# Manifold Learning Visualization of Network Traffic Data \*

Neal Patwari npatwari@eecs.umich.edu hero@eecs.umich.edu

Alfred O. Hero III

Adam Pacholski apachols@umich.edu

University of Michigan Dept. of Electrical Engineering and Computer Science 1301 Beal Avenue, Ann Arbor, MI, USA

### **ABSTRACT**

When traffic anomalies or intrusion attempts occur on the network, we expect that the distribution of network traffic will change. Monitoring the network for changes over time, across space (at various routers in the network), over source and destination ports, IP addresses, or AS numbers, is an important part of anomaly detection. We present a manifold learning (ML)-based tool for the visualization of large sets of data which emphasizes the unusually small or large correlations that exist within the data set. We apply the tool to display anomalous traffic recorded by NetFlow on the Abilene backbone network. Furthermore, we present an online Java-based GUI which allows interactive demonstration of the use of the visualization method.

# **Categories and Subject Descriptors**

C.2 [Computer-Communication Networks]: Network Monitoring. General Terms: Algorithms, Measurement. **Keywords**: Internet traffic anomaly detection & forensics. data mining.

### INTRODUCTION

Statistical intrusion and anomaly detection methods allow networks to be monitored for attacks for which attack signatures have not yet been developed. However, the huge quantity and high-dimensionality of internet traffic data are significant challenges which research must overcome in order to achieve high reliability and low false-alarm rates. Recently, subspace-based analysis of traffic data by Lakhina, Crovella, and Diot [1, 2] has shown that high-dimensional Abilene traffic measurements can be well-represented within a very low-dimensional subspace. The 'curse-of-dimensionality' can be avoided when high-dimensional data can be represented well in a low-dimensional subspace.

In this paper, we use a manifold learning method to take very high-dimensional NetFlow traffic measurements from the Abilene backbone network and reduce their dimensionality to two dimensions. The resulting 2-D 'map' of the measurements provides a means for visualization of the relationships which exist in a set of traffic data. These relationships may be spatial, eg., between measurements taken across a backbone network or between IP addresses, autonomous systems (AS), or origin-destination (OD)-flows; temporal, eg., measurements taken at different times; or between different applications, as indicated by port numbers.

Such visualization is complementary to detection methods which rely on dimensionality reduction. Subspace-based detection [2] has been successfully used infer the presence of spatial traffic distribution anomalies in network-wide traffic measurements. This inference is done by quantifying the amount of traffic which cannot be represented within a low-dimensional subspace. The method we present in this paper allows the visualization of the traffic which can be represented within a low-dimensional subspace. Furthermore, this work uses a non-linear dimensionality reduction method rather than a linear subspace method.

# 1.1 Related Work

Much research in dimensionality reduction and internet traffic visualization are foundational to this paper.

### 1.1.1 Dimensionality Reduction

Classical multi-dimensional scaling (MDS) and principal components analysis (PCA) are perhaps the most commonly recognized dimensionality reduction methods, but they assume that the high-dimensional data points lie on a linear subspace (for example, on a 2-D plane) of the highdimensional space. Manifold learning algorithms [3, 4, 5] are more general - data may lie on a curved subspace. For example, a 2-D manifold could be a portion of a sphere, or a 'swiss roll'. When high dimensional data points are random but highly correlated with their close neighbors, data points don't tend to fall into linear subspaces. As a result, manifold learning methods can be more effective than linear methods like MDS.

Sketching is a dimensionality reduction method which projects data onto random linear lower dimensional subspaces, in a way that preserves, inter-data distances, approximately, with high probability [6]. Sketching has been used to dramatically reduce the number of dimensions necessary to store multi-dimensional histograms [7]. Since the visualization presented in this paper calculates distances between multi-dimensional histograms, sketching could be used to reduce the storage and communication complexity of a distributed implementation of the proposed method. Sketching is not tested in this preliminary, centralized implementation.

<sup>\*</sup>This research was supported in part by National Science Foundation ITR Grant No. CCR-0325571.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGCOMM'05 Workshops, August 22–26, 2005, Philadelphia, PA, USA. Copyright 2005 ACM 1-59593-026-4/05/0008 ...\$5.00.

Clustering of packet traffic data [8] is a data reduction method in which unusually large concentrations of traffic in feature space are automatically identified and reported. Dimension reduction might help improve clustering performance by reducing the 'curse of dimensionality'.

In this paper, we use the manifold learning method, distributed weighted MDS (dwMDS) method [5]. Its key features are: (1) An algorithm which allows for a distributed implementation across the network with minimal communication requirements, (2) Consideration of prior information, which allows use of a 'typical' map to be used as a baseline, thus allowing for easy comparison of data maps over time, (3) A weighted cost function that allows neighbor relationships that are believed to be more accurate to be weighted more heavily, and (4) A majorization-based optimization such that each iteration is guaranteed to improve the value of the cost function.

### 1.1.2 Visualization

We believe that data visualization will complement statistical detection methods, and help provide information to help a human moderator make a decision regarding whether or not an anomaly has occurred, and if so, to determine its temporal and spatial characteristics.

Visualization methods have found much use in network monitoring. For example, visualization of flows by application over time is commonly done using FlowScan [9]. Monitoring the number of flows over time is an excellent tool to identify DoS attacks. However, as attacks (such as the Slammer worm) exploit smaller user populations, even obscure services' traffic must be monitored. Dimensionality reduction is a means to monitor, separately, hundreds or thousands of traffic statistics but to minimize the complexity of the information display. Notably, the 'Spinning Cube of Potential Doom' [10] graphically shows port scans and worm activity in 3D by plotting recent packets' IP source and destination address and destination port. Such visualization is also presented by NVisionIP [11], a multi-faceted visualization tool which can filter traffic and 'zoom in' to various network scales.

Graph visualization techniques have been used to show the physical connections that exist in a network. Visualizations of the global internet, such as CAIDA's Skitter plot [12], are important statements about the interconnectivity of the global network. Other tools developed at CAIDA, such as Otter and Walrus, provide 2-D and 3-D visualization of network graphs. The visualization method presented in this paper provides information not just about the connections that exist, but also the traffic correlations that exist. Connection distances match correlation - when correlation between two nodes is high, they are plotted close together.

#### 1.2 Framework

We frame the dimensionality reduction problem as a sensor data localization problem. In this framework, 'sensors' are the hardware or software which record data, for example, on each router in a backbone network. The traffic data which they record can be of arbitrarily high dimension. For example, rather than counting the grand total number of flows (just one dimension), they would count the total number of flows from each source IP address (up to 2<sup>32</sup> dimensions). The key to understanding a particular sensor data map is to answer three questions:

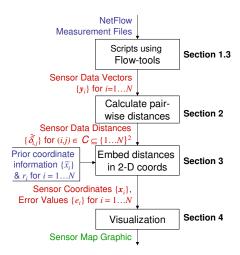


Figure 1: Flow chart of data visualization from Net-Flow data input to data map output.

- 1. Where are the sensors?: Sensors can be 'located' at physical computers, i.e., at backbone routers, or at IP addresses; or they can be 'located' at less physical concepts such as source or destination ports, or particular time periods. A sensor attached to a particular source port monitors only traffic which matches its source port, and a sensor attached to a time monitors only traffic which arrives during that time period.
- 2. What traffic statistic is recorded?: Sensors might measure flows, packets, or octets, or some combination.
- 3. What are the dimensions?: Sensors can divide traffic by source or destination IP address, port, or AS; time period; link or router; or some combination. Traffic statistics are then recorded for each dimension separately.

For example, in Section 4.1, sensors are located at backbone routers, recording the number of flows from each source IP address. As another example, in Section 4.2, sensors are located on destination port numbers, recording the number of flows from each source IP address. As a final example, in Section 4.3, sensors are located on backbone routers, recording the total number of packets received in each of the past T 5-minute time intervals.

We note that this framework can be used to describe the measurements in [1], in which sensors were located at all 10-minute time intervals over the course of a week, and sensors measured total octets on each link across the (Sprint-Europe or Abilene) backbone network.

We denote the data measured at sensor i as  $\mathbf{y}_i$ , where  $i \in \{1, \dots, N\}$ , where N is the total number of sensors. The high-dimensional vector  $\mathbf{y}_i$  is defined by,

$$\mathbf{y}_i = [y_i(l_1), y_i(l_2), \dots, y_i(l_{|\mathcal{L}|})],$$
 (1)

where  $\mathcal{L}$  is the set of possible dimensions (see #3 above) with  $|\mathcal{L}|$  elements, and  $l_k \in \mathcal{L} \forall k$ . In many cases,  $y_i(l_k) = 0$  for most of its elements  $l_k$ , thus we store  $\mathbf{y}_i$  as a sparse vector. For example, the set of possible IP addresses is much larger than the set of IP addresses observed in a particular traffic stream.

We use the flow-tools package created by Mark Fullmer [13] to process NetFlow files in order to generate the data for the vectors  $\{\mathbf{y}_i\}$ .

ATLA	CHIN	DNVR
130.14.24.0, 1545	129.25.0.0, 13913	129.25.0.0 14331
131.247.224.0, 1487	141.89.48.0, 8738	207.46.104.0 12142
128.61.64.0, 1197	207.46.104.0, 3708	207.46.248.0 7198
198.32.152.0, 1147	204.179.120.0, 3520	207.68.176.0 4968
164.111.192.0, 1139	203.250.224.0, 3441	64.4.16.0 4156
131.247.232.0, 1098	207.46.248.0, 3300	207.68.168.0 3707
	•	
:	;	:

Table 1: Example data: Top few lines of  $\{y_i\}$  for 3 Abilene routers, flows by source IP (last 11 bits zeroed) for 5 minutes ending 20 Jan 2005 01:00 UTD.

### 2. MEASUREMENT DISTANCES

Next, distances between  $\{y_i\}_i$  are calculated. In this paper, we normalize each data vector such that its sum is one:

$$\tilde{\mathbf{y}}_i = \frac{\mathbf{y}_i}{\|\mathbf{y}_i\|_1} \tag{2}$$

where  $\|\mathbf{y}_i\|_1$  is the  $L_1$  norm, *i.e.*, the total traffic measured at sensor i. The value  $\tilde{\mathbf{y}}_i(l)$  thus is the fraction of traffic measured in dimension l. We use this normalization to prevent two sensors from being measured to be far apart solely because they experience different absolute traffic levels. For example, on Abilene, the Atlanta router generally records about 1/3 of the Indianapolis router's traffic. Without normalization, these two routers would always be far apart. With normalization, the distance between routers depends on the distribution, rather than the total quantity of traffic.

Let  $\delta_{i,j}$  be the distance between the measurement vectors from sensors i and j. In this paper, we use the Euclidean distance between  $\tilde{\mathbf{y}}_i$  and  $\tilde{\mathbf{y}}_j$ ,

$$\delta_{i,j} = \left[ \sum_{l \in |\mathcal{L}|} (\tilde{\mathbf{y}}_i(l) - \tilde{\mathbf{y}}_j(l))^2 \right]^{\frac{1}{2}}$$
 (3)

Other distance metrics, such as histogram intersection [14], may be desirable in some cases, and this is a topic for future research. Euclidean distance is a simple, symmetric and finite distance metric used here as a proof of concept.

### 2.1 Neighbor Selection

Using the set  $\{\delta_{i,j}\}$ , the sensor neighbor relation is determined. Intuitively, pairs of sensors which are 'close' to each other will consider each other to be neighbors. We use the K-nearest neighbors method to determine the neighbor set, which we denote  $\mathcal{C} \subset \{1,\ldots,N\}^2$ , In K-nearest-neighbors, a pair i and j are neighbors if either j is one of the K nearest neighbors of i. Specifically, i and j are neighbors if  $\tilde{\delta}_{i,j}$  is one of the K smallest of either set  $\{\delta_{i,k}\}_{k\neq j}$  or set  $\{\delta_{k,j}\}_{k\neq j}$ .

### 2.2 Constant Map Size

We note that the *shape* of the map is important, but the *size* of the map is not very informative. In fact, if the size of the map changes, it will make it more difficult to view

it over time. To prevent this, we normalize distances by a constant,

$$\tilde{\delta}_{i,j} = \delta_{i,j} \frac{D_{\mathbf{x}}}{D_{\mathbf{y}}} \tag{4}$$

where  $D_{\mathbf{y}} = \sum_{(i,j) \in \mathcal{C}} \delta_{i,j}$  is the sum of all of the distances between neighbors, and  $D_{\mathbf{x}}$  is a given, desired distance sum, for some  $D_{\mathbf{x}} \in \mathbb{R}$ .

### 3. COORDINATE EMBEDDING

Using the distances  $\{\tilde{\delta}_{i,j}\}$ , for neighbor pairs  $(i,j) \in \mathcal{C}$ , we next calculate 2-D coordinate embedding. In the dwMDS method, this means that we find the 2-D coordinates  $\{\mathbf{x}_i\}_{i=1...N}$  which best represent the pair-wise distances, and the prior coordinate information, in a weighted-least-squares sense. Specifically, we find the 2-D vectors  $\{\mathbf{x}_i\}_{i=1...N}$  which minimize the following cost S,

$$S = 2\sum_{i=2}^{N} \sum_{j=1}^{i-1} w_{i,j} \left[ \tilde{\delta}_{ij} - \|\mathbf{x}_i - \mathbf{x}_j\| \right]^2 + \sum_{i=1}^{N} r_i \|\mathbf{x}_i - \overline{\mathbf{x}}_i\|^2$$
 (5)

where  $w_{i,j}$  is the weight assigned to the link (i,j),  $r_i$  is the weight assigned to the prior coordinate information for sensor i, and  $\overline{\mathbf{x}}_i$  is the prior coordinate given for sensor i. First, we discuss the assignment of  $w_{i,j}$ , and then, we discuss the assignment of prior coordinate information.

# 3.1 Assignment of Weights

Weights  $w_{i,j}$  allow pairs which are perceived to have a less accurate measured distance to be down-weighted. Inspired by the weighting frequently used in locally weighted regression methods (LOESS) [15], we use a weighting  $w_{i,j}$  that is a decreasing function of measured distance  $\tilde{\delta}_{i,j}$ :

$$w_{i,j} = \begin{cases} \exp\left\{-\tilde{\delta}_{i,j}^2/h_{i,j}^2\right\}, & \text{if } (i,j) \in \mathcal{C}, \\ 0, & \text{otherwise,} \end{cases}$$
 (6)

where  $h_{i,j} = \max_k \{\tilde{\delta}_{i,k}, \tilde{\delta}_{k,j}\}$ . This choice of  $w_{i,j}$  is symmetric and equalizes the (nonzero) weight distribution in all sensors. In the examples shown in [5], LOESS weighting results in higher precision than an alternative equal-weighting scheme, *i.e.*,  $w_{i,j} = 1$  for all neighbor pairs (i, j).

# 3.2 Prior Coordinate Information

Prior coordinate information is an option in the dwMDS method. If there is no prior information for sensor i, we set  $r_i = 0$ . If all sensors have  $r_i = 0$ , the calculated output map can arbitrarily be translated, rotated, and flipped without affecting its cost S. The purpose of the map is to show the relationships between sensors' data, and as such, translation, rotation, and flipping do not change the meaning of the map. However, if the user will view many maps in sequence over time, it would be confusing if each subsequent map was nearly identical but with arbitrary rotation. Prior coordinate information  $(r_i > 0)$  is a means to provide a stable orientation for a sensor map. For example, for the sensor maps presented in Section 4.1, we use the mean coordinates plotted in Fig. 2(a) as the prior coordinates  $\{\overline{\mathbf{x}}_i\}$ .

For  $0 < r_i < \infty$ , the choice of  $r_i$  determines how sensitive the sensor's location is to changes in sensor data. A high  $r_i$ , that is,  $r_i \gg \sum_j w_{i,j}$  will mean that node i will be estimated to be very close to  $\overline{\mathbf{x}}_i$  except in extreme circumstances. In Section 4.1, we use a low  $r_i$ ,  $i.e.r_i = 10^{-3}$ , which

is  $\ll \sum_{j} w_{i,j}$ , so that the prior coordinates are used *only* to remove translational and rotational degrees of freedom. The shape of the sensor map is not affected when such a low  $r_i$  is used. The web applet presented in Section 4.3 allows  $r_i$  to be set by the user; so testing at different  $r_i$  is possible.

# 3.3 Algorithm Overview

The dwMDS algorithm is a distributable, iterative algorithm to minimize S and find the 2-D coordinate embedding. In this paper, we only present a broad overview of the dwMDS algorithm, and refer to [5] for more details.

The dwMDS is distributable because the global cost S in (5) is separable by sensor, i.e.,  $S = \sum_{i=1}^N S_i$ , where  $S_i$  is sensor i's contribution to the cost. Sensor i can minimize  $S_i$  using only distances  $\tilde{\delta}_{i,j}$  between itself and its neighbors j, and its neighbors' most recent 2-D coordinate estimate. This minimization is done using a simple majorization method, which guarantees non-increasing cost in each round of the optimization. After sensor i minimizes  $S_i$ , it updates its neighbors with its new coordinate estimate. Then, it passes to the sensor i+1, which performs its own minimization. The algorithm visits each sensor in turn, and may perform several rounds until convergence.

We use a centralized implementation of the dwMDS algorithm, in which the data is collected and optimized at a single processor [16]. This is used to demonstrate the visualization method's capabilities, but the distributed implementation discussed here and in [5] is future work. We limit ourselves to 2-D sensor maps in this paper, but the dwMDS method could produce higher dimensional coordinates, if desired.

#### 3.4 Error Metric

We also define an error value  $e_i^2$  for sensors  $i \in \{1, ... N\}$ :

$$e_i^2 = \sum_{(i,j)\in\mathcal{C}} w_{i,j} \quad \tilde{\delta}_{ij} - \|\mathbf{x}_i - \mathbf{x}_j\|^2$$
 (7)

The value of  $e_i$  quantifies the information lost in the 2-D representation of sensor i's coordinates. It represents the quantity of measurement distances  $\{\tilde{\delta}_{i,j}\}_j$  which are not represented by the 2-D coordinates of itself,  $\mathbf{x}_i$ , and its neighbors,  $\{\mathbf{x}_j\}$ . The value  $e_i^2$  is analogous to the residual value in PCA and can be used to help decide whether or not the sensor's data is anomalous [1]. In this paper, we 'shade' each sensor as a function of  $e_i$ : low  $e_i$  is white, and high  $e_i$  is black. Examples are presented in the following section.

#### 4. CASE STUDIES

In this section we show several examples of the visualization method on Abilene traffic data. We use NetFlow data recorded during January 2005, available from the Abilene Observatory [17]. Note that, for privacy reasons, only the most significant 21 bits of IP addresses are available - the last 11 bits are zeroed out. Furthermore, NetFlow data is sampled at 1/100, so for each packet reported here, there were 99 more unrecorded.

### 4.1 Router Maps

When sensors are routers in a backbone network, the 'router map' shows the spatial characteristics and correlations of the routers' traffic, rather than just the connectivity of the routers. In the following examples, sensors are

routers, and we measure flows in a 5-minute period, separated by source IP address. (Example vectors are shown in Table 1.) The size of each sparse data vector is limited to 1000 – flows not from the top 1000 source IP addresses (/21) for a particular router are ignored. For this reason, source IP addresses are typically lost if they have less than 10-50 flows. We set the number of neighbors to K=5.

# 4.1.1 Typical Activity

First, we characterize 'typical' router map behavior by calculating router maps for the four-week period 02-Jan to 29-Jan. Since a router map is calculated for each 5 minute period, we calculate a total of 4\*7\*24\*60/5 = 8064 maps. The sample mean and covariance of  $\mathbf{x}_i$  over the 8064 maps is shown in Fig. 2(a). Although there are certainly attacks active during this 4-week period, averaging maps over a long period of time provides intuition about the typical behavior for the router map.

The typical router map in Fig. 2(a) both makes sense geographically and describes typical traffic patterns seen on Abilene. Much of Abilene traffic is East-West or West-East, and Northern routers (especially DNVR, KSCY, IPLS, and CHIN) bear much of this traffic. These routers have very correlated sensor data because a significant proportion of Abilene OD-flows pass through all of them.

#### 4.1.2 06-Jan-2005

First, we show that the router map changes very dramatically with very large traffic changes. On Thursday, 6-Jan-2005, during the 5-minute period ending at 17:55 UTD, NetFlow recorded on the CHIN router a total of  $9\times 10^4$  single-packet flows (of 40-byte packets) from two source IP addresses in Taiwan to a small range of destination IP addresses in Hungary. This volume is about 25% of the typical flow volume on CHIN. The anomalous traffic was observed on CHIN and no other router, thus distances between sensor data recorded at CHIN and other routers are unusually high, and the 2-D coordinates for CHIN must be distant from all the others. Also, because of the normalization done to calculate  $\tilde{\delta}$  from  $\delta$ , the rest of the map distances have shrunk to compensate. This is shown in Fig. 2(b).

# 4.1.3 20-Jan-2005

On Thursday, 20-Jan at 01:00 UTD, there are a large number (14,000) of 29-byte UDP packets from a 129.25.0.0 (Drexel U.) source IP address sent to a 131.252.120.0 (Portland State U.) destination (see Table 1). The packets are from source port 3095 or 3096 to a wide range of destination ports > 1024. These packets travel through the WASH, NYCM, CHIN, IPLS, KSCY, DNVR, and STTL (Northern) routers, but not through any other routers. Distances between the listed Northern routers and the other Southern routers are unusually high. In the router map shown in Fig. 2(c), there is a clear split in the map between the two sets of routers.

# 4.2 Port Maps

Next, we present 'port maps', which can be used to visually detect ports which exhibit anomalous traffic. For a port map, each sensor is attached to a destination port (only sees flows which match the port) and record the number of flows per source IP address. We would expect that attacks or scanning activity would exhibit very different distributions

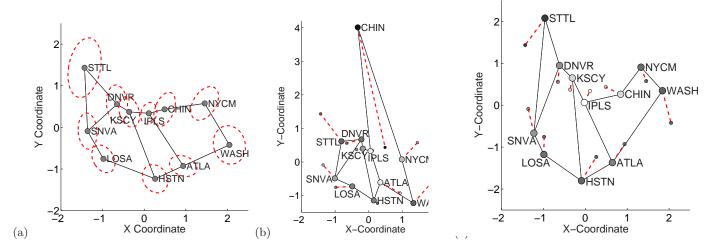


Figure 2: (a) Mean ( $\bullet$ ) and 1- $\sigma$  uncertainty ellipse (---) of router maps from 2-Jan to 29-Jan. Maps during (b) port scan on 6-Jan 17:55 and (c) attack on 20-Jan 01:00, show router coordinates ( $\bullet$ ) connected (---) to the mean ( $\cdot$ ) from (a), and shaded by error value  $e_i$ . All figures show Abilene backbone links (—).

of source IP addresses as compared to ports that aren't the subject of attacks or scans. The port map can be calculated using only the data at one particular router (Fig. 3 uses data from IPLS).

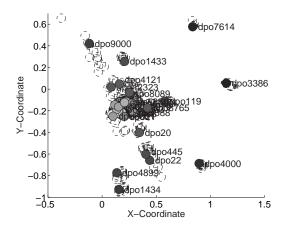


Figure 3: Destination port map for 01-Jan-2005 at 3:35 UTD (•dpo#) along with past 1 hour map history (dotted line circles). Sensors are attached to the top 30 destination ports (by total flows) and measure number of flows per source IP address.

Here, we attach sensors to the top 30 ports (by number of flows). We measure flows per source IP address within the current 5 minute time period. Then, the port map is calculated with K=15 and no prior  $(r_i=0 \text{ for all } i)$ . For comparison, we include port maps calculated over the past hour, which have been rotated and translated to line up with the most recent map, for ease of comparison. Fig. 3 shows an example for 1-Jan-2005. While most ports fall very close to each other near the origin, a few ports have very different source IP address distributions, thus are placed far from the center and have a higher  $e_i$  (as indicated by the darker shade of their circle). Notably, destination ports with known

vulnerabilities, 9000, 1433, 7614, 3386, 4000, 445, 22, 4899, and 1434 are mapped far from the center.

Many of these ports are mapped to be 'abnormal' over a wide range of time (days), but destination port 7614 appeared on the port map rarely. Its appearance at 03:55 corresponds to a scan originating from source IP 211.82.216/21 to destination port 7614 using single, 48-byte packet flows. The port map emphasizes such a small (480 flows out of  $1.2 \times 10^5$  total) scan because 100% of port 7614 traffic has an otherwise unrecorded source IP address.

# 4.3 Map Web Applet

We have available a Java-based applet which displays temporal and spatial traffic on the Abilene backbone [16], as shown in Fig. 4. The applet has a graph of the traffic for each router, and can plot traffic levels from 21 different sets of ports grouped by application.

In addition to displaying the traffic by port and time, the applet calculates a router map for each time. In this router map, the sensors are routers, and they measure total packets, by port and 5-minute time interval. The user can select which groups of ports, and the number of time intervals to use as dimensions. Rather than normalizing to calculate  $\tilde{\mathbf{y}}_i(l)$  as described by (2), we use  $\tilde{\mathbf{y}}_i(l) = \mathbf{y}_i(l) - \overline{\mathbf{y}}_i(l)$ , where  $\overline{\mathbf{y}}_i(l)$  is the median of the past  $T_m$  time samples, where  $T_m$  is also user-adjustable. Using a filtered traffic stream emphasizes the changes that occur over time. When traffic changes over time, the router map shows where (which routers) the change is most dramatic. Details are available online [16].

# 5. FUTURE WORK AND CONCLUSION

We have introduced a visualization tool which can aid in the discovery of traffic anomalies in high-dimensional Net-Flow data. We urge the reader to visit the online repository to view other sensor maps and to utilize the web-based visualization applet [16]. Clearly, manifold learning methods can be used to provide information about relationships that exist in sets of network traffic data.

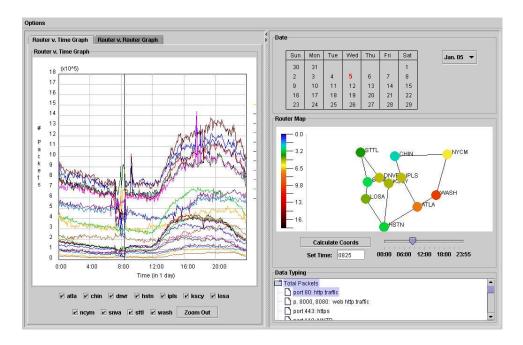


Figure 4: Total traffic and port 80 traffic on 05-Jan-2005, displayed using the visualization applet [16]. The router map is calculated for 08:25 UTD, during scheduled maintenance of the CHIN-IPLS link, during which traffic drops at CHIN and IPLS and increases dramatically at the HSTN and ATLA routers.

Future work will attempt to automate the detection and classification process, similar to the application of subspace-based detection methods in [1, 2]. Using sketching would further adapt the method for distributed implementation. Other distance metrics may emphasize similarities in traffic distributions; and other manifold-learning methods such as Isomap [3] should be tested. It is hoped that router maps, port maps, and other types of sensor maps may together serve as a step-by-step investigation aid, iteratively helping to locate a traffic anomaly in a very high-dimensional spatial and temporal data space.

# Acknowledgements

The authors thank Panna Felsen, who initiated the visualization applet as part of the NASA SHARP program.

# 6. REFERENCES

- A. Lakhina, M. Crovella, and C. Diot, "Diagnosing network-wide traffic anomalies," in ACM SIGCOMM '04, Aug. 2004.
- [2] ——, "Characterization of network-wide anomalies in traffic flows," in *ACM/SIGCOMM Internet Measurement Conference*, Oct. 2004.
- [3] J. B. Tenenbaum, V. de Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, pp. 2319–2323, Dec 2000.
- [4] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by local linear embedding," *Science*, vol. 290, pp. 2323–2326, Dec 2000.
- [5] J. A. Costa, N. Patwari, and A. O. Hero III, "Distributed multidimensional scaling with adaptive weighting for node localization in sensor networks," *IEEE/ACM Trans. Sensor Networks*, submitted May 2004, (revised Jan. and May 2005). [Online]. Available: http://www.eecs.umich.edu/~hero/comm.html
- [6] A. C. Gilbert, Y. Kotidis, S. Muthukrishnan, and M. J. Strauss, "QuickSAND: Quick summary and

- analysis of network data," DIMACS, Tech. Rep. 2001-43, Nov. 2001. [Online]. Available: http://www.math.lsa.umich.edu/~annacg/ps.files/quickdimacstr.ps
- [7] N. Thaper, S. Guha, P. Indyk, and N. Koudas, "Dynamic multidimensional histograms," in ACM SIGMOD 2002, June 2002, pp. 428–439.
- [8] C. Estan, S. Savage, and G. Varghese, "Automatically inferring patterns of resource consumption in network traffic," in ACM SIGCOMM'03, Aug. 2003, pp. 137–148.
- [9] D. Plonka, "FlowScan: A network traffic flow reporting and visualization tool," in *Large Installation* System Administration (LISA) Conference 2000, Dec. 2000, pp. 305–317.
- [10] S. Lau, "The spinning cube of potential doom," Communications of the ACM, vol. 47, no. 6, pp. 25–26, June 2004.
- [11] K. Lakkaraju, W. Yurcik, and A. J. Lee, "NVisionIP: Netflow visualizations of system state for security situational awareness," in ACM VizSEC/DMSEC04, Oct. 2004, pp. 65–72.
- [12] Cooperative Association for Internet Data Analysis (CAIDA). http://www.caida.org/.
- [13] M. Fullmer and S. Romig, "The OSU Flow-tools package and Cisco NetFlow logs," in *Large* Installation System Administration (LISA) Conference 2000, Dec. 2000, pp. 291–303.
- [14] M. J. Swain and D. H. Ballard, "Color indexing," *Int. J. Comp. Vision*, vol. 7, no. 1, pp. 11–32, 1991.
- [15] W. Cleveland, "Robust locally weighted regression and smoothing scatterplots," *J. American Statistical Assoc.*, vol. 74, no. 368, pp. 829–836, 1979.
- [16] Map-tools online supplement. http://www.engin.umich.edu/~npatwari/mnd05.
- [17] Internet2: Abilene Observatory. http://abilene.internet2.edu/observatory/.