

# Behavior-Based Network Traffic Synthesis

Yingbo Song, Salvatore J. Stolfo, and Tony Jebara

**Abstract**—Modern network security research has demonstrated a clear necessity for open sharing of traffic datasets between organizations - a need that has so far been superseded by the challenges of removing sensitive content from the data beforehand. Network Data Anonymization is an emerging field dedicated to solving this problem, with a main focus on removal of identifiable artifacts that might pierce privacy, such as usernames and IP addresses. However, recent research has demonstrated that more subtle statistical artifacts may yield fingerprints that are just as differentiable as the former. This result highlights certain shortcomings in current anonymization frameworks; particularly, ignoring the behavioral idiosyncrasies of network protocols, applications, and users. Network traffic synthesis (or simulation) is a closely related complimentary approach that, while more difficult to execute accurately, has the potential for far greater flexibility. This paper leverages the statistical-idiosyncrasies of network behavior to augment anonymization and traffic-synthesis techniques through machine-learning models specifically designed to capture host-level behavior. We present the design of a system that can automatically learn models for network host behavior across time then use these models to replicate the original behavior, to interpolate across gaps in the original traffic, and demonstrate how to generate new diverse behaviors. Further, we measure the similarity of the synthesized data to the original, providing us with a quantifiable estimate of data fidelity.

## I. INTRODUCTION

Open access to quality data is the foundation upon which nearly all modern engineering research disciplines are built upon. The availability of openly distributed and standardized datasets is a prerequisite for scientific measurements, verifications, and comparisons of existing technologies. In the field of networking research, however, creating publicly releasable datasets is met with many hurdles. First and foremost, as an information-exchange medium, network data contains the personal communications of individuals, and use such data carries privacy concerns. Furthermore, the scale of typical enterprise networks implies a scale of thousands, if not tens of thousands of users, which has in the past made artificial simulation difficult to achieve.

This paper presents a new method for synthesizing network traffic that can quantifiably mimic real-user behavior. While different types of network traffic generators have been revealed in the past, each serving distinct roles, to the best of our knowledge, ours is the first system which was designed to agnostically mimic user behavior by automatically *learning* behavioral profiles per user, then using these profiles to synthesize new traffic, in a way that is quantifiably similar to the original. In this manner, our system learns to mimic the behaviors of a network as a whole based upon individual members. This work is an extension of our recent results on behavior profiling and anonymization using statistics-preserving anonymity by crowds [1].

Y. Song is with the Department of Computer Science, Columbia University, New York, NY 10027-7003. Phone: +1 212-939-7078, E-mail: yingbo@cs.columbia.edu

S. J. Stolfo is with the Department of Computer Science, Columbia University, New York, NY 10027-7003. Phone: +1 212-939-7080, E-mail: sal@cs.columbia.edu

T. Jebara is with the Department of Computer Science, Columbia University, New York, NY 10027-7003. Phone: +1 212-939-7079, E-mail: jebara@cs.columbia.edu

In the networking and security research community, there is a significant need for publicly available research data. A 2008 survey by Mirkovic showed that out of a total of 144 papers published in Special Interest Group on Data Communication (SIGCOMM) and Internet Measurement Conference (IMC) in 2006 and 2007, 49 of these efforts had utilized network traces in their evaluations, but only 10 had used publicly available datasets [2]. The cause was cited to be a lack of high quality standardized and openly available data. It is because of this need that organizations such as OpenPacket [3], and the more recent U.S. Department of Homeland Security-sponsored PREDICT [4], were recently created - to facilitate the sharing of information between research organizations. However, sharing network data has proven to be difficult, at best. A packet capture of all network traffic, for example, would include web traffic showing which websites users visited, where they transfer files to and from, the locations of their email, banking, and other private accounts, as well as any credentials not protected by encryption. In addition to personal information, the disclosure of network profiles such as vulnerability fingerprints in existing machines, firewall policies, details of existing security services, location of database and other sensitive servers, and network infrastructure in general, can all lead to unintended negative consequences for the releasing party. Network trace anonymization (NDA) is an emerging field that is dedicated to solving this problem. NDA technologies aim to facilitate the exchange of authentic network data by providing the methods by which releasing parties can remove any potentially sensitive artifacts from the data beforehand. Ideally, the anonymized data can then be shared without fear of breaching privacy. Network traffic *synthesis*, or simulation, is a closely related complimentary approach which is more difficult to execute but has the potential for greater flexibility.

In a recent paper, we proposed an anonymization system that utilized machine-learning to train behavioral profiles for network hosts. Then, by inter-mixing traffic among similarly behaving members, we derive a pseudo statistical-mixed-net that provided anonymity-by-crowds, while preserving the statistical properties within the data [1]. In this paper, we extend this work to demonstrate how, given these behavior profiles, synthetic network data can be generated that mimic the properties of the original traffic. We show how this type of synthetic data can be used for many purposes, such as producing realistic network traffic, reshaping behavior, augmenting anonymization technologies, as well as simulating networks. In addition to anonymization, this technology has potential use in areas such as network measurement, and simulation environments such as the DARPA National Cyber Range (NCR) project [5].

The main novelty of our paper is the design of a system that can automatically learn models for network host behavior across time, then use this model replicate that behavior. We show how to manipulate the models to generate diverse behaviors. Further, we can measure the similarity of the new data to the original, providing us with a quantifiable estimate of data fidelity.

## II. BACKGROUND AND RELATED WORKS

Network-data generation has been an interesting research topic for some time, for many reasons. Most focus on network benchmarking, routing testing, firewall, IDS, and other security related testing. As recently as 2008, great interest has emerged in the network traffic anonymization (NDA) field which studies the intersections between the needs of network measurement and security community balanced with the privacy issues with using data corresponding to human users. This section briefly covers an overview of related topics in this specialized area, and describes our novel approach in building systems which automatically models real-user activity and mimics this in synthesized data.

The 2008 survey by Mirkovic [2] and other similar studies [6], make clear the existence of a deficiency for publicly available large network-traffic datasets. This need has resulted in the introduction of organizations such as the U.S. Department of Homeland Security-sponsored PREDICT [4] which aim to facilitate the distribution of network data for research. Unlike data used in other disciplines, raw network traffic data often include sensitive information that cannot be released without breach of user privacy. A packet capture of all network traffic, for example, would include web traffic showing which websites users visited, where they transferred files to and from, the locations of their e-mail, banking, and other private accounts, as well as any credentials not protected by encryption. In addition to personal information, the disclosure of network profiles such as vulnerability fingerprints in existing machines, firewall policies, details of existing security services, location of database and other sensitive servers, and network infrastructure in general, can all lead to unintended negative consequences for the releasing party.

In a recent paper, we proposed an anonymization system that utilized machine learning to learn unique behavioral profiles for each host on a network, then by inter-mixing the traffic amongst similarly behaving hosts, we obtain a pseudo statistical mixed-net which provides anonymity by crowds, while preserving the useful statistical properties of the individuals. In this complementary paper we study the topic of statistical-behavior-based traffic synthesis for network data generation. We augment modified data with synthetic models, and aim to further improve NDA by fuzzing the original data to dilute it of sensitive information.

The most well known examples in artificial traffic generation come from the network measurement disciplines. Often, network engineers need to test things such as routing paths, load balancer implementation, firewall policies, and appliances used to maintain the network. In these disciplines the goal is to be able to generate a wide range of data at a fast pace. Breakpoint Systems [7] and Spirent [8] are examples of this role. Traffic generators from these companies can send 40 Gb/s of data at a router, hitting every port, using every protocol, setting every TCP/IP flag, and any enumerations thereof. The second most commonly observed need for traffic generation is for security testing. Some of the previously mentioned systems are capable of generating exploit traffic to test IDS implementations under heavy load. This area is filled with specialist tools that test specific security problems. These tools generate traffic as sort of a

by-product to their main intended use, and in that sense they can be considered malicious traffic generators, to an extent. Well known tools in this area are NMap port scanner [9], the Nessus web vulnerability scanner [10], the Metasploit exploitation framework [11], and the list goes on for very long.

In the area of precision network problem diagnosis, specialized packet crafting tools include hping [12] and Nemesis [13]. Ostinato is a GUI-based packet assembler that aims to work like Wireshark [14] in reverse. Network Expect [15] is a powerful tool that allows one to craft and manipulate complicated TCP/IP sessions using a high level scripting language. The difference between these methods and ours is that they serve to generate specific classes of data for specific roles and do not attempt to mimic existing network traffic behavior and synthesize data for the purpose of simulation.

In 2008, a broad agency announcement was sent out by DARPA describing the goal to create a national testing platform for cyber-warfare technologies. From this emerged the National Cyber Range (NCR) project [5], a multi-million dollar effort to build a system which simulates a pseudo-Internet with thousands of nodes sending and receiving realistic looking traffic. Our proposed method fits well into this category as a technology which can be deployed on an actual network and used to automatically learn host behaviors, then the statistical models are transferred to a separate network where the synthesis algorithms mimic the network hosts that they have learned to emulate for network simulation purposes.

## III. NETWORK TRAFFIC MODELS

Given a packet capture from an enterprise network, it's very likely that one can glean content signatures that uniquely identify specific individual hosts. For example, by observing web traffic we might notice specific usernames being transmitted from a machine to a foreign web server. A collection of such signatures can distinguish one user from another and help form a content-based fingerprint for that user. A related question would then be, is it possible to obtain a signature for a user based solely on *statistical* information (or "behavior") represented by meta-data: by tracking rates of packet exchanges, frequencies, and port activities? In a recent paper describing our IMPACT anonymization framework [1] we demonstrated that this is indeed possible. With significant accuracy we can identify users based entirely on this meta-data, without payload content. In this prior work our algorithms were used to cluster similarly behaving hosts in an offline transform that simulates the effects of a statistics-preserving anonymous VPN with multiple in and multiple out aggregation points. In this paper, we extend our models to perform data generation for simulation. Here, our focus is on synthesizing network data that is faithful to the trained behavioral models, in order to replicated targeted behavior. Our algorithm operates on data in the form of Netflow-like statistics. Specifically, we need the following information:

{TIME, SRC, SRCPORT, DST, DSTPORT, #PKTS}

A record for a user's behavior activity would contain hundreds to thousands of these entries, representing behavior throughout the day. Given the port numbers we can estimate which services the user executed. This information is easily

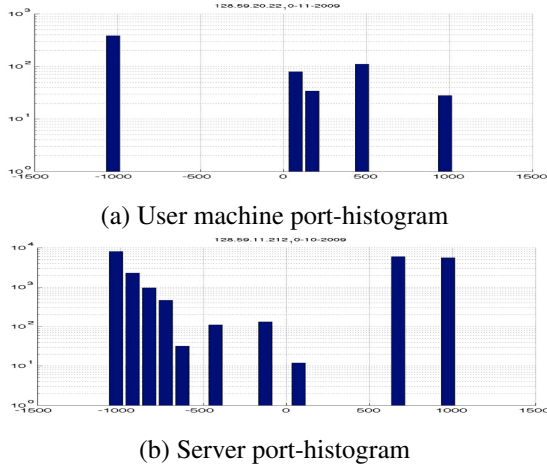


Fig. 1. Log-scale port histograms. The first half of the x-axis ( $x < 0$ ) represents Received-On ports; the second, Sent-To ports. Notice the distinct behavioral patterns.

extracted from Cisco Netflow records as well as packet traces. Measurable traffic data reflect a complex interaction of different inputs throughout the many layers of the TCP/IP protocol stack, as well as the state of the machine: a user clicks on a URL on a webpage, this causes his browser to dispatch a message to a foreign server in the form of an HTTP request, which is transmitted in a TCP stream, that is broken into independent IP datagrams, which are encapsulated in Ethernet frames, and then passed to the local gateway for routing. At each of these stages/layers, measurements of quantifiable behavior, such as volume and velocity of the packets, reflects user behavior as well as implicitly taking into account the state of the network stack and the operating system. In practice, the two most used all-purpose logging methods are packet captures and Netflow captures. It is for these reasons that we have chosen this statistical representation to build our behavior models upon.

### A. Statistical representations of behavior

Figure (1) shows the emergence of distinguishable behavior when examining TCP statistics. This plot shows the log-scale plot of the port-traffic histogram of two difference machines on Columbia’s network, using traffic captured across the span of one day. Here, we see a user machine making requests to various service ports on a foreign machine, and receiving responses back on its ephemeral port. The server machine on the right receives several distinct types of incoming requests and makes three main types of outbound connections: a low-port connection, most likely SSH, an outbound request for a port in the range of 600, and replies on ephemeral ports (in this figure, ephemeral ports are coalesced together on port 1025). Our statistical model is time-series analogue of this behavioral snapshot, one that tracks this system-behavior *across time*. Given sufficient data sampling, our results show that very distinct *statistical fingerprints* of host-identity emerge. This reflects the natural behavioral idiosyncrasies that lay hidden in plain-sight within network data.

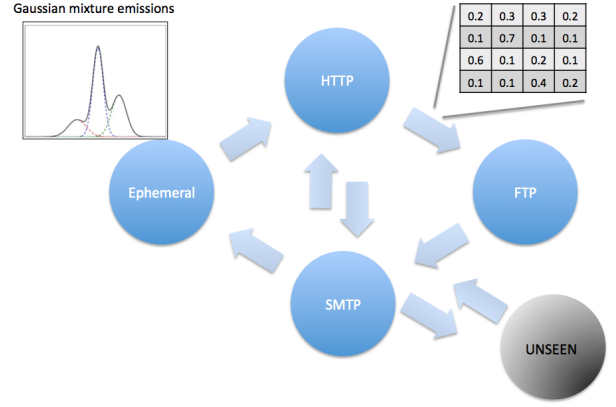


Fig. 2. Protocol-transition Markov model diagram.

### B. Protocol-Transition Markov Model

This section provides a concise description of our model and focuses on our extensions. We refer the reader to our previous paper for more detailed background and derivation details [1]. The PT-Markov model tracks the dynamics of a system as a transition between specific application services represented by distinguishable protocols ( $FTP \rightarrow HTTP \rightarrow SMTP, \dots$ ). Let  $s_t$  represent the “state” of the system at time-step  $t$ ; representing the active service at that time. Let  $x_t$  represent the distribution on the volume of data observed at that time-step. The overall interaction is modeled as a single-step Markov model:

$$p(x|s, \theta) = p(s_1)p(x|s_1) \prod_{t=2}^T p(s_t|s_{t-1})p(x_t|s_t). \quad (1)$$

In the above equation,  $p(s_t|s_{t-1})$  models the likelihood of transitions between states (services) and  $p(x|s)$  is the “emissions” model which represents the volume of data observed when the system was in that state. For this volume distribution, we use a scalar Gaussian mixture model of form:

$$p(x|\theta) = \sum_{i=1}^M \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left\{ -\frac{(x - \mu_i)^2}{2\sigma_i^2} \right\}. \quad (2)$$

A mixture model allows us to implicitly capture an estimate of the different types of traffic exchanges observed under a particular protocol. For example, the model state for port 80 would track a distribution for the size of websites visited; the state for port 25 would track a distribution for the lengths of emails sent, and so forth. The above emissions model can be easily estimated using Expectation Maximization (EM). Since the value of the state is known, estimation of the state-transition probability table  $T(a, b) = p(s_t|s_{t-1})$  consists of simply tracking the transition-frequencies as follows:

$$T(a, b) = \frac{\#(s_i = b, s_{i-1} = a)}{\#(s_i = b, s_{i-1} = a) + \#(s_i \neq b, s_{i-1} = a)}. \quad (3)$$

Here,  $\#(s_i = b, s_{i-1} = a)$  is a function that counts the number of times state  $b$  follows state  $a$  in the training data (how many times a user transitions from FTP to HTTP, for example).

Putting these together we obtain the full probability distribution  $\mathbf{s} = \{s_1, \dots, s_T\}$ ,  $\mathbf{x} = \{x_1, \dots, x_T\}$  and  $\Theta = \{\theta_1, \dots, \theta_M\}$ :

$$p(x_0, s_0) = 1 \quad (4)$$

$$p(\mathbf{x}|\mathbf{s}, \Theta) = \prod_{t=1}^T T(s_{t-1}|s_t) \sum_{i=1}^M \pi_{s,i} \mathcal{N}(x_t|\theta_{s,i}). \quad (5)$$

$\pi_{s,i}$  is the mixing proportion for the GMM emission model for state  $s$ , and  $\sum_i \pi_{s,i} = 1, \forall s$ . Further, for each emissions model of any given state (whose parameter is represented by  $\theta_i$ ), we need to track table of source and destination port pairings. Let  $b_{i,j} = p(\text{dst} = j|\text{src} = i)$ , this tracks the distribution of source ports which sent traffic to the destination port represented by this state. Table  $\mathbf{b}$  is estimated in the same manner as  $T(a, b)$ . This gives us a final set of parameters that make up the model  $\theta = \{\mathbf{b}, \pi_1, \mu_1, \sigma_1, \dots, \pi_M, \mu_M, \sigma_M\}$ .

### C. Statistical features and implementation details

If we track transition of services based on port values then, in theory, we could have a unmanageable  $65,536 \times 65,536$ -sized transition matrix. In practice, this is not the case given the design of TCP. There are three types of ports defined in TCP, referred to as the “well known,” “registered,” and “ephemeral” ports. Each type of port has its own predetermined range. “Well-known” ports extend from port 0 to port 1024 and are used by very common network services, such as SSH. The “registered” range extends from 1,025 to 49,151 and is used by third party applications; several Bittorrent clients uses port 6900, for example. Finally all other services with no registered ports uses a randomly selected port from the “ephemeral” range of 49,152 to 65,536. Given this implementation design, two things become apparent. First, since the “well-known” range tend to experience more stable and dense traffic activity it is best to maintain a one-to-one feature mapping. That is, features relating to SSH should not be combined with features relating to FTP. Second, since the “ephemeral” range maintains no correlation between the specific service and the port used, we can collapse all activity in this range into a single bin.

Feature extraction for the “registered” range is not well defined for several reasons. First, TCP/IP implementations are not consistent across all platforms and overlapping port registration does occur. In addition, different operating systems will allocate ports differently based on implementation, version, and the third-party application behavior. Third, it is not a strict requirement for services to use the ports that they are registered to. In our work, we use a histogram to represent this range and bin services based on their numerical proximity. If we have  $49,151 - 1,025 = 48,126$  registered ports, and a 100-bin histogram, we would accumulate activity for all ports between 1025-1124 into the first bin of the histogram, 1125-1224 in the second bin, and so on. The results in this paper were derived from models using 20-bin histograms; chosen because port activity is typically very sparse, absent the presence of port-scanning activities. Most hosts typically use only a handful of services (such as SMTP and HTTP); these reside mostly in the “well known” range. Sparsity is particularly pronounced in the “registered” range. Using large histograms would induce higher

training-data requirements on the machine-learning algorithms and makes accurate parameter estimation more ill-posed.

When measuring volume, network traffic can fluctuate dramatically for particular hosts even when measuring the same services. This is due to the exchange of protocol-control and data messages, which can induce a large variance in the statistical model and reduce modeling accuracy. To compensate, we use a log-squashing function on the volume feature. Let  $x$  represent the volume of activity that a particular service observes for a given session, instead of measuring  $x$  in our models we measure  $\log(x)$ . Finally, our model tracks the distribution of port pairings. For each destination port we track an independent set of associated source ports from where we have observed traffic. A problem arises when a state is encountered in testing that was not seen in training, for example the appearance of HTTP traffic where only SMTP was observed during training, then no entry for that protocol would exist in the trained model. This is solved by adding a single dummy state to represent all protocols unseen during training, with a fixed-value set for both entry- and exit-transition and emission probabilities. Full evaluations which demonstrate that our model can distinguish between hosts with upwards of 80% accuracy are provided in [1].

## IV. SYNTHESIS

This section describes our algorithm for synthesizing Netflow from our statistical model and describe our extensions into packet-trace synthesis. We describe the technical details of our algorithm and provide full pseudocode for our functions.

### A. Synthesizing Netflow

A few simple mathematical notations needs to be introduced before describing the system. Let  $s \sim \mathcal{M}(p_1, p_2, \dots, p_k)$  denote a random draw from a multinomial distribution specified by the normalized ratios parameter  $\mathbf{p} = \{p_1, p_2, \dots, p_k\}$ . Meaning, if we had an alphabet of  $k$  characters with the respective probabilities of appearance  $\mathbf{p}$  where  $\sum_i p_i = 1$ , as the number of draws grows the normalized proportion of selected values approaches  $\mathbf{p}$ . Let  $x \sim \mathcal{N}(\theta_i)$  denote sampling a scalar quantity from the Gaussian mixture with parameter  $\theta_i = \{\mu_{i,1}, \sigma_{i,1}, \dots, \mu_{i,m}, \sigma_{i,m}\}$  for mixture size  $m$ . The pseudo code for our synthesis algorithm is given as follows:

SYNTHESIZE-NETFLOW( $\theta, n$ )

1 **return** FEATURES-TO-DATA(SYNTHESIZE-STATS( $\theta, n$ ))

SYNTHESIZE-STATS( $\theta, n$ )

```

1  ▷ Use the frequency estimate to set the initial state
2   $d_0 \sim \mathcal{M}(\pi_1, \pi, \dots, \pi_k)$ 
3  ▷ Src/Dst port model used to estimate the pairing
4   $s_0 \sim \mathcal{M}(b_{d_0,1}, b_{d_0,2}, \dots)$ 
5   $x_0 \sim \mathcal{N}(\theta_{s_0})$            ▷ Sampling the volume
6  ▷ The rest is generated by induction
7  for  $t \leftarrow 1$  to  $n - 1$  :
8      do  $d_t \sim \mathcal{M}(a_{d_{t-1},1}, a_{d_{t-1},2}, \dots, a_{d_{t-1},k})$ 
9           $s_t \sim \mathcal{M}(b_{d_t,1}, b_{d_t,2}, \dots)$ 
10          $x_t \sim \mathcal{N}(\theta_{s_t})$ 
11 return [ $\mathbf{s}, \mathbf{d}, \mathbf{x}$ ]
```

DATASET	GENUINE	SYNTHETIC	RANDOM SAMPLED
WEST POINT BORDER	-4.1445E+10	-6.6203E+10	-8.3211E+10
NSA CAPTURE	-5.2935E+10	-6.5322E+10	-6.6581E+10
LBNL ENTERPRISE	-3.9719E+10	-6.8987E+10	-8.8169E+10
COLUMBIA UNIV.	-3.0717E+10	-3.1106E+10	-5.5863E+10
GEORGE MASON UNIV.	-1.3528E+10	-3.5596E+10	-3.7741E+10
UMICH. MERIT	-2.0368E+10	-4.1504E+10	-5.2505E+10

TABLE I

LOG-LIKELIHOOD VALUES OF GENUINE, SYNTHETIC, AND RANDOMLY SAMPLED DATA. AVERAGE OF 10 TRIALS REPORTED.

Our synthesis method uses the service (port) transition probabilities, learned during training, as parameters in a multinomial distribution. This is then used to sample the sequence of service states in a traffic stream. This gives us a maximum-likelihood transition between destination ports, in a series. For each destination element, recall that we have a separate model for the distribution of source ports in the  $\mathbf{b}$  parameter. Given the strong correlation between source and destination ports, joint-modeling for source and destination pairings is redundant. That is, we do not need to explicitly model  $p(s_i, d_i)$  as opposed to our current  $p(s_i|d_i)$ . Finally, volume information  $x_i$  is drawn from the Gaussian mixture parameter which is unique to each service model  $i$ . Function SYNTHESIZE-STATS simulates traffic patterns using our feature representation. This result is then converted back to actual port values by the following function.

FEATURES-TO-DATA( $\mathbf{s}, \mathbf{d}, \mathbf{x}$ )

```

1  [ $\mathbf{s}', \mathbf{d}', \mathbf{x}'$ ]  $\leftarrow \emptyset$ 
2   $n \leftarrow \text{LENGTH}(\mathbf{s})$ 
3  for  $t \leftarrow 1$  to  $n$  :
4      do  $s'_t \leftarrow \text{BIN-TO-PORTS}(s_t)$ 
5          $d'_t \leftarrow \text{BIN-TO-PORTS}(d_t)$ 
6          $x'_t \leftarrow \exp(x_t)$ 
7  return [ $\mathbf{s}', \mathbf{d}', \mathbf{x}'$ ]

```

BIN-TO-PORTS is a function that reverses the feature-extraction transformation performed on the dataset during training and remaps the data from features back to their normal representation. For example,  $d_i = 1045$  (a port-bin) then  $d_i$  gets remapped to a random ephemeral value between the range of [49152, ..., 65536].

BIN-TO-PORTS( $x, h$ )

```

1  if  $x < 1026$  :                                $\triangleright$  Well known ports
2      then return  $x$ 
3  elseif  $x > (1025 + h)$  :                        $\triangleright$  Ephemeral ports
4      then return  $\text{RAND}(1, \dots, 16384) + 49151$ 
5   $z \leftarrow (49152 - 1025)/h$                     $\triangleright$  Registered ports
6  return  $1025 + (x - 1025) \times h + \text{RAND}(1, \dots, h)$ 

```

### B. Quantifying Similarity

We quantify the fidelity of the synthesized data by evaluating the log-likelihood of the generated data over the trained models. A faithful reproduction would yield a higher log-likelihood score than a poorly synthesized sequence. The likelihood score is calculated using Equation (5) in the previous section.  $\log(p(\mathbf{s}, \mathbf{x}|\Theta))$  is used to avoid underflow errors. We com-

pare the log-likelihood scores of synthetic data with the original (genuine) data as well as data piece together by randomly selecting subsets of other hosts’s original traffic. The latter case, represented in the table by “Random” is meant to represent a naive way of simulating data. We chose to randomly sample real data vs generating completely random data as this gives us a more challenging performance baseline to compare against given that we are using considering real traffic. Table (I) shows this likelihood comparison. All results are average of 10 randomized trials. As expected, our synthesized data consistently exhibit higher loglikelihood than the randomly sampled data, but lower than the genuine set, which represents an upper bound. This means that our method of synthesizing data is more accurate than if one were to randomly piece together sessions of *actual* traffic from other hosts in the same dataset.

### C. Interpolation and Behavior Shaping

Behavior interpolation can be achieved by using a switched-Markov-model. In this case, multiple behavioral models are connected using a new external state variable. During the chain of data synthesis, this variable switches between multiple models probabilistically under a user-specific distribution. A simple method is to provide a ratio parameter and sample from a multinomial, as was done previously. SYNTHESIZE-MIXED-NETFLOW provides the pseudocode for this procedure.

SYNTHESIZE-MIXED-NETFLOW( $\Theta, n, \mathbf{p}, c$ )

```

1  [ $\mathbf{s}, \mathbf{d}, \mathbf{x}$ ]  $\leftarrow \emptyset$ 
2  for  $t \leftarrow 1$  to  $\lfloor n/c \rfloor$  :
3      do  $i \sim \mathcal{M}(\mathbf{p})$ 
4         [ $\mathbf{s}_t, \mathbf{d}_t, \mathbf{x}_t$ ]  $\leftarrow \text{SYNTHESIZE-NETFLOW}(\theta_i, c)$ 
5         [ $\mathbf{s}, \mathbf{d}, \mathbf{x}$ ]  $\leftarrow \text{APPEND}([\mathbf{s}_t, \mathbf{d}_t, \mathbf{x}_t])$ 
6          $t \leftarrow t + 1$ 
7  return [ $\mathbf{s}, \mathbf{d}, \mathbf{x}$ ]

```

Data synthesis using switched-models involves generating blocks of traffic of length  $c$  for each model, and combining them. The result is that individual sections of the data will be statistically similar to different profiles. In practice,  $c$  can be randomized at each iteration within the loop to avoid block patterns. One might consider why  $c$  should not be set as 1. The reason for this is that the PT-Markov model measures *transitional* probability. Therefore, if  $c = 1$ , in the worst case where we have two completely different behavior models, we can have a thrashing scenario where the model-switch occurs after each generated entry. This would, in theory, yield a result that is dis-

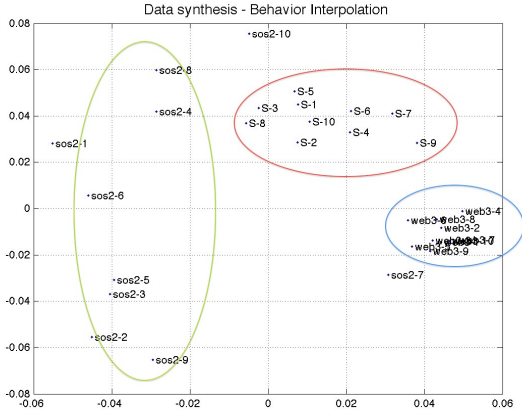


Fig. 3. Behavior interpolation: Two clusters of behavior profiles are embedded (“sos” and “web”) along with a new class of behavior (“S”) that interpolates between these two models.

similar to both of the initial models. This thrashing effect is mitigated if two models are similar to each other, meaning they share a certain amount of overlap in ports used, thus allowing the models to transition between each other more smoothly. In practice,  $c$  can be adjusted based on the model similarities.

We use a *visual* similarity-measurement technique to quantify the fidelity of the synthesized data in this interpolation setting. This is done by plotting points in a 2-D plane in a way such that each point corresponds to a synthesized data sequence, and the distances between points in this 2-D plane is proportional to the distances between the corresponding traffic samples in the original space. This technique is known as “embedding,” and we have developed methods for network-data embedding in our prior work [1]. In this experiment, we began with two classes of behavior: one modeled after a web server (denoted by “web”) and the other modeled after a general-computation server (denoted by “sos”). A switched model is used to synthesize samples of their interpolated behavior which we labeled the “S” class. Visually, we can confirm that the synthesized data indeed falls between these two distinct classes of behavior. The exact method which we used to map this abstract behavior manifold is known as Kernel Principle Component Analysis (KPCA). The kernel used, in this case, is a variation of the Probability Product kernel [16] derived specifically for the PT-Markov model. For more details we refer the reader to our other paper which describes this work in more detail [1].

#### D. Behavior Regression

In addition to behavior interpolation and shaping, we can use these models to perform data regression to fill-in missing segments of data. Ideally, we would like the synthesized data to match with the actual data along the boundaries so that transition probabilities are properly maximized. While this is not so difficult if we use a simple one-step Markov model as the basic setting of our model uses, problems would arise in the future if we were to extend our model and use a multi-step Markov dependency assumptions where we have  $p(s_t | s_{t-1}, s_{t-2}, \dots, s_{t-k})$ . In this case aligning the boundaries becomes more difficult. Given this we propose a more general

solution that can work for any variation on the behavior models. An efficient randomized algorithm is to simply generate many solutions and pick the one that yields the highest likelihood given the choice of probability model  $p(s, \mathbf{x} | \Theta)$ . This takes the form of the following algorithm:

NETFLOW-REGRESSION( $\theta, \mathbf{s}, \mathbf{d}, \mathbf{x}, c$ )

```

1   $n \leftarrow$  SOME LARGE NUMBER
2   $D \leftarrow [\infty, \infty, \dots, \infty]$  of length  $n$ 
3   $[\mathbf{s}', \mathbf{d}', \mathbf{x}'] \leftarrow$  SYNTHESIZE-NETFLOW( $\theta, n$ )
4  for  $t \leftarrow 4$  to  $n - c - 4$  :
5      do  $D_t = \sum_{i=t-3}^t \|d_i - d'_i\| + \sum_{i=t+c}^{t+c+3} \|d_i - d'_i\| \dots$ 
6           $+ \sum_{i=t-3}^t \|s_i - s'_i\| + \sum_{i=t+c}^{t+c+3} \|s_i - s'_i\|$ 
7       $j \leftarrow$  FIND-INDEX-OF(MIN( $D$ ))
8       $[\mathbf{s}^*, \mathbf{d}^*, \mathbf{x}^*] \leftarrow [s'_{j, \dots, j+c-1}, d'_{j, \dots, j+c-1}, x'_{j, \dots, j+c-1}]$ 
9  return  $[\mathbf{s}^*, \mathbf{d}^*, \mathbf{x}^*]$ 

```

Here a long sequence is synthesized and a sweep is performed to find the indices where the boundaries are most similarly aligned between the original data and the synthetic data, as measured using a simple Euclidean norm. The values between the boundaries within the synthetic data are then taken to patch the missing gap of size  $c$  within the original.

26910	4765	46	26910	4765	46
25156	62631	116	25156	62631	116
80	5732	184	<b>80</b>	<b>7178</b>	<b>143</b>
443	63653	52	<b>443</b>	<b>45850</b>	<b>222</b>
80	52282	419	<b>80</b>	<b>31877</b>	<b>124</b>
25443	60829	46	<b>443</b>	<b>34456</b>	<b>72</b>
443	45368	92	<b>25514</b>	<b>4390</b>	<b>133</b>
25186	59341	95	<b>443</b>	<b>59685</b>	<b>124</b>
27473	50041	92	<b>26136</b>	<b>58498</b>	<b>117</b>
80	64955	116	<b>26772</b>	<b>60792</b>	<b>279</b>
25161	4145	819	25161	4145	819

Fig. 4. Filling in missing data: (Left) Shaded entries represents missing data. (Right) Synthesized entries.

Figure (4) provides an example of this regression method. Here, we have pre-trained a model for a particular host. We took a previously unseen segment of traffic from this host and removed a portion of it, this is represented by the shaded values on the left side of the figure. On the right side of the figure, denoted in bold font, we show the sequence of entries generated using our regression technique. Visually, we can confirm the similarity between the real and synthesized data.

#### E. Extensions into Packet-Trace Synthesis

Our algorithm can generate statistical Netflow-like data, for the form seen in Figure (4). This data encapsulates the statistical behavioral profile for any given host on a network. For many research purposes this is all that is needed; in many network-traffic measurement problems, for example, one is interested in finding bottlenecks, heavy hitters, rates, and patterns among traffic. Given that the synthesized traffic faithfully recreates the same statistical distribution as the original traffic, this is sufficient for many network research tasks that only require statistics. However, certain, research requires packet data. Recently, as an extension of this work, we have produced a systematic method to generate full TCP sessions driven by these statistical profiles. The synthesized packet streams contain fully session-

izable packets. This implies many properties: proper TCP handshakes, proper header values such as window-sizes, as well as sequence numbers which are consistent to the path MTUs *etc.* Our method is implemented as a custom pseudo-TCP software stack. The size, duration, length, direction, timing and other statistical characters of newly generated packet sessions are specified by the synthesized statistical data, and the software stack uses this information to craft packets which are consistent with this profile and in a manner that conforms to RFC specifications for TCP. The full details of this procedure will be explained in a forthcoming paper.

#### F. Memory cost and runtimes

The memory footprint of our models is small when implemented using efficient underlying data structures. The largest variable within the model is the transition-probabilities table  $T$ . This size is upper bounded by the size of the list of potential port features. In our implementation, we had 1025 well known ports, 20 registered-ports bins, 1 ephemeral-port bin, which yields a  $1046 \times 1046$ -sized transition table. If double precision floating point storage is used for all values then we have a roughly 8Mb maximum storage requirement per model. However, the transition-probabilities table is typically very sparse, because most hosts exercises only a small subset of potential network protocols, unless a host is experiencing some sort of port-scanning activity. In practice, for a typical host which uses a dozen network services, the storage requirement using sparse-matrix implementation is roughly 12Kb per host. This means that on modern computing platforms it is easily possible to compute profiles for networks containing tens of thousands of hosts and keep all parameters in memory simultaneously.

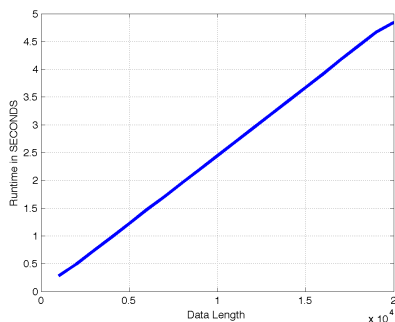


Fig. 5. SYNTHESIZE-NETFLOW maintain  $O(n)$  growth. Large traffic set with 20,000 session entries took less than five seconds to synthesize.

Figure (5) shows the runtime of our synthesis algorithm. As expected, our algorithm grows linearly in  $O(n)$ . The experiment shown in Fig. (5) shows runtimes for synthesizing the behavior of a host which had used 16 network services. The figure shows that synthesizing a relatively large 20,000 element Netflow record required less than five seconds. Our computation was performed on a 64-bit 2.66Ghz processor with six cores. We believe our relatively un-optimized research-oriented code can be improved in production settings and significantly higher performance can be achieved.

## V. CONCLUSION

In this paper we have described a model for representing network behavior at the host level, for purposes of network traffic simulation. We showed how to extend these models and use them to synthesize realistic looking traffic that mimic the original hosts. We demonstrated how to modify these models to manipulate the resulting synthetic data, by adjusting the parameters, and interpolating between different models to simulate adjustable new behavior. We further showed how to use our models to perform data regression, to fill in “missing” gaps of traffic. Methods to quantify the fidelity of these algorithms are presented, and backed up by experiments which demonstrate favorable performance of our algorithm. To the best of our knowledge, we have presented the first system that can be deployed to a foreign network environment, automatically learn the behavior of hosts on the network, then synthesize new traffic which simulates this behavior on a separate network. at both the Netflow and packet levels.

## ACKNOWLEDGMENTS

We thank Prof. Angelos Stavrou and Brian Schulte at GMU, and Kyle Creyts at Merit Network for providing us with valuable data (UMich. dataset.) This research was sponsored by Department of Homeland Security, SPAWAR Contract No. N66001-09-C-0080, Privacy Preserving Sharing of Network Trace Data (PPSNTD) Program and a DURIP Instrumentation grant from AFOSR (FA 99500910389).

## REFERENCES

- [1] Y. Song, S. J. Stolfo, and T. Jebra, “Markov models for network-behavior modeling and anonymization,” Columbia University, Technical Report cues-029-11, June 2011.
- [2] J. Mirkovic, “Privacy-safe network trace sharing via secure queries,” in *Proceedings of the 1st ACM workshop on Network data anonymization*, 2008.
- [3] R. Bejtlich, J. Cummings, and S. Parsell, “Openpacket.org: a centralized repository of network traffic traces for researchers,” 2011. [Online]. Available: <https://www.openpacket.org>
- [4] PREDICT, “PREDICT: The Protected Repository for the Defense of Infrastructure Against Cyber Threats,” <http://www.predict.org>, 2011. [Online]. Available: <http://www.predict.org>
- [5] DARPA, “National Cyber Range (NCR) program,” DARPA-BAA-08-43, 2008.
- [6] M. Allman and V. Paxson, “Issues and etiquette concerning use of shared measurement data,” in *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*, ser. IMC ’07. New York, NY, USA: ACM, 2007, pp. 135–140. [Online]. Available: <http://doi.acm.org/10.1145/1298306.1298327>
- [7] BreakingPoint Systems, “Breakingpoint network traffic generation systems,” 2011. [Online]. Available: <http://www.breakingpointsystems.com/>
- [8] Spirent, “Spirent hypermetrics 40/100g module,” 2011. [Online]. Available: <http://www.spirent.com/>
- [9] G. Lyon, “Nmap - free security scanner for network exploitation & security audits,” 2011. [Online]. Available: <http://nmap.org>
- [10] Tenable Security, “Nessus: The network vulnerability scanner,” 2011. [Online]. Available: <http://www.tenable.com>
- [11] Metasploit, “Metasploit Project,” 2010, <http://www.metasploit.com>.
- [12] S. Sanfilippo, “hping: active network security tool,” 2006. [Online]. Available: <http://www.hping.org/>
- [13] J. Nathan, “Nemesis: Packet injection tool suite,” 2011. [Online]. Available: <http://nemesis.sourceforge.net/>
- [14] W. Foundation, “wireshark: network protocol analyzer,” 2011. [Online]. Available: <http://www.wireshark.org/>
- [15] E. Paris, “Network expect,” 2010. [Online]. Available: <http://netexpect.org/>
- [16] T. Jebra, R. Kondor, and A. Howard, “Probability product kernels,” *Journal of Machine Learning Research*, vol. 5, pp. 819–844, 2004.