

SLIC: Self-Learning Intelligent Classifier for network traffic



Dinil Mon Divakaran*, Le Su, Yung Siang Liao, Vrizlynn L. L. Thing

Department of Cyber Security and Intelligence, Institute for Infocomm Research (I²R), Agency for Science, Technology and Research, 1 Fusionopolis Way, #21-01 Connexis (South Tower), Singapore 138632, Singapore

ARTICLE INFO

Article history:

Received 27 February 2015

Revised 30 July 2015

Accepted 27 August 2015

Available online 4 September 2015

Keywords:

Classification

Learning

Traffic

Protocol

Network

ABSTRACT

Internet traffic classification plays an important role in the field of network security and management. Past research works utilize flow-level statistical features for accurate and efficient classification, such as the nearest-neighbor based supervised classifier. However, classification accuracy of supervised approaches is significantly affected if the size of the training set is small. More importantly, the model built using a static training set will not be able to adapt to the non-static nature of Internet traffic. With the drastic evolution of the Internet, network traffic cannot be assumed to be static. In this paper, we develop the concept of 'self-learning' to deal with these two challenges. We propose, design and develop a new classifier called Self-Learning Intelligent Classifier (SLIC). SLIC starts with a small number of training instances, self-learns and rebuilds the classification model dynamically, with the aim of achieving high accuracy in classifying non-static traffic flows. We carry out performance evaluations using two real-world traffic traces, and demonstrate the effectiveness of SLIC. The results show that SLIC achieves significant improvement in accuracy compared to the state-of-the-art approach.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

Classification of application protocols in the ever changing Internet traffic is a fundamental problem in network security and management. The phenomenal growth of the Internet not only gives rise to new and different kinds of applications, but also leads to dynamic changes in network traffic over time. Accurate traffic classification is of paramount use in a number of network security problems such as, intrusion detection, detection and prevention of various network attacks (like Denial of Service attacks), anomaly detection, lawful inspection, cybercrime forensic analyses, etc. Traffic classification also plays an important role in network management, for example in traffic prediction and Quality of Service (QoS) provisioning. Network administrators find accurate classifier an indispensable tool for securing and

managing their networks. The failure of the naïve and conventional port-based traffic classification in accurately identifying network application protocols has lead to the development of different solutions that exploit the flow-level statistical properties in network traffic.

Over the past many years, the research community has explored various machine learning techniques, such as supervised classification algorithms, unsupervised classification algorithms (clustering), and the combination of the two, namely semi-supervised classification, to provide accurate and efficient solutions for statistical features based flow-level traffic classification [1]. In a supervised traffic classification approach, for example the parametric Naive Bayes (NB), a number of labeled traffic flows called training instances are used to build a model for predicting the unlabeled flows. However, supervised traffic classification faces two fundamental issues.

1. Dependence on labeled training instances: One issue with supervised methods is, they are heavily dependent on the training data. When the number of available labeled training instances is small, or selected in a non-representative

* Corresponding author. Tel.: +6590554094.

E-mail addresses: divakaran@i2r.a-star.edu.sg, dinil.d@gmail.com (D.M. Divakaran), lsu@i2r.a-star.edu.sg (L. Su), liauys@i2r.a-star.edu.sg (Y.S. Liao), vriz@i2r.a-star.edu.sg (V.L. L. Thing).

manner, the classification accuracy could be severely affected. Hence, most existing classifiers require large number of labeled data to achieve a high classification accuracy. On the contrary, in practice, manually labeling traffic flows for training purposes is a time-consuming task and requires strong expertise, and can yet be challenging for new applications and encrypted flows. It is thus motivating to provide an efficient classification method which could produce high accuracy with small number of training instances.

2. Inability to deal with non-static traffic: An observation from existing classification approaches (and not just supervised classification) comes to our attention. There is no utility of a test instance to the classifier: once a test instance (usually, a traffic flow) has been predicted by the model and subsequently assigned a label, it is no longer used in any further predictions. In other words, the classification model and thus its accuracy depends on the selection of the *static* set of labeled training instances. However, intuitively, an “intelligent” classification system should not be restricted to this static training set. Instead it should be able to *self-learn* from the predicted instances, to continuously update and evolve its training set as well as classification model. The aim of this learning process is not only to improve the prediction accuracy of the future instances, but also to steadily adapt to the evolution of the Internet traffic. Such an approach will bring significant improvement in accurately classifying continuously changing network traffic. Indeed, dealing with non-static data is identified as one of the 10 most challenging problems in data mining research [2]. With the prolific increase in the number of users connecting to the Internet, and the advancements being made in the connecting technologies, applications evolve rapidly over time [3]. This in turn leads to dynamic changes in Internet traffic. The need to deal with evolving network traffic is particularly important for classifiers which start with small number of training samples that might not be representative enough for application classes.

Some recent works have proposed new traffic classification algorithms, which utilize flow correlation information for more accurate classification with only a small number of labeled training instances [4,5]. However, designing a new classifier that not only starts with a small training set, but also learns on its own (self-learning), and yet achieves high accuracy is a challenging task. We take up this challenge in our work here, and introduce the concept of *self-learning* to solve the above mentioned two issues. We take our first step by coming up with a new classifier called Self-Learning Intelligent Classifier, or SLIC. This system is a combination of two main parts: a standard supervised classifier for predicting flow label, and a decision maker that decides and selects potentially useful flows for (re-)building the training set and the model. While previous solutions use aggregation of flows, or bag-of-flows (Section 3.1), to improve prediction accuracy, we use this concept for both prediction and self-learning. We show that SLIC performs significantly better than a k -Nearest Neighbor (k -NN) based classifier that also exploits correlation in traffic flows. This k -NN based classifier was recently developed to work with small number of labeled data, and was shown to have high accuracy [4]. The performance improvement achieved by SLIC is much more valuable when seen in the light that, unlike existing classifiers, SLIC is

dynamic in nature and is capable of learning and building a model when faced with realistic dynamic traffic.

With this work, we make the following major contributions:

- We propose a new direction of research in traffic classification, where the challenge is to build a classifier that requires only a small number of labeled data to start with, and yet can learn and adapt on its own to classify application protocols in dynamically changing network traffic with high accuracy.
- We formally propose and develop a novel network traffic flow classification system, called Self-Learning Intelligent Classifier (SLIC), that learns and evolves on its own, to improve the overall classification accuracy over time.
- We carry out extensive evaluation of SLIC using two real-world traffic traces. We compare SLIC with a k -NN based classifier, where the only difference between the two classifiers is that the latter works in a traditional way with no self-learning process.

The remainder of the paper is organized as follows. In Section 2, we provide a brief review of related works on network traffic classification. Section 3 briefs on relevant preliminary concepts. We give an overview of SLIC in Section 4, and develop the proposed system in Section 5. Section 6 presents performance evaluation and results.

2. Related work

Over the last decade, a large number of research works have been proposed on the application of machine learning techniques to statistical feature-based network traffic flow classification. Each traffic flow is represented as a vector of different features such as packet-size statistics (min, max, mean, etc.), connection-size statistics, inter-packet time, etc. We briefly review these works. For more details, interested readers may refer to the survey by Nguyen and Armitage [1].

In supervised classification, a model is built using labeled training flows, which is later used to predict or map future unlabeled flows into one of the classes. Among the pioneering works under this category, Moore and Zuev [6] uses supervised Naive Bayes techniques to classify network traffic flows based on statistical features. By extending [6], Auld et al. [7] utilized Bayesian neural networks for achieving more accurate results. Naive Bayesian classifier was also used to identify Skype connections with encrypted data [8]. Kim et al. evaluated a set of widely used supervised classifiers, including Naive Bayes, Naive Bayes Kernel Estimation, Bayesian Network, Neural Networks, k -Nearest Neighbors (k -NN) and Support Vector Machines (SVM) [9]. The two parametric classifiers, SVM and Neural Networks, and the non-parametric k -NN were the top and closely performing classifiers. Roughan et al. used LDA (Linear Discriminant Analysis) to perform low-error class-of-service mapping of network traffic into interactive, bulk data transfer, streaming and transactional classes [10].

Recently, a series of work by Zhang et al. [4,5,11] utilized flow correlation information to further improve the classification accuracy with small number of training instances based on k -NN and Naive Bayes classifier. As our solution also

exploits flow correlation, we compare SLIC with [4] in our study.

In contrast to supervised methods, the unsupervised approach, known as clustering, takes as input a set of unlabeled flows and automatically groups them into different clusters, and subsequently assigns a test instance into its nearest cluster. McGregor et al. [12] proposed an approach that utilized expectation-maximization algorithm to form clusters. Zander et al. used AutoClass [13] to group traffic flows, and Bernaille et al. [14] explored k -means clustering mechanism and mapped each cluster to an application through payload analysis tool. Erman et al. [15] provided a comparison among three clustering algorithms, namely k -means, DB-SCAN and AutoClass. Their study showed clustering could produce high-purity clusters when the number of allowed clusters is set to be much larger than the number of real applications. In [16], the authors apply GMM (Gaussian Mixture Models) and HMM (Hidden Markov Models) to effectively perform early identification of application class at the beginning of a TCP connection. There are also attempts to use a hybrid combination of flow statistical feature-based clustering algorithms with other approaches. For example, Wang et al. integrated clustering algorithm with payload-based signature matching method [17].

However, there is one significant drawback of clustering-based approaches. As mentioned above, in order to obtain high-purity clusters, the allowed cluster number needs to be set much higher than the number of real applications; it thus creates difficulty in mapping large number of clusters to small number of applications. Without the knowledge of each application, this drawback is hard to overcome. Moreover, the fundamental difference between SLIC and clustering-based approach is that SLIC is still a supervised learning-based method, and naturally the hassle of mapping clusters to applications is not applicable to our system. In their seminal work, Erman et al. [18] proposed a hybrid approach by combining the ideas of both supervised and unsupervised methods, which embedded a small amount of labeled flows into large amount of unlabeled ones. The hybrid approach first performs an unsupervised clustering process and groups unlabeled flows into small clusters, subsequently with the aid of the labeled flows, clusters are then mapped into applications. Following this work, Li et al. proposed a semi-supervised classification based on support vector machine [19], and Wang et al. [20] utilized flow-level side information instead of labeled samples during clustering. Very recently, Grimaudo et al. [21] proposed a learning-based semi-supervised approach using k -means clustering, to adapt to the changes in the network traffic. However their approach does not consider small size flows, besides facing the traditional clustering disadvantages, for example, of mapping large number of clusters to application classes. One drawback common to the semi-supervised approaches is that they need to aggregate a large number of unlabeled flows to produce high purity clusters. On the contrary, SLIC being a supervised approach, the above mentioned constraint does not apply to it, and can therefore be easily deployed for online classification.

A particular concept called online machine learning [22,23] is worth a special mention here. It is a supervised learning approach where data is available in a sequential

order. Although the ultimate goal for online machine learning is similar to the traditional supervised classification, which is to determine a mapping of data points to class labels, there is one major difference. In online learning, the training set and its corresponding prediction model are updated after the arrival of every new data (as mentioned earlier, the data is available in a sequential order), while in traditional approach, one has access to the entire and fixed training set at once.

At first glance, online machine learning may appear to be quite similar to our proposed SLIC system: both have a dynamic training set and a continuously evolving prediction model aiming for better accuracy. However, for online machine learning, once a new data point arrives and has been predicted, its true class label is *revealed* and used for updating the training set and prediction model, before the next data arrives. This is a rather restrictive scenario as the true label of the data might not be available immediately after prediction, or even not available at all. SLIC, on the other hand, performs self-learning and updates the model *without* any knowledge of the true label. That is, the ground truth of predicted instances are never known to SLIC. Under this challenging scenario, SLIC needs to be designed in such a way that enables it to select good representatives from predicted data for learning and improving the classifier.

3. Preliminary

In this section, we briefly recall some preliminary concepts that are used throughout this paper.

3.1. Bag-of-flows (BoF)

An important concept that we adopt for our solution is bag-of-flows — an aggregation of flows or connections based on correlation. As is well-known, a flow is usually defined by the five-tuple of source and destination IP addresses, source and destination ports and protocol. In conventional classification, traffic flows are taken into the classifier as individual and independent instances one by one. However, as mentioned and shown in [4], the correlation information between an aggregation of flows can significantly improve the classification accuracy, especially when only small training sets are available. The correlation information used, or more concretely, the way to construct a correlated aggregation of flows, is based on the “three-tuple heuristic” — in a certain period of time, the flows sharing the same three-tuple {destination IP address, destination port, protocol} form an aggregation.

The concept of using the three-tuple heuristic for flow correlation has been considered in several previous works [20,24,25]. We adapt the same definition and notation as [4], and refer to a correlated aggregation as bag-of-flows (BoF).

Definition 1. A BoF is a set of traffic flows localized in time, and having the same three-tuple of {destination IP address, destination port, protocol}.

In subsequent content, we denote a BoF by

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

where \mathbf{x}_i is a single flow represented as a feature vector. Note, the size of a BoF is not fixed; rather the size depends on how

the flows are aggregated to form BoFs, and might vary from one BoF to another.

During classification, each flow \mathbf{x}_i in \mathbf{X} runs through a classifier and is given a class label. Once every flow in \mathbf{X} obtains a class label, the label of BoF is decided based on some criteria, for example, majority voting; and subsequently every instance in \mathbf{X} is assigned the same class label. In the following, we also use *bags* to refer to BoFs.

3.2. Euclidean distance

Our SLIC is built based on a k -NN classifier, which calculates the distances between the feature vectors of a test instance and its labeled neighbors, to decide on the label for the test instance. Equally importantly, we use a distance metric in the learning phase of SLIC (as explained in Section 5.2). For both purposes, we adopt the commonly used Euclidean distance defined below.

For two vectors \mathbf{u}, \mathbf{v} with the same dimension n , the Euclidean distance is computed as,

$$\text{dist}(\mathbf{u}, \mathbf{v}) = \sqrt{\sum_{i=1}^n (u_i - v_i)^2}, \quad (1)$$

where u_i and v_i denote the i th component of vectors \mathbf{u} and \mathbf{v} , respectively.

3.3. Centroid

We also define the concept of “application class centroid” in the proposed SLIC model. Informally, in an n -dimensional space, a centroid of a set of points is the arithmetic mean in all coordinates of all points in the set. Formally, given a set \mathbf{S} of n -dimensional vectors \mathbf{s}_j , where $j \in [1, |\mathbf{S}|]$, the centroid $\mu_{\mathbf{S}}$ of this set \mathbf{S} is calculated as,

$$\mu_{\mathbf{S}} = \frac{1}{|\mathbf{S}|} \sum_{j=1}^{|\mathbf{S}|} \mathbf{s}_j$$

where the addition and average is taken component-wise.

4. System model overview

In this section, we provide an overview of our proposed system. Fig. 1 depicts the system. Incoming traffic flows are pre-processed to construct bag-of-flows. SLIC consists of two important components - the classifier and the decision

maker. The classifier builds a model based on the training set, and classifies BoFs. The classifier in SLIC is different from the traditional classifier in two ways:

- The model used for predicting the application class of a flow is constantly updated; and this process continues as long as SLIC decides to learn.
- In SLIC, the model is built and predictions are made on individual flows. But, the final label assigned to a flow is based on the label assigned to the bag it belongs to. This is described in Section 5.2.

The second and the most important part of SLIC is the decision maker. The decision maker decides which among the test instances are good enough to be inducted into the training set. One naïve approach to design a self-learning classifier would be to take every predicted instance into the training set and rebuild the model based on the updated training set. However, this approach suffers from serious drawbacks: (i) this greatly affects the efficiency of the classifier, as it has to perform the process of rebuilding and regenerating a model for the classifier on a per-bag basis; (ii) wrongly predicted BoFs will be added into the training set, consequently leading to a negative impact on the prediction accuracy of future instances. Instead of improving the prediction accuracy, the classifier might potentially perform even worse by inducting the wrongly predicted instances. The possibility of such a negative impact is particularly high when the initial training set size is small.

An ideal decision maker should be able to select good candidates for rebuilding the model. By setting certain decision criteria, the system needs to decide which test instance, after being predicted and assigned a class label, is potentially good for improving the prediction. Such candidates should be inducted into the training set for rebuilding the model. Designing such decision criteria is the core part of our contribution. On one side, inevitably, we should allow the system to make mistakes occasionally, but not too many. The portion of wrongly predicted BoFs that pass the decision criteria, and have subsequently been added into the training set, should be small and make no significant adverse effect on the overall accuracy. On the other hand, the decision criteria should not be too strict such that only a few instances are inducted; this would make self-learning an extremely slow process.

To minimize error, we design the decision criteria to be dependent only on the initial training set, which is a static and small set of labeled data with ground truth. The test instances that pass the decision criteria are inducted to

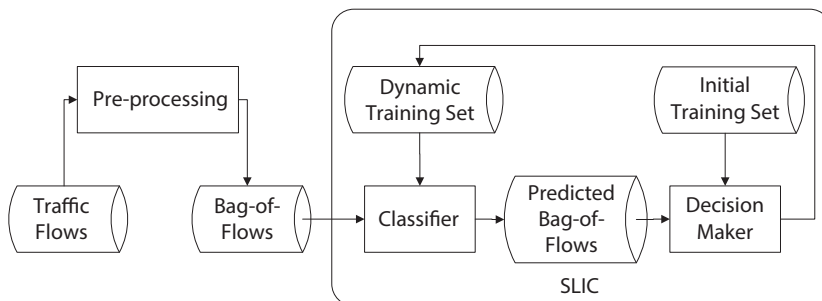


Fig. 1. SLIC system model.

Table 1

Notations.

\mathcal{I}	Initial (labeled) training set
\mathcal{D}	Dynamic training set; $\mathcal{D} \supset \mathcal{I}$
\mathcal{T}	Testing set
A	Set of application classes
N	Maximum (permitted) size of dynamic training set \mathcal{D}

gradually build up the dynamic training set. Each time the dynamic training set inducts a set of test instances, the model is rebuilt. We detail the system in the following section below.

5. Self-Learning Intelligent Classifier (SLIC)

SLIC builds its first model using a small set of labeled data, which we call the initial training set \mathcal{I} . Let A be the set of application classes, with $|A|$ denoting the number of classes. Each element of \mathcal{I} is a tuple (\mathbf{x}, a) , where \mathbf{x} is a flow (feature vector) and $a \in A$ is an application class. For each $a \in A$, $\mathcal{I}_a \subset \mathcal{I}$ is the subset of the initial training data belonging to class a . Obviously, \mathcal{I}_a s form a partition of \mathcal{I} . In implementation, \mathcal{I}_a does not maintain label information (as all labels are the same, and it is the label (a) ; therefore an element of \mathcal{I}_a is not a tuple, but a flow of class a).

The dynamic training set, denoted by \mathcal{D} , is a superset of the initial training set \mathcal{I} . SLIC rebuilds its model as and when \mathcal{D} changes. SLIC uses k -NN as the underlying classifier; therefore, the computational complexity of SLIC is dependent on the size of the training set. For this reason, the rebuilding of the model stops when the size of the dynamic training set \mathcal{D} reaches a pre-defined limit N . The testing set is denoted by \mathcal{T} . As mentioned before, each instance in the testing set is a bag-of-flows. Table 1 lists the notations that are commonly referred in this paper.

In the following, we develop SLIC. After discussing the training phase in Section 5.1, we describe the classifier and the decision maker in Section 5.2.

5.1. Training

Algorithm 1 gives the steps for training and building a model for the classifier in SLIC. Comments in the algorithms are preceded by ‘##’.

Algorithm 1 train(m , trace).

```

1: for each  $a \in A$  do
2:    $\mathcal{I}_a \leftarrow$  Initialize training set for class  $a$  with  $m$  unique
     flow instances from trace ##  $\mathcal{I}_a = \{\mathbf{x}_1^a, \mathbf{x}_2^a, \dots, \mathbf{x}_m^a\}$ 
3:    $\mu_a \leftarrow$  centroid( $\mathcal{I}_a$ )
4: end for
5:  $\mathcal{I} \leftarrow \bigcup_{\forall \mathbf{x} \in \mathcal{I}_a, \forall a \in A} (\mathbf{x}, a)$ 
6:  $\mathbf{M} \leftarrow$  build_kNN_model( $\mathcal{I}$ )
7: return
```

The first step in the **for** loop segregates the given traffic trace with labeled data to form an initial training set for each application class \mathcal{I}_a ($a \in A$) consisting of m unique flows. In line 3, we compute centroids of application classes (refer

Section 3.3); centroids later (in Algorithm 5) play an important role in deciding whether a test instance is ‘good’ enough to be inducted into the dynamic training set \mathcal{D} . Line 5 constructs the initial training set by taking the union of all tuples. The initial training set \mathcal{I} is then passed to the function `build_kNN_model` that builds a model based on the k -NN classifier. This model is used for predicting the test instances later, in Algorithm 2.

Algorithm 2 slic(\mathbf{M} , \mathcal{T} , \mathcal{I}).

```

1:  $\mathcal{D} \leftarrow \mathcal{I}$  ## Initialize dynamic training set
2:  $b^s \leftarrow 0$  ## Initialize size of inducted batch of bags
3: for each  $\mathbf{X} \in \mathcal{T}$  do
4:    $p\_label \leftarrow$  classifier( $\mathbf{M}$ ,  $\mathbf{X}$ )
5:   if (decision_maker( $\mathbf{X}$ ,  $p\_label$ ) == True) then
6:     store_batch( $\mathbf{C}$ ,  $\mathbf{X}$ ,  $p\_label$ )
7:     increment  $b^s$ 
8:     if  $b^s == b^{\max}$  then
9:       induct( $\mathcal{D}$ ,  $\mathbf{C}$ )
10:       $\mathbf{M} \leftarrow$  rebuild_kNN_model( $\mathcal{D}$ )
11:       $b^s \leftarrow 0$ 
12:     end if
13:   end if
14: end for
15: return
```

Strictly speaking, there is no training phase for a k -NN based classifier, as it is a non-parametric approach. Thus, the `build_kNN_model` function merely aggregates the labeled instances and forms the training set. However, for ease of exposition, we still name this process as training phase.

Algorithm 1 is executed only once, but the model \mathbf{M} built in this algorithm will be rebuilt multiple times later.

5.2. Classifier and decision maker

Once the initial model is built, SLIC performs the prediction of test instances and improves the model on its own, as per Algorithm 2. After providing a high level description of Algorithm 2 (slic) below, we discuss the underlying algorithms.

5.2.1. slic

In Algorithm 2, SLIC begins by initializing the dynamic training set \mathcal{D} to the initial training set \mathcal{I} . Looping over all test instances, in line 4, the function `classifier` (Algorithm 3)

Algorithm 3 classifier(\mathbf{M} , \mathbf{X}).

```

1: label_count: Array[1, ..., |A|]
2: for each  $\mathbf{x} \in \mathbf{X}$  do
3:    $p \leftarrow$  kNN_predict( $\mathbf{M}$ ,  $\mathbf{x}$ )
4:   increment label_count[ $p$ ]
5: end for
6: return  $\arg \max_p \text{label\_count}$  ## Break ties randomly
```

assigns the test instance \mathbf{X} a class label based on the model \mathbf{M} . Subsequently, in line 5, the function `decision_maker` (details given in Algorithm 5) decides if the predicted instance

should be inducted into the dynamic training set. If the decision is positive, the test instance is considered as an element of the candidate induction set, denoted as \mathbf{C} , and stored in a batch (line 6). When the batch size reaches a maximum value b^{\max} , a subset of the stored batch of instances is inducted in line 9. The algorithm `induct` is described in [Algorithm 4](#). Af-

Algorithm 4 `induct(\mathcal{D} , \mathbf{C})

---`

```

1:  $l \leftarrow \arg \min_a |\mathbf{C}_a|$  ## application class with min. no. of in-
   ducted instances
2:  $\max\_induct \leftarrow (1 + y\%) \times |\mathbf{C}_l|$ 
3:  $\mathbf{B} = \emptyset$  ## Set of instances that will be inducted
4: for each  $a \in A \setminus l$  do
5:   if  $|\mathbf{C}_a| > \max\_induct$  then
6:      $\mathbf{B} \leftarrow \mathbf{B} \cup \{\text{randomly selected } \max\_induct \text{ no. of in-}$ 
        $\text{stances from } \mathbf{C}_a\}$ 
7:   end if
8: end for
9:  $\mathcal{D} \leftarrow \mathcal{D} \cup \mathbf{B}$ 
10: Reset  $\mathbf{C}$  and  $\mathbf{B}$  to  $\emptyset$ 

```

ter inducting new instances into the dynamic training set \mathcal{D} , the model \mathbf{M} has to be regenerated, and this is achieved by the function `rebuild_knn_model` in line 10.

Observe that the model built in the training phase uses data with true labels (line 6, [Algorithm 1](#)); but the model rebuilt in the testing phase here uses data with predicted labels. One advantage of k -NN is that, similar to the `build_knn_model` in [Algorithm 1](#), rebuilding is inexpensive as it merely adds more data to the existing set. Next, we provide detailed descriptions of the three important algorithms used in [Algorithm 2](#), namely, `classifier`, `induct` and `decision_maker`.

Remark. For simplicity, [Algorithm 2](#) assumes a finite testing set; but, it is straightforward to modify the algorithm to work on a continuously flowing stream of data.

5.2.2. classifier

Recall, a test instance is a bag-of-flows (BoF). In [Algorithm 3](#) for `classifier`, the class label of a test instance \mathbf{X} is predicted in line 3. For deciding on the application class of a BoF, we use k -NN to predict each of the flows in the given BoF. The label to which the maximum number of flows belong to, is returned as the label of the BoF (line 6). Here we use the simple majority voting for assigning label to the BoF, and ties are broken randomly. Note, the k -NN prediction in line 3 of [Algorithm 3](#) is based on the specific model \mathbf{M} which is not static. As explained later, the model \mathbf{M} is regenerated multiple times until the dynamic training set inducts no more instances.

5.2.3. induct

The most important decision on inducting a test instance is taken by the function `decision_maker`; however, this in itself is not sufficient to guarantee good performance. A common problem that impedes performance in learning systems is the imbalance in class distribution [2]. Imbalanced or skewed class distribution means that the number of data in one class is significantly different from other classes. The

performance of k -NN too is known to be affected by this problem, and there are solutions to tackle this [26,27]. Thus, when SLIC inducts new test instances into the training set, [Algorithm 4](#) ensures that the resulting dynamic training set is not considerably imbalanced. More precisely, the algorithm `induct` controls the imbalance in the dataset.

In [Algorithm 4](#), line 1 finds the label that has the minimum number of induction candidates in the given batch \mathbf{C} . Set this label to be l ; and the number of flows with this label (classified as l) is $|\mathbf{C}_l|$. The following steps ensure that the difference in the number of data in any two classes is within $y\%$ of $|\mathbf{C}_l|$. Therefore, the maximum number of instances that is inducted in any class with label a , $a \neq l$, is $(1 + y\%) \times |\mathbf{C}_l|$ (line 2). For a class with label a , if the number of candidates ($|\mathbf{C}_a|$) is greater than that can be inducted, the algorithm in line 6 selects the necessary number of instances from \mathbf{C}_a randomly. We note that there are other ways of choosing from the candidates. For example, one approach is to rank the candidates according to their distances to the class centroid, and select $(1 + y\%) \times |\mathbf{C}_l|$ with the shortest distance. However, we decide to have a simple solution of choosing the candidates randomly. The computational complexity of induction for each bag is $\mathcal{O}(1)$, and hence that of a batch of candidates is $\mathcal{O}(b^{\max})$. SLIC being a continuously learning process, the induction will not stop until the maximum number of training instances allowed is reached.

5.2.4. decision_maker

The algorithm `decision_maker` is the most important and distinguishing feature of SLIC. It decides whether a given (and predicted) test instance is a good candidate that can be inducted into the dynamic training set \mathcal{D} . Observe that a test instance has already been classified as belonging to one of the application classes before facing the decision on induction. It is worth noting that this label is not provided as ground truth, but rather the one predicted by SLIC; therefore, the label need not be correct. This differs from the online learning approaches in the literature (refer to [Section 2](#)). For SLIC to improve its classification performance, the challenge is to induct test instances that are “good”.

Our design of the criteria for inducting a given test instance is inspired by the k -means clustering [28]. At a high level, we measure the distance between the centroid of a given test instance (recall, a test instance is a BoF) to centroids of all classes, and induct the instance only if it is sufficiently near to the centroid of the predicted class. The distances to the centroids of different application classes are used to estimate the goodness of a test instance. The detailed steps are given in [Algorithm 5](#). We explain the steps.

It is important to note here that the centroids of application classes are computed using the initial training set \mathcal{I} (see line 3 of [Algorithm 1](#)). That is, centroid μ_a of an application class a , $a \in A$, is dependent only on its initial training set \mathcal{I}_a , and is therefore pre-computed as well as static. The first conditional statement in [Algorithm 5](#) checks if the size of the dynamic training set exceeds the limit specified by N ; and if so, no more instance will be inducted. The size of \mathcal{D} may still go beyond N after a round of induction. This is not a problem, because the batch size limit b^{\max} can be used to ensure that exceeded size of \mathcal{D} is not, say, 5% more than N .

Algorithm 5 decision_maker(\mathbf{X} , p_label).

```

1:  $\mu_a$  : Centroid of flows  $\mathbf{x}$ ,  $\mathbf{x} \in \mathcal{I}_a$ 
2: if size( $\mathcal{D}$ ) >  $N$  then
3:   return False
4: end if
5:  $\mu_{\mathbf{x}} \leftarrow \text{centroid}(\mathbf{X})$ 
6: for each  $a \in A$  do
7:    $d[a] \leftarrow \text{dist}(\mu_a, \mu_{\mathbf{x}})$  ## refer Eq. (1)
8: end for
9:  $l \leftarrow \arg \min_a d$ 
10: if  $l \neq \text{p\_label}$  then
11:   return False
12: end if
13: sort( $\mathbf{d}$ ) ##  $\mathbf{d}[1] = \mathbf{d}[l]$ 
14: if  $\theta \times \mathbf{d}[1] > \mathbf{d}[2]$  then
15:   return False
16: end if
17: return True

```

As a test instance is a bag-of-flows, we also compute its centroid; this is done in line 5 of Algorithm 5. Next, in the **for** loop of the algorithm, using Eq. (1) we compute the Euclidean distance of the centroid of the given test instance \mathbf{X} to the centroid of each of the classes. There are two conditions that need to be satisfied for a BoF instance to be inducted. First, the label of the application class that is nearest to the BoF instance should be the same as the label p_label assigned to the BoF by SLIC (lines 9–12). Second, in the case where the first condition is satisfied, the distance to the second nearest application class should be θ times more than the distance to the nearest class (lines 13–16). Only if both criteria are satisfied, will SLIC decide to induct the test instance (line 17). We refer to θ as the *distance factor*.

Though Algorithm 5 consists of a loop and a sort, they run in constant times as the number of applications and features are constant values. Therefore, the complexity of decision_maker is dominated by line 5, which is linear in the number of flows in a bag. However, as the maximum size of bags is also fixed at the beginning of classification, and is typically small, the algorithm takes constant time to decide on a given test instance.

We now conclude the descriptions of all algorithms in the proposed SLIC system.

6. Performance evaluation

In this section, we evaluate SLIC over two real-world traffic datasets. We compare the performance with the BoF-based nearest neighbor classifier proposed in [4], which was shown to outperform other traditional classifiers (such as SVM, Neural Nets and C4.5 Decision Tree). The prediction of a bag-of-flows instance based on the majority of the decisions on flows (constituting the bag) is one of the specific solutions in [4]. As this naturally appears like an extension of k -NN to bag-of-flows classification, in the following, we refer to this classifier simply as k -NN.

SLIC is designed to function as an online traffic classifier; however, it is not possible to evaluate classification solutions in the absence of ground truth, as is the case in online

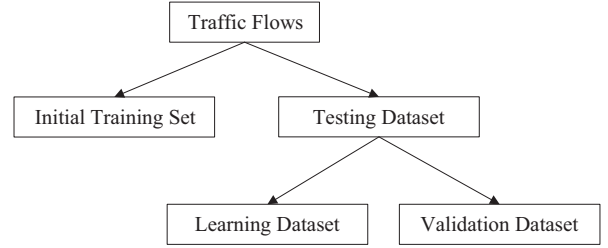


Fig. 2. Partitioning of traffic trace into different datasets.

testing. Therefore, in the next section, we present a few necessary changes that need to be applied to SLIC for offline performance evaluation. The information on the datasets and the metrics used for evaluations are given in Sections 6.2 and 6.3, respectively. After discussing the settings for the experiments in Section 6.4, we present and discuss the results in Section 6.5.

6.1. Offline testing

As opposed to an online scenario where there is continuous stream of network traffic flows, only finite number of flows collected as static traffic traces (detailed in Section 6.2) are available for offline evaluation. In this section, we explain the modifications applied to the algorithms described in Section 5 to effectively perform offline evaluation, and argue that these changes will have no significant impact over the online performance of SLIC.

6.1.1. Modification 1

The first modification is to create the testing set. In traditional classifiers, creating a testing set is not a concern, as once the training instances are randomly selected from the traces, the remaining data can be used for testing. However, SLIC is a continuously learning classifier that inducts predicted instances from the testing dataset. For unbiased prediction, the inducted instances have to be removed from the testing data, thereby shrinking its size. In effect, if the traditional segregation is applied, the evaluation of SLIC will be over a reducing testing data, making the evaluation difficult and even inaccurate. Such a problem is not seen in existing works — as there is no learning phase in them, the classification is performed over the testing set only once. To overcome this problem, in (each run of) an experiment, we first randomly select the training samples from the traces to form the initial training set \mathcal{I} . The remaining data is then randomly split into two disjoint and smaller sets, which we call *learning* and *validation* datasets. Refer to Fig. 2 for an illustration of the partitioning of a traffic trace.

The learning dataset is dynamic, and SLIC inducts test instances only from the learning dataset for self-learning. The inducted test instances are removed from the learning set; therefore, the number of instances in the learning set reduces over time. The validation set (comparable with the standard “testing set”) remains static for the entire execution; and it is used for the sole purpose of evaluating the classifiers, (BoF-based) k -NN and SLIC. In the beginning of an execution, before any test instance is inducted (i.e., before the learning begins), the validation set is evaluated for classification results

using the initial training set \mathcal{I} . As this happens before SLIC comes into play, this classification corresponds to k -NN and the results obtained are thus due to the (BoF-based) k -NN classifier. Then onwards, SLIC inducts data from the learning dataset and rebuilds its model over iterations. The last evaluation on the validation dataset using the final model gives the classification results for SLIC.

Two additional points need to be noted for the above mentioned testing setup. First, to validate that the learning and validation datasets show no particular bias, at the start of an experiment, the initial training set is also used to predict the flows in the learning dataset. The accuracy over the learning dataset is compared with the one obtained on the validation dataset. Our experiments show, the two accuracies vary only within 1% difference, which is normal due to randomness. This justifies the partitioning is valid. Second, the partitioning of the testing dataset does not necessarily need to be on an equal-sized basis. We also tested on a 60% learning – 40% validation partition, and the performance showed no significant difference. We believe, for a large dataset, as long as the splitting is not too extreme (such as 90–10% partition), both datasets would have sufficient flows to make them representative.

6.1.2. Modification 2

The second modification is on the simulation of the batch learning process. As mentioned earlier in Section 5.2, BoFs are inducted in batches. Once a batch of test instances is inducted, the model used for prediction is regenerated. Since a newly regenerated model is (presumably) better than its predecessor, there is a higher chance of predicting the right class labels of newly arriving BoF instances. This eventually increases the chance of inducting a good predicted instance. One can visualize this to work smoothly in an online setting, where there is continuous stream of traffic flows arriving at the classifier. However, in offline testing, there are only a finite number of flows (and hence BoFs), and the improvement attained by going through this finite test data might be limited.

To simulate an online scenario, we introduce *learning iteration* – SLIC iterates over the learning set multiple times. At the beginning of each iteration, when the predicted BoFs run through the decision criteria, all BoFs that satisfy the constraints are stored as a batch of candidates for induction. At the end of the iteration, Algorithm 4 (induct) is used to select a set of BoFs from this batch. These selected bags are inducted into the dynamic training set \mathcal{D} , and the model updated. At the same time, the inducted BoFs (or equivalently, individual flows in the bags) are deleted from the learning dataset. This completes one learning iteration. In the next round, SLIC iterates on the remaining learning dataset, to check if it can induct more BoFs with its new model, and thereby continue learning. This continuing process perfectly simulates online classification – bags that could have potentially been rejected in one iteration may pass the decision criteria in the next iteration, due to an improved model. The modification also simulates the impact of the learning process. However, as there is no continuously flowing data stream, we need to add additional stopping criteria for the learning process. In the online setting, the induction (and hence model regeneration) stops when the size of the

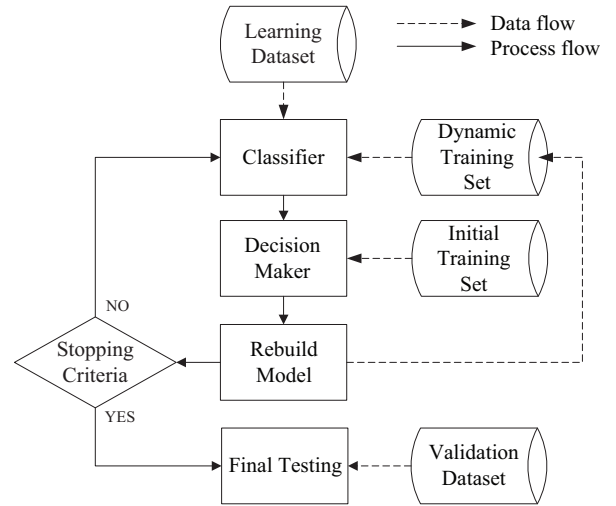


Fig. 3. Illustration of offline version of SLIC.

dynamic training set \mathcal{D} reaches the predefined limit N . In the offline setting, having just this condition might lead to large number of iterations and longer execution times. Hence, for offline setting, we decide to stop the iterations when one of the following conditions is met:

- when the size of the dynamic training set reaches N ;
- when the number of iterations over the learning dataset reaches η , a pre-defined threshold;
- when no test instance is inducted in an iteration over the learning dataset.

Note that if the third condition is met, the model does not change, and hence there is no point iterating again. When the iteration stops due to one of the above three conditions, the learning in SLIC comes to an end, and the last generated model becomes the final one. We highlight, this modification is purely for offline testing; in online classification, the learning has no reason to stop until the maximum training size of N is reached (and this limit is for computational efficiency). Fig. 3 illustrates the working of SLIC with the above mentioned modifications.

6.2. Traffic traces

For testing, we use two large-scale real network traffic flow traces that were recently used to evaluate new classifiers [4,5]. The *isp* trace is a full payload traffic dataset collected at a 100 Mbps edge link of a small-medium ISP located in Melbourne, Australia during Nov–Dec 2010 [20]. It contains around 200,000 traffic flows. The *wide* trace was taken from a US–Japan trans-Pacific backbone line with a 150 Mbps Ethernet link that carries commodity traffic for the WIDE organization [29]. The trace was collected at 72-hour basis from 18–20 March 2008 and 30 March–4 April 2009. The *wide* trace contains around 1.6 million flows. Since the number of flows across applications is highly skewed in both traces (with differences in orders of magnitude), we extracted flows in such a way so as to ensure that the maximum number of flows in any class is not more than five times of any other classes.

Table 2Features used for the two datasets, *isp* and *wide*.

<i>isp</i>	<i>wide</i>
c2s maximum packet size	c2s maximum packet size
c2s minimum inter-packet time	c2s packet size standard deviation
s2c minimum inter-packet time	c2s volume of bytes transferred
s2c maximum packet size	s2c maximum packet size
s2c minimum packet size	s2c packet size standard deviation
s2c packet size standard deviation	s2c volume of bytes transferred

After this, protocols with insufficient flows were not considered for evaluation.

The ground truth for the flows in both traces are obtained via deep packet inspection (DPI) tools, such as I7-filter [30] and Tstat [31]. The *isp* trace has seven application protocols, each with sufficient number of flows for meaningful testing: BitTorrent, HTTP, IMAP, POP3, SMTP, SSH, XMPP. The *wide* trace has been categorized into five application classes, four of which are the protocols BitTorrent, DNS, HTTP, and FTP, and the fifth one is an aggregate class *Mail*. The Mail class combines the protocols POP3 and SMTP. The reason for combining the two mail protocols into one single application is because POP3 on its own does not have sufficient data in *wide* compared to other protocols. Moreover, as mail protocols are quite similar, grouping them into a larger Mail category is justified.

We consider bidirectional flow features for classification. In both traces, the same set of 20 statistical features are available. The features can be broadly seen to be as one of the two types, namely client-to-server (c2s) and server-to-client (s2c). These features cover wide range of properties of real network traffic flows, for example number of packets/bytes transferred, minimum/maximum/average values of packet sizes and inter-packet times, and their standard deviations. However, not all features mentioned above are used for testing. Feature selection is an important data pre-processing step to remove redundant and irrelevant features. Feature reduction is also useful in reducing the running time of the solutions, especially for the nearest-neighbor-based classifiers due to the “curse of dimensionality”. Random forest (RF) has been shown to be an efficient and accurate approach for finding the important features [32]. We run RF on both traces and on each trace RF produces a list of (*feature, importance*) pairs. We select six features with the highest importance from these two lists and use them for evaluations. The six selected features for *isp* and *wide* traces are given in Table 2.

6.3. Metrics for evaluation

To evaluate the performance of SLIC, we use two commonly and widely adopted metrics in classification, namely overall accuracy and F_1 score. The overall accuracy is the ratio of the sum of correctly predicted BoFs (across all application classes) to the size of the population of BoFs. The F_1 score for a class is defined as:

$$F_1 \text{ score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

Precision and recall for an application class are defined as:

$$\text{precision} = \frac{\# \text{True Positive}}{\#(\text{True Positive} + \text{False Positive})};$$

$$\text{recall} = \frac{\# \text{True Positive}}{\#(\text{True Positive} + \text{False Negative})}.$$

Precision gives the fraction of correctly predicted instances of all those predicted for (and as) a class; recall is the fraction of correctly predicted instances of the true instances of a class.

6.4. Settings

We use the unweighted k -NN. We observed that $k = 1$ and $k = 3$ give the highest accuracies, with overlapping confidence intervals; refer to Appendix A. As an inducted bag-of-flows brings in multiple flows into the dynamic training set, intuitively a value of k that allows voting is a natural choice. We therefore set the value of k to three in our experiments. We utilize the Python *scikit-learn* library [33] to implement the k -NN classifier as well as for data pre-processing such as feature selection using Random Forests. The value of N , the maximum size of the dynamic training set \mathcal{D} , is set to 2500 flows for experiments on both traces, and this value is independent of the number of application classes. As the training set size influences the efficiency of k -NN classifier, a larger N slows down the experiments. The threshold to the number of iterations (over the learning dataset), η , is set to five; we further discuss on this parameter in Section 6.5.3.

Another important parameter which is worth a special mention is the size of BoFs. For the upper bound of the bag-sizes, having a smaller value is preferred, as it reduces the waiting time required to aggregate flows into bags. We set the upper bound of bag-size to be 10 for both traces, and refer to the Appendix B for experiments and observations leading to this value. The bag-size range for *isp* trace is set to 3–10, and for *wide* it is set to 1–10. Later, in Section 6.5.3, we experiment the performance by varying the minimum bag-size.

Experiments are distinguished by the number of instances per application class, m , used for building the initial training set \mathcal{I} . Note, $|\mathcal{I}_a| = m, \forall a \in A$. The number of experiments is therefore equal to the number of values that m takes; and the list of values taken by m is [5, 10, 20, 30, 40, 50]. Each experiment is run 100 times, sufficient enough to get 95% confidence interval on accuracy. In each run of the experiment, the initial training set is chosen randomly. The learning and validation datasets are also created randomly for each run.

6.5. Results

We present the results on the *isp* and the *wide* traces. Each point on a plot is the mean value of 100 runs. 95% confidence intervals are marked for the figures plotting accuracies. For the quantitative analyses below, we use rounded values.

6.5.1. Trace *isp*

The new parameter introduced in the design of SLIC is θ , the distance factor in Algorithm 5. For the *isp* trace, we experimented with different values of the distance factor; Fig. 5(a)

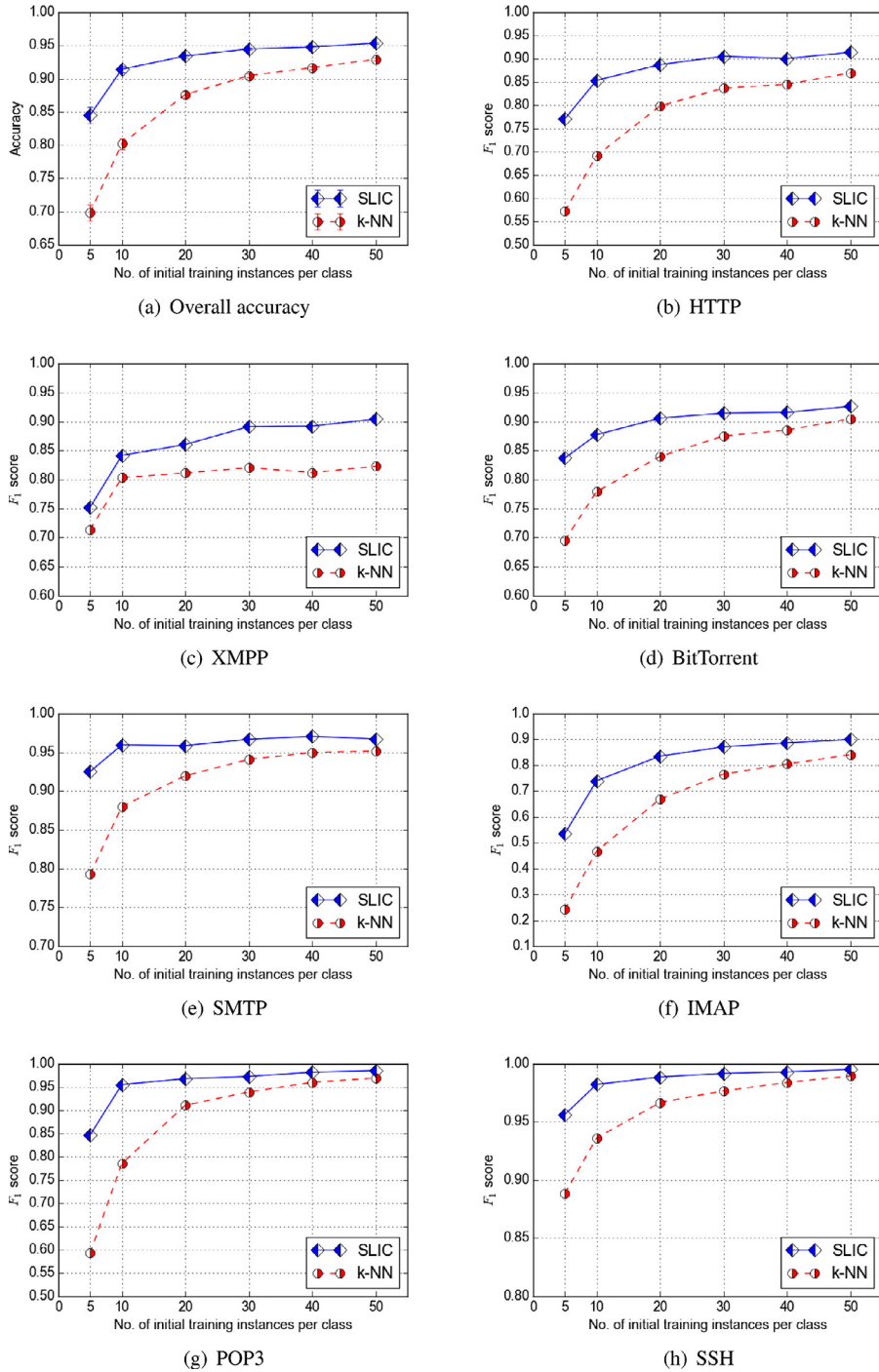


Fig. 4. Classification results on *isp* trace. Overall accuracy with 95% confidence interval is plotted in (a); (b)–(h) plot the F_1 scores for the different protocols. The metrics are plotted against m , the number of training instances (flows) per application class.

plots the corresponding accuracies. The accuracies obtained with $\theta = 1$ and $\theta = 2$ are similar, and higher than those obtained with other values of θ . For the following detailed set of experiments on the *isp* trace, we set the distance factor to one. This means, if a bag predicted as class a finds the centroid of \mathcal{I}_a nearest to it (among all applications \mathcal{A}), the bag is inducted into the dynamic training set \mathcal{D} .

Fig. 4(a) gives the overall classification accuracy achieved on the *isp* trace with different number of training instances per class (m). The bag-size range was set to 3–10. We observe that the (BoF-based) k-NN, which does not learn using the predicted data, achieves a moderate overall accuracy of 70% when it uses five instances per class for training (i.e., in this experiment $|\mathcal{I}| = m \times |\mathcal{A}| = 35$). With self-learning at

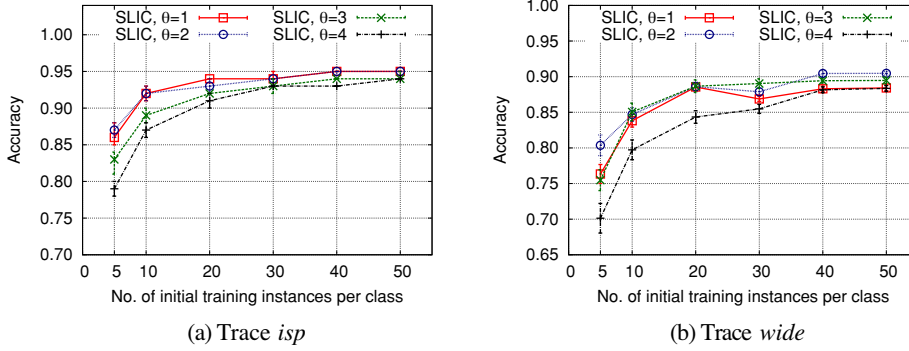


Fig. 5. Accuracy, for different values of distance factor θ .

its core, SLIC is seen to improve this performance significantly. The accuracy achieved by SLIC is 85%, which is 21% relatively higher than that achieved by k -NN. While the accuracies achieved by both the classifiers increase with increasing m , SLIC further improves the accuracies achieved by BoF-based k -NN. In particular, SLIC is able to achieve more than 91% accuracy with just 10 training instances per class.

The remaining figures, Fig. 4(b)–(h), plot the F_1 score for the seven protocols. For HTTP, observe that the best classification result obtained using k -NN is 87%; and this required a large initial training size of 50 per class. SLIC achieves a similar 86% F_1 score with just 10 training instances per class. The best F_1 score achieved for XMPP by k -NN is 82%, whereas SLIC is able to achieve close to 90% F_1 score, thereby increasing the classification performance by almost 10%. For BitTorrent, SMTP, IMAP, POP3 and SSH, SLIC brings in significant performance improvement over k -NN when the initial training data is small. The F_1 score for POP3, with $m = 5$, improved from a mere 59% to an impressive 85% due to SLIC. Similarly, for SMTP, the F_1 score improved by almost 18% from 79 to 93% with five initial training instances. For SSH (in Fig. 4(h)), we also observe that SLIC achieves close to the best classification result of 98% with just 10 training instances per class.

6.5.2. Trace wide

As mentioned earlier, the *wide* trace has been categorized into five application classes, namely HTTP, DNS, BitTorrent, Mail and FTP. Also recall that we set the minimum bag-size to be one, thereby allowing bags of all sizes to be considered for classification. Fig. 5(b) plots the accuracies for different values of the distance factor. We observe that $\theta = 2$ performs best among the different values of θ . For the experiments that follow, the distance factor is set to two.

Fig. 6(a) plots the overall accuracy for k -NN and SLIC. In comparison to the *isp* trace, both the classifiers achieve lower accuracies on the *wide* trace. Classification of bags of sizes lower than three might be the reason for this; to verify we experiment on the effect of bag-sizes in the next section and provide further analysis.

With five instances per class ($m = 5$), the improvement in accuracy brought about by SLIC over k -NN is over 26%. From Fig. 6(a), we observe that SLIC achieves an overall accuracy of 80%. Also note that with increasing m , SLIC is still able to bring in further improvements in accuracies over k -NN.

Fig. 6(b)–(f) plot the F_1 scores for the different protocols. The F_1 scores achieved by SLIC for HTTP is always above 90%, while this is not the case with k -NN. For DNS, the performance achieved by SLIC is very high for all sizes of initial training set; for $|m| = 10$, SLIC achieves an F_1 score of 93%, which is higher than that achieved by k -NN for $|m| = 50$. With five instances per class, the relative improvement in accuracy achieved by SLIC over k -NN for BitTorrent is more than 50%. The accuracy of SLIC for Mail is moderate when initial training sets are small in size; yet the relative improvement over k -NN is huge — the F_1 score increases from 0.35 to 0.62 due to SLIC, for $|m| = 5$. For FTP, though SLIC improves the classification performance attained by k -NN, the absolute F_1 score reaches 80% only when $m = 40$. We suspect this may be due to the less number of test instances available under the FTP class — FTP had 85% lesser number of BoFs than most other classes.

6.5.3. Discussions

Bag-size: As mentioned earlier, BoF size plays an important role for the improvement of classification accuracy. Bag-size dictates the time (the pre-processor of) the classification system needs to wait in order to gather sufficient number of flows to form bags. This waiting time can be too long, possibly indefinite in the worst case. Consider the case where a particular user in an enterprise network visits a suspicious site whose IP address keeps changing. Unless other users visit the same site at almost the same times, the bags formed will always be of size one. Thus, it is important for a classification system to be able to predict on a single flow (or equivalently, bags of size one).

Motivated by the above reasoning, we further evaluate the classification performance on *isp* trace when the bag-size is varied. We set the number of initial training instances per class to be five. As bags of smaller sizes (say, with one flow) are inevitable, we evaluate the accuracies of the classifiers by varying the minimum bag-size in this set of experiments. Note that, when the minimum bag size is one, bags of all sizes are considered for testing; whereas if the same is set to three, bags of sizes one and two are ignored. The upper bound on the bag-size is set fixed to 20.

Fig. 7 shows the experiment results. The X-axis marks the minimum bag-size for each set of experiments. From the figure we make the following observations:

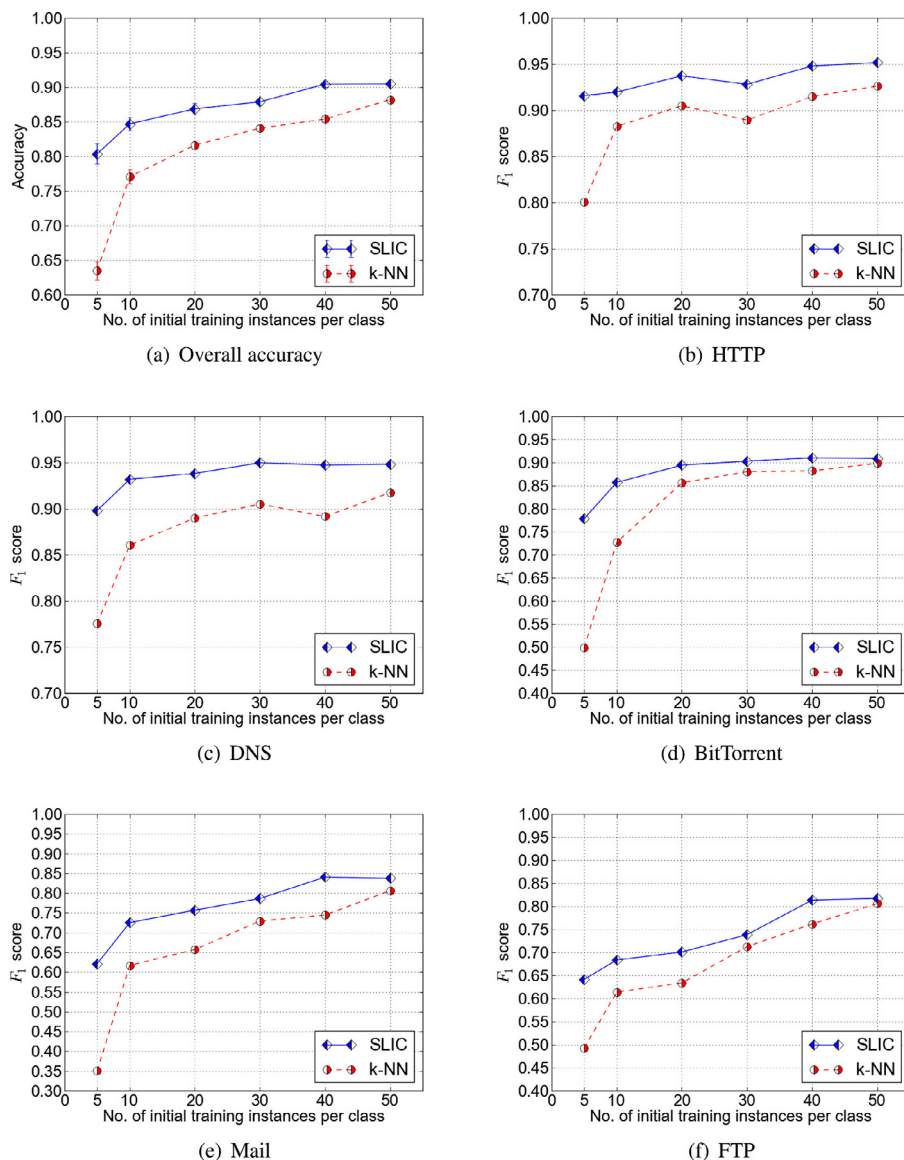


Fig. 6. Classification results on *wide* trace. Overall accuracy with 95% confidence interval is plotted in (a); (b)–(f) plot the F_1 scores for the different classes. The metrics are plotted against m , the number of training instances (flows) per application class.

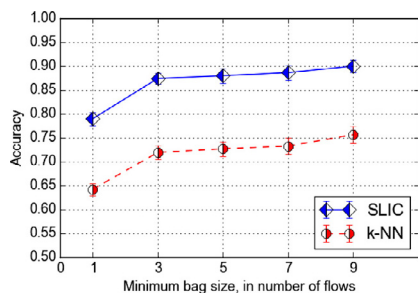


Fig. 7. Accuracy, when the minimum bag-size is varied.

- Regardless of the lower bound, SLIC consistently achieves significant improvement in accuracy, compared to the BoF-based k -NN classifier – SLIC achieves 20% or higher accuracy. Recall that the number of initial training instances is only five per class.
- With a lower bound of one for the bag-size, the accuracy achieved by SLIC is 79%. This is higher than the best accuracy obtained by k -NN (75%, when the minimum bag-size is nine). This again shows the strength of SLIC: instead of setting a high lower bound and missing out flows in small-sized bags in order to achieve better accuracy, SLIC, with the lowest bound of one, ignores no flow, and yet achieves high accuracy. This makes SLIC a practical on-line classifier, as it does not have to wait indefinitely, or even for long time, to be able to predict the class labels accurately.

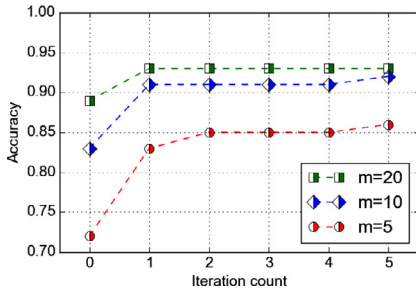


Fig. 8. Improvement in accuracies over iterations, for SLIC.

- Bags of sizes one and two indeed affect the classification results of both the classifiers. We observe a decrease of roughly 10% for SLIC, and 12% for k -NN between when the minimum bag-size is reduced from three to one. This is the reason for the diminished improvement for *wide trace* in comparison to the *isp trace*.

Number of iterations (η): The number of iterations η was set to five to allow enough iterations for SLIC to learn and stabilize. Fig. 8 plots the improvement in accuracies per iteration over the validation datasets of *isp trace*. These plots are for different sizes of the initial training set (5, 10 and 20). The accuracy for iteration zero is obtained using the initial training set, and without any self-learning. Observe that SLIC stabilizes quite quickly, achieving significant improvement in the first iteration.

Classification time: We sampled the time taken to classify flows, for both k -NN and SLIC. The experiments were performed on a desktop machine running 64-bit Intel Core 2 Quad CPU Q9550 @2.83GHz, with 8GB of RAM. The mean time to classify a flow for the BoF-based k -NN is 0.337 ms; whereas SLIC took a mean time of 0.358 ms to classify a flow. As SLIC does additional task of decision making in comparison to k -NN, the higher classification time for SLIC is expected. For the sampled experiments, we observed that the relative increase in time for SLIC was 6.4%, which is acceptable.

7. Conclusions and future directions

This work proposed a new direction to solve the two fundamental issues in supervised traffic classification: (i) learning using small number of labeled data, and (ii) adapting to non-static network static. We proposed and developed SLIC, a classification system that starts with a small set of labeled data and self-learns by inducting selected test instances (without ground truth) into its training set and rebuilds the model for classification. In short, SLIC evolves on its own and achieves significant improvement over the initial prediction accuracy achieved without self-learning. We demonstrated the effectiveness of SLIC by evaluating it against a BoF-based k -NN on two real-world traffic traces.

This is our first step towards developing self-learning traffic classification solutions. With this work, we see a number of interesting directions to advance the research on self-learning classifiers. We discuss a few important ones here. As SLIC is built on k -NN classifier, the size of the dynamic training set is constrained by computational complexity. This

means, if SLIC is to be continuously learning, the growing training size would slow down the classifier. There are two ways to overcome this limitation. One way is to continuously prune the dataset to remove obsolete instances from the training set, thereby making space for more inductions; and this process can be repeated infinitely. As long as the size of the dynamic training set is large enough to allow instances of a class to be a proper representation of the population, this solution will be effective. We envision such a pruning process will allow SLIC to achieve its best performance with an initial training set of just five instances per class. A second interesting approach is to develop self-learning over a parametric classifier such as Naive Bayes. As the computational complexity of Naive Bayes classification is independent of the size of the training set, the self-learning system has freedom to induct any number of good candidates. Another potential research direction is to adapt SLIC to detect and classify even unknown application protocols. Detection of zero-day protocols is not possible using the current supervised approach, but the core self-learning concept of SLIC can be used to develop a new solution.

Acknowledgment

We sincerely thank Dr. Jun Zhang and his research team from Deakin University, Australia, for generously sharing their labeled traffic traces with us for performance evaluation.

Appendix A. Experimenting with different values of k

Fig. 9(a) and (b) plot the accuracies for k -NN and SLIC, respectively, for different values of k used. We highlight that, for each experiment and a given value of k , both k -NN and SLIC started with the same set of initial training data. For example, consider the experiments for an initial training set size of five per class, and value of k set to nine. Here, using the same initial training set as k -NN, SLIC achieved an accuracy of 73%, while k -NN achieved only 52% accuracy. Though the trend of decreasing accuracy with increasing k is same with both classifiers, the accuracy achieved by k -NN decreases drastically; for $m = 5$, the accuracy reduced from 80% to a mere 52%, when the value of k was increased from one to nine. Accuracy of SLIC, on the other hand, decreases much more gradually. The differences in accuracies of SLIC is more prominent for experiments with smaller initial training set. Among those values of k that allow voting to decide the class of an instance, we observe that $k = 3$ performs the best.

Appendix B. Upper bound of bag-sizes

Both BoF-based k -NN and SLIC classifies bag-of-flows (BoFs) into application classes. The time to prediction is largely governed by the time the classifier needs to wait before the bag-of-flow to be formed. In real deployment, the system waits for a given interval of time during which the flows are processed and formed into bags. The longer the interval, the larger the size of the bags. However, a shorter interval reduces the time to classify the bags aggregated in the interval. Therefore, a shorter interval is preferred if the smaller-size bags do not affect the accuracy considerably. We

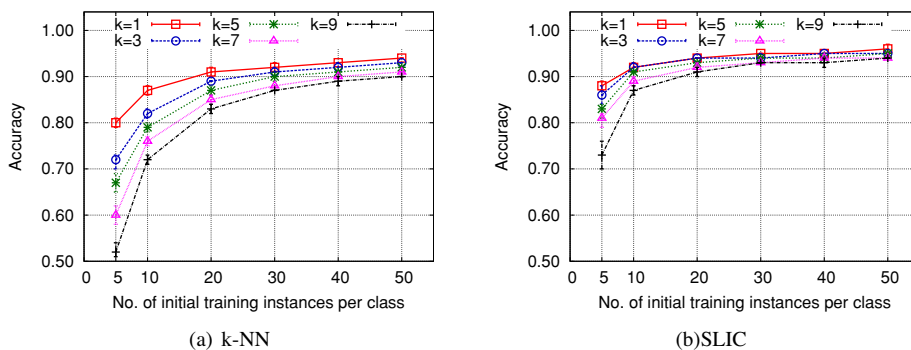


Fig. 9. Classification accuracies for different values of k .

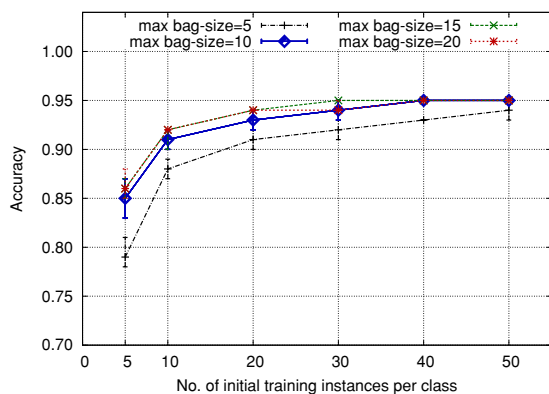


Fig. 10. Classification accuracies, when upper bound of bag-size range is varied.

conducted experiments to find the right value for the upper bound of the bag-sizes. Fig. 10 plots accuracy for different values of maximum bag-size; we fixed the lower bound to three.

The smallest value for the maximum bag-size we tested is five; and from the figure, we observe that the accuracy is least for this setting. The accuracies achieved for higher values of upper bound – 10, 15 and 20 – are very close with overlapping confidence interval. Unless otherwise stated, among these three values, we choose the smallest value of 10 as the upper bound for the bag-sizes.

References

- [1] T. Nguyen, G. Armitage, A survey of techniques for internet traffic classification using machine learning, *IEEE Commun. Surv. Tutor.* 10 (4) (2008) 56–76.
- [2] Q. Yang, X. Wu, 10 challenging problems in data mining research, *Int. J. Inf. Technol. Decis. Mak.* 5 (4) (2006) 597–604.
- [3] Cisco Visual Networking Index: Forecast and Methodology, 2013–2018, 2014. http://www.cisco.com/c/en/us/solutions/collateral/service-provider/ip-ngn-ip-next-generation-network/white_paper_c11-481360.pdf.
- [4] J. Zhang, Y. Xiang, Y. Wang, W. Zhou, Y. Xiang, Y. Guan, Network traffic classification using correlation information, *IEEE Trans. Parallel Distrib. Syst.* (TPDS) 24 (1) (2013) 104–117.
- [5] J. Zhang, C. Chen, Y. Xiang, W. Zhou, Y. Xiang, Internet traffic classification by aggregating correlated naive Bayes predictions, *IEEE Trans. Inf. Forensics Secur.* (TIFS) 8 (1) (2013) 5–15.
- [6] A.W. Moore, 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*, 2005, pp. 50–60.
- [7] T. Auld, A. Moore, S. Gull, Bayesian neural networks for internet traffic classification, *IEEE Trans. Neural Netw.* 18 (1) (2007) 223–239.
- [8] D. Bonfiglio, M. Mellia, M. Meo, D. Rossi, P. Tofanelli, Revealing skype traffic: when randomness plays with you, *ACM SIGCOMM Comput. Commun. Rev. (CCR)* 37 (4) (2007) 37–48.
- [9] H. Kim, K. Claffy, M. Fomenkov, D. Barman, M. Faloutsos, K. Lee, Internet traffic classification demystified: myths, caveats, and the best practices, in: *Proceedings of the ACM CoNEXT*, 2008, pp. 11:1–11:12.
- [10] M. Roughan, S. Sen, O. Spatscheck, N. Duffield, Class-of-service mapping for QoS: a statistical signature-based approach to ip traffic classification, in: *Proceedings of the 4th ACM SIGCOMM Conference on Internet Measurement (IMC)*, 2004, pp. 135–148.
- [11] J. Zhang, C. Chen, Y. Xiang, W. Zhou, A. Vasilakos, An effective network traffic classification method with unknown flow detection, *IEEE Trans. Netw. Serv. Manag.* 10 (2) (2013) 133–147.
- [12] A. McGregor, M.A. Hall, P. Lorier, J. Brunskill, Flow clustering using machine learning techniques, in: *Proceedings of the 5th International Workshop Passive and Active Network Measurement (PAM)*, 2004, pp. 205–214.
- [13] S. Zander, T. Nguyen, G. Armitage, Automated traffic classification and application identification using machine learning, in: *Proceedings of the IEEE Conference on Local Computer Networks (LCN)*, 2005, pp. 250–257.
- [14] L. Bernaille, R. Teixeira, I. Akodkenou, A. Soule, K. Salamatin, Traffic classification on the fly, *ACM SIGCOMM Comput. Commun. Rev. (CCR)* 36 (2) (2006) 23–26.
- [15] J. Erman, M. Arlitt, A. Mahanti, Traffic classification using clustering algorithms, in: *Proceedings of the SIGCOMM Workshop on Mining Network Data (MineNet)*, 2006, pp. 281–286.
- [16] L. Bernaille, R. Teixeira, K. Salamatin, Early application identification, in: *Proceedings of ACM CoNEXT*, 2006, pp. 6:1–6:12.
- [17] Y. Wang, Y. Xiang, S. Yu, An automatic application signature construction system for unknown traffic, *Concurr. Comput. Pract. Exp.* 22 (13) (2010) 1927–1944.
- [18] J. Erman, A. Mahanti, M. Arlitt, I. Cohen, C. Williamson, Offline/realtime traffic classification using semi-supervised learning, *Perform. Eval.* 64 (9–12) (2007) 1194–1213.
- [19] X. Li, F. Qi, D. Xu, X. song Qiu, An internet traffic classification method based on semi-supervised support vector machine, in: *Proceedings of IEEE International Conference on Communications (ICC)*, 2011, pp. 1–5.
- [20] Y. Wang, Y. Xiang, J. Zhang, S. Yu, A novel semi-supervised approach for network traffic clustering, in: *Proceedings of 5th International Conference on Network and System Security (NSS)*, 2011, pp. 169–175.
- [21] L. Grimaudo, M. Mellia, E. Baralis, R. Keralapura, SeLeCT: self-learning classifier for internet traffic, *IEEE Trans. Netw. Serv. Manag.* 11 (2) (2014) 144–157.
- [22] M. Oppor, A Bayesian approach to on-line learning, in: D. Saad (Ed.), *On-line Learning in Neural Networks*, Cambridge University Press, 1998, pp. 363–378.
- [23] S.olla, O. Winther, Optimal perceptron learning: as online Bayesian approach, in: D. Saad (Ed.), *On-line Learning in Neural Networks*, Cambridge University Press, 1998, pp. 379–398.
- [24] J. Ma, K. Levchenko, C. Kreibich, S. Savage, G.M. Voelker, Unexpected means of protocol inference, in: *Proceedings of the 6th ACM SIGCOMM Conference on Internet Measurement (IMC)*, 2006, pp. 313–326.
- [25] M. Canini, W. Li, M. Zadnik, A.W. Moore, Experience with high-speed automated application-identification for network-management, in: *Proceedings of the 5th ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS)*, 2009, pp. 209–218.

- [26] Y. Li, X. Zhang, Improving k nearest neighbor with exemplar generalization for imbalanced classification, in: Proceedings of the 15th Pacific-Asia Conference on Advances in Knowledge Discovery and Data Mining, PAKDD'11, 2011, pp. 321–332.
- [27] H. Dubey, V. Pudi, Class based weighted k -nearest neighbor over imbalance dataset, in: Proceedings of the 17th Pacific-Asia Conference on Advances in Knowledge Discovery and Data Mining (PAKDD), 2013, pp. 305–316.
- [28] J.A. Hartigan, M.A. Wong, Algorithm AS 136: a k -means clustering algorithm, Appl. Stat. 28 (1) (1979) 100–108.
- [29] K. Cho, K. Mitsuya, A. Kato, Traffic data repository at the WIDE project, in: Proceedings of the Annual Conference on USENIX Annual Technical Conference (ATEC), 2000, p. 51.
- [30] I7-filter, <http://i7-filter.clearfoundation.com>.
- [31] TCP Statistic and Analysis Tool, <http://tstat.tlc.polito.it/index.shtml>.
- [32] L. Breiman, Random Forests, Mach. Learn. 45 (1) (2001) 5–32.
- [33] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: machine learning in Python, J. Mach. Learn. Res. 12 (2011) 2825–2830.



Dinil Mon Divakaran is a Scientist at the Cyber Security & Intelligence department in the Institute for Infocomm Research (I²R), Singapore. Prior to this, for two years he was a Research Fellow at the Department of Electrical and Computer Engineering in the National University of Singapore (NUS). He has also worked as an Assistant Professor in the School of Computing and Electrical Engineering at the Indian Institute of Technology (IIT) Mandi. He carried out his Ph.D. at the INRIA team in ENS Lyon, France. He holds a Master degree from IIT Madras, India. His research interests and experiences are in the areas of network

security and QoS delivery. His research works revolve around applications of machine learning, game theory, queueing theory and probabilistic models, as well as the study of optimization problems and design of heuristics, all in the broad area of computer networks.



Le Su received the B.Sc. degree from Nanyang Technological University, Singapore in 2010, majored in Mathematical Sciences. He completed the Ph.D. degree from the same school in 2014 with research focus on the area of applied cryptography, especially on pairing-based cryptographic scheme design. He is now working as a Research Scientist in Cyber Security & Intelligence department, Institute for Infocomm Research, A*STAR, Singapore. His current research interests include network security, machine learning, privacy-preserving technique, security intelligence, and applied cryptography.



Yung Siang Liao received the B.Eng. degree in Electrical Engineering from the National University of Singapore, Singapore in 2012. He is currently a research engineer with the Cyber Security and Intelligence Department at the Institute for Infocomm Research (I²R), Singapore. His current research interests include security intelligence and machine learning.



Vrizlynn Thing leads the Cyber Security & Intelligence (CSI) R&D Department at the Institute for Infocomm Research, A*STAR, Singapore. The department focuses on digital forensics, cybercrime analysis, cyber security & intelligence, and mobile security research and technology development. She is also an Adjunct Associate Professor at the National University of Singapore (School of Computing), and the Singapore Management University (School of Information Systems). She has over 13 years of security and forensics R&D experience with in-depth expertise in various aspects of cyber security, forensics & intelligence. Her research

draws on her multidisciplinary background in computer science (Ph.D. from Imperial College London, United Kingdom), and electrical, electronics, computer and communications engineering (B.Eng. and M.Eng. by Research from Nanyang Technological University, Singapore). During her career, she has taken on various advisory and consultancy roles for national level cyber security initiatives. She also participates actively as the Lead Scientist of collaborative projects with industry partners such as MNCs, SMEs and the government agencies. More information about her work can be found at www.i2r.a-star.edu.sg/vriz.