Soule, and K. Salamatián, "Traffic classification on the fly," *ACM SIGCOMM Computer Communication Review*, vol. 36, no. 2, pp. 23–26, 2006.
- [27] S. Yahalom, "Available at <http://www.codeproject.com/articles/20501/tcp-session-reconstruction-tool>," 2012.
- [28] C. Partridge, "Jacobson on tcp in 30 instructions." *Message-ID: 1993Sep8*, vol. 213239, 1993.
- [29] T. Nguyen and G. Armitage, "A survey of techniques for internet traffic classification using machine learning," *Communications Surveys & Tutorials, IEEE*, vol. 10, no. 4, pp. 56–76, 2008.
- [30] J. Frijters, "Ikvm, an implementation of java for mono and the .net framework [computer software and documentation]," 2004.
- [31] M. J. Islam, Q. M. J. Wu, M. Ahmadi, and M. A. Sid-Ahmed, "Investigating the performance of naive-bayes classifiers and k-nearest neighbor classifiers," *Convergence Information Technology, International Conference on*, vol. 0, pp. 1541–1546, 2007.
- [32] D. Zuev and A. Moore, "Traffic classification using a statistical approach," *Passive and Active Network Measurement*, pp. 321–324, 2005.
- [33] I. Witten and E. Frank, *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann Pub, 2005.
- [34] G. John and P. Langley, "Estimating continuous distributions in Bayesian classifiers," in *Proceedings of the eleventh conference on uncertainty in artificial intelligence*, vol. 1. Citeseer, 1995, pp. 338–345.
- [35] T. Korting, "C4. 5 algorithm and multivariate decision trees," *Image Processing Division, National Institute for Space Research-INPE Sao Jose dos Campos-SP, Brazil*.
- [36] K. Singh and S. Agrawal, "Comparative analysis of five machine learning algorithms for ip traffic classification," in *Emerging Trends in Networks and Computer Communications (ETNCC), 2011 International Conference on*. IEEE, 2011, pp. 33–38.
- [37] J. Quinlan, *C4. 5: programs for machine learning*. Morgan kaufmann, 1993, vol. 1.
- [38] C. McCarthy and A. Zincir-Heywood, "An investigation on identifying ssl traffic," in *Computational Intelligence for Security and Defense Applications (CISDA), 2011 IEEE Symposium on*, april 2011, pp. 115–122.
- [39] Y. Freund and R. Schapire, "Experiments with a new boosting algorithm," in *MACHINE LEARNING-INTERNATIONAL WORKSHOP THEN CONFERENCE-*. MORGAN KAUFMANN PUBLISHERS, INC., 1996, pp. 148–156.
- [40] R. Jain, *The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling*. Wiley - Interscience, ISBN:0471503361., New York, NY, April, 1991.
- [41] J. Naous, G. Gibb, S. Bolouki, and N. McKeown, "Netfpga: reusable router architecture for experimental research," in *Proceedings of the ACM workshop on Programmable routers for extensible services of tomorrow*. ACM, 2008, pp. 1–7.