

# Putting Security in Context: Visual Correlation of Network Activity with Real-World Information

W.A. Pike, C. Scherrer, and S. Zabriskie

**Abstract** To effectively identify and respond to cyber threats, computer security analysts must understand the scale, motivation, methods, source, and target of an attack. Central to developing this situational awareness is the analyst's world knowledge that puts these attributes in context. What known exploits or new vulnerabilities might an anomalous traffic pattern suggest? What organizational, social, or geopolitical events help forecast or explain attacks and anomalies? Few visualization tools support creating, maintaining, and applying this knowledge of the threat landscape. Through a series of formative workshops with practicing security analysts, we have developed a visualization approach inspired by the human process of contextualization; this system, called NUANCE, creates evolving behavioral models of network actors at organizational and regional levels, continuously monitors external textual information sources for themes that indicate security threats, and automatically determines if behavior indicative of those threats is present on a network.

## 1 Introduction

Visualization can have a central role in helping computer security analysts understand the changing state of their systems. But to take action on the basis of what they see in visual displays, analysts must be able to do more than just perceive anomalous changes; they must be able to *explain* them. Analysts also need the capacity to be proactive. Ideally, they should be able to identify potential security threats before an attack is incipient. Explaining and predicting security events often require more than just network- or host-centric information. External information such as exploit and vulnerability reports from other organizations, security advisories and even news stories are critical in helping analysts put anomalies in context.

---

W.A. Pike, C. Scherrer, and S. Zabriskie

Pacific Northwest National Laboratory, MSIN K7-28, P.O. Box 999, Richland, WA 99352, USA,  
e-mail: william.pike@pnl.gov, chad.scherrer@pnl.gov, sean.zabriskie@pnl.gov

We introduce a security visualization and analysis approach called NUANCE, which results from extensive engagement with practicing analysts on computational methods that help them put their network observations into real-world context. NUANCE abstracts packet-level data to higher-level behavioral models for each IP address and group of addresses (such as an organizational subnet) observed on a network. Through a new heterogeneous data integration approach, NUANCE then visually fuses these behavioral models with contextual information on current threats and exploits from open sources. Our approach automatically collects this contextual information and determines to which actors or groups observed on a network (e.g., particular IP addresses, organizational units, or geographic areas) it is relevant. Creating, maintaining, and applying contextual knowledge of the threat landscape is the analyst's stock in trade, and NUANCE automates and scales up this practice. By seeking out context, NUANCE results in a lower information processing burden on the analyst and can direct the analyst's attention to activities most worthy of scrutiny.

We also describe the process of working with practicing analysts to develop and refine our visualization techniques. Through a user-centered process that included a series of formative design and evaluation workshops, we produced a set of visual interfaces that help analysts identify and explain off-normal activities. These workshops motivated our research by framing the questions that analysts ask in the course of investigating anomalies and the work practices that visualization systems need to support.

## 2 Related Work

Current network monitoring tools can produce potentially overwhelming volumes of data (Conti et al., 2006). To improve the detection of malicious events, and make best use of innate human abilities to integrate observations into explanations, there is a critical need for visual analysis techniques with which security professionals can efficiently interrogate this data. Moreover, visualization can be coupled with more sophisticated data pre-processing steps such that the human analyst is brought into the analysis loop at the appropriate point; rather than creating synoptic visual displays of all activity, for instance, we can first employ data reduction techniques that generate cognitively appropriate visual displays of the highest-value information.

### 2.1 *The Importance of Maintaining Context*

Human cognition is fundamentally about putting information in context, and it is context that allows analysts to make judgments about the meaning and importance of observed events. Just as real-world context guides pre-attentive selection of salient

features in the physical environment (Barsalou et al., 1993), contextual selection should extend to visualization environments. To help analysts deal with volumes of alerts and warnings, context-aware systems can triage these events in the same way a human would.

In an examination of the work practices of security analysts, Goodall et al. (2004) found that analysts often use online contextual information sources to keep up on the rapidly changing security landscape. Thompson et al. (2006) suggest that text information, including data from public information sources, be included in visualization applications for network security. We came to a similar finding during our own workshops (see Sect. 3), where analysts reported that there was more useful information online about potential threats than they had time or capacity to process.

## ***2.2 Visualizing Packets and Flows***

Many contemporary visualization approaches for network analysis problems focus on node connectivity and traffic flow (e.g., Paulson, 2004; Krasser et al., 2005) or on firewall alert and packet-level visualization (e.g., Conti et al., 2006). NVisionIP (Lakkaraju et al., 2004) exemplifies visualization of the lowest level of network flow data, providing a broad overview of traffic volume. Erbacher et al. (2002) produce an aggregate display that emphasizes higher-level behaviors, such as traffic load and connection patterns to individual systems on a network, and demonstrate that visualization can help analysts detect “interesting” cases quickly.

A primary concern with synoptic visualization of traffic, however, is that subtle nuances in traffic rate, destination, or type can go unnoticed. Kafadar and Wegman (2006), for instance, simply characterize “exotic” traffic as that with particularly high IP address or port frequencies, despite the risk that threats may also be carried in rare, sparse traffic. To address these subtleties, Wegman and Marchette (2003) call for “evolutionary graphics” capable of better communicating traffic change over time.

Ko et al. (1993) provide a foundation for tracking the temporal evolution of network actors, attempting to find the same actor in reports from multiple network resources. Our approach is to aggregate low-level log data to broader behaviors that characterize an actor (malicious or otherwise) over time – while allowing the end user to drill-down to individual records when needed. This technique reduces, in part, the data overload problem by increasing the level of abstraction at which network data can be analyzed. Our behavior models can also help express the temporal patterns and periods in traffic that are important – yet generally overlooked – in anomaly detection (Cordella et al., 2005). Leckie and Yasinsac (2004) take a similar behavior-centered approach, characterizing the activity of actors on a network based on the spread of their sessions during six time periods throughout the day; our technique, detailed in Sect. 4, uses a more precise continuous activity model.

## 2.3 *Visualizing Correlated Activity*

In response to the heterogeneous nature of data produced by distributed network sensors, attention has recently been devoted to visual environments that synthesize disparate information into coherent views. SecureScope (D'Amico and Larkin, 2001), for instance, helps situate cyber attacks in organizational and geographic context by visualizing features such as the mission components affected and the geographic location of critical assets. VisAlert (Foresti et al., 2006) visually links the “what, where, and when” of alerts from network sensors, although it depends wholly on the user to detect malicious activities. Hertzog (2006) visualizes relationships between users, applications, ports, and external hosts using parallel coordinate plots; large numbers of outgoing connections are grouped effectively in these graphics, allowing the analyst to readily detect cases where an application is communicating on an unusual port.

Some work on detecting correlates for cyber attacks is motivated by the subtle nature of insider threat. For instance, the vast majority of insider attacks demonstrate extensive advanced planning (Randazzo et al., 2004). Such cases exemplify the need for contextual information – in this case, information on behavioral changes, information access patterns, and so on – that can help detect nascent threats before they are executed. To this end, Stolfo et al. (2003) evaluate email flow characteristics (such as attachment frequency and time of use) to detect changing social cliques and policy violations. A natural extension of this work would apply text analysis techniques, as we do, to identify themes in traffic content that indicate concern.

From a data reduction perspective, there are existing techniques for correlation of multiple logs for improved anomaly and intrusion detection; EMERALD (Porras et al., 2002), for instance, fuses events at the alert level using Bayesian networks. Such approaches provide a lower-volume, higher-value information stream to the analyst. However, automated correlation of structured log data with unstructured contextual information remains a research challenge, and it is this problem that NUANCE helps address. Furthermore, despite recent advances in visualizing the state of network activity, most visual tools are still reactive – they are only effective at discovering attacks that are incipient. But there is a valuable role for visualization in helping analysts forecast potential events of concern in advance of their observation on the network; a system that processes online discussions can help analysts find out about security events that others are experiencing, for instance. Making use of real-time, open source reporting is especially important in effective notification of, and response to, zero-day attacks.

## 3 Technical Approach

To establish requirements for visual aids to contextualization, we engaged a group of seven practicing security professionals from our organization in brainstorming, design, and evaluation workshops over the course of a year. The goal was to

understand the kinds of “world knowledge” analysts apply, the nature of the external information sources they use, and their desiderata for visual interaction and discovery mechanisms. While the lessons derived from these sessions drove the design requirements for the NUANCE tools we present in Sect. 4, they also provide general guidelines for visual analysis in computer security applications. Below, we summarize the requirements derived from these sessions, which fell into two categories: understanding context and understanding behavior.

### ***3.1 “I Just Want to Know Where to Focus My Time”***

The primary concern that emerged from our workshops was that analysts need help looking beyond anomalies; they want to detect unwanted activity, which they distinguished from the merely anomalous. Said one participant, “my work is exploratory and creative by nature, not engineered or planned.” Signature-based approaches to anomaly and intrusion detection, while an important element of their arsenal, did not always support the need for open-ended discovery. And while the analysts had tools at their disposal that supported unstructured analysis, one challenge for them was knowing where to start looking; they needed the tools to incorporate “focusing” aids that suggested high-value exploration paths.

One of the primary drivers of exploration, we found, was the world of external information that analysts used to keep aware of the changing security landscape. Typically, analysts monitored sources like the SANS Internet Storm Center and US-CERT advisories, but each also defined an idiosyncratic set of online resources that they found helpful in identifying activities to be on the lookout for. Known malicious actors and traffic could be netted by existing signature-based tools, but they were concerned over learning about new exploits for which signatures did not exist.

Our analysts suggested that one technology that would benefit them was a way to associate key words in the online text sources they read with events in the log files they explore. They wanted to be able to quickly find, for event patterns of interest in their logs, any available information online that could help explain them. They suggested that “off-network” text information is even more important in the area of insider threat. Much of the forewarning of insider events was described as “soft”: reports from other staff or changes in the style, content, or pattern of communication. As with external threats that might begin with a reconnaissance phase, insider activities are marked by a number of elements that happen before the event, none of which would be captured in the logs of current monitoring tools. The same content-based approach to associating terms in external text sources with network events could be used for insider threat mitigation if it is applied to network traffic content. While insider threat detection is outside the scope of the present work, the basic techniques we develop could be extended to this area.

A motivating example from finance emerged from one of our workshops to illustrate the role of context in visual analysis. A company’s stock price is a summary signal for the behavior of the people and institutions that own it. Changes in the

signal reflect changes in behavior – buying and selling. News stories about that company can serve as both *indicators* and *drivers* of changes in the stock price signal. As an indicator, for instance, a news article about a company's woes might coincide with an observed downward trend in its stock price. As a driver, an article about the unexpected release of a new product might motivate buyers to purchase the stock, raising its price. And just as a market analyst cannot effectively act if he or she only sees one or the other component of the system (either the price history or the contextual articles), a security analyst needs both signal and context to understand the state and trend of the network. In particular, the better the analyst can understand the relationships between drivers and the observed signal, the better decisions he or she can make; ideally, the analyst can identify contextual events likely to influence the signal before they are reflected in that signal. To proactively respond to potential threats, the security analyst needs to be aware of the current threat landscape and must understand the relationships between those threats and his or her own systems. Visualization can be effective in communicating associations between the two.

### 3.2 “We Need to Organize Our Hay into Smaller Piles”

The second major area of concern for analysts in our workshops was developing a greater understanding of the nature of the behaviors on their networks. Unsatisfied with tools that gave them a broad overview of the “haystack” that was their network, our analysts were looking for data reduction techniques that increased the information density of their visual displays. This “smaller haystacks” approach, if coupled with contextual information that could at least point them to the right haystack, as they said, would result in efficiency gains in finding the needles of malicious activity.

The first approach to data reduction that the analysts wanted to see was the ability to reduce massive amounts of transaction level data down to a set of behaviors that represent the trends in those transactions over time. These behavior-based views would help them see what normal activity from an individual IP address, group, or location looked like, and would therefore let them detect off-normal conditions more easily. A visual display should show the analyst whether current activity for given actor is within its expected behavior, but can do so at a high level that does not clutter the display with lower-level transaction information (although this should be available through drill-down).

Once the baseline models had reduced large amounts of transaction records to a more succinct behavior, the analysts were looking for tools that would organize these behaviors into “cliques” that would help them detect distributed but related activities. These cliques would be produced at multiple levels of granularity, from the individual IP address, to the organization (company, university, and so on), and region (city, province, country). Organizational cliques can help detect coordinated activities that derive from a location such as an internet café, even if the particular IP addresses involved change over time. In addition to creating cliques on the basis of

shared network attributes, cliques should be created that reflect common behaviors regardless of their location or organization of origin.

Understanding both behavior and context helps analysts achieve what they stated as their ultimate goal, which was to be more proactive in their work. A challenge for visualization, they told us, was to give them the flexibility to be both predictive and reactive; to look ahead for events they *might* experience and to look backward for forewarning they could have had about events they *did* experience.

### 3.3 Behavior Modeling

The NUANCE visualization approach that emerged from the requirements gathered during our analyst workshops fuses behavioral analysis with contextual information. This section describes the resulting behavior modeling and visualization technique; Sect. 5 introduces our contextual analysis methods.

Rather than build visualizations for individual network transactions, our work takes the approach that what is important is anomalous or malicious *behaviors*, which may be manifested in a series of transactions over time (from seconds to days or even years). NUANCE uses a hierarchical modeling approach to represent actor behavior. At the most granular level is the IP address; a behavioral model can be constructed for every IP and group thereof. NUANCE defines group membership through geographic and organizational attributes retrieved through who is and gazetteers, although custom groups could be created by the analyst. (We note that NUANCE itself does not reconcile spoofed addresses, although it can make use of existing traceback techniques). Groups can also be created at the port level. Sample groups might include “China”, “Zhejiang Province”, “XYZ University” and “FTP traffic from XYZ University”. When a new actor appears on a network (e.g., an IP address that has not been observed before), the hierarchical modeling approach allows it to be assigned baseline behavioral models based on its organization and region until sufficient history has been observed to generate a model of its own.

Because baseline activity for “normal” behavior is not assumed to be the same for all IP addresses, the analyst can examine behavioral profiles aggregated to any group of users – important when trying to detect external threats where the attacker changes IP address over time. We can look at behavior from a particular place, for instance, and note when that origin location’s behavior changes, regardless of which actors are responsible for the change (since some of the actors may be “new” IP addresses).

#### 3.3.1 Model Definition

The NUANCE behavior modeler consists of parsing, statistics, and curve-fitting components. The parser receives real-time network data (such as from NetFlow logs, although any log data can be used) and summarizes each record as a timestamp and a list of groups to which the transaction belongs. Group membership is assigned by

rules such as “Is this packet from a .edu domain?”, “Is this an HTTP packet?”, or “Is this an HTTP Packet from ABC Co.?” The statistics thread monitors output from the parser and maintains an array of sufficient statistics for each IP address and group (hereafter we use the term “actor” to describe both unique IP addresses and groups thereof). For a given actor, our statistical model expresses the expected traffic rate over time as a periodic function. The assumption of periodicity suggests the use of a Fourier series, and to ensure that the expected traffic rate is never negative, we use an exponentiated Fourier series.

Let  $k$  be a vector of periods of interest, and let  $K$  be the least common multiple of the  $k_i$ 's; we currently use  $k = (6, 8, 12, 24)$ , in units of hours, although  $k$  can include any time periods of interest, from minutes to years. Our model is then that session time modulo  $K$  is random, with a density function we will now describe.

For any real-valued parameter vectors  $\alpha$  and  $\beta$ , we can write the series

$$\varphi_{\alpha\beta}(t) = \sum_i \left( \alpha_i \cos \frac{2\pi t}{k_i} + \beta_i \sin \frac{2\pi t}{k_i} \right).$$

Since  $\exp\{\varphi_{\alpha\beta}(t)\}$  is always positive, we can normalize it to arrive at a density function

$$f_{\alpha\beta}(t) = \frac{\exp\{\varphi_{\alpha\beta}(t)\}}{I_{\alpha\beta}},$$

where  $I_{\alpha\beta}$  is the normalizing constant

$$I_{\alpha\beta} = \int \exp\{\varphi_{\alpha\beta}(t)\} dt.$$

Every choice of  $\alpha$  and  $\beta$  leads to a density  $f_{\alpha\beta}(t)$ . Now suppose for a given actor we observe transactions starting at times  $\{t_1, \dots, t_N\}$ , and we wish to choose  $\alpha$  and  $\beta$  to best fit the observed data. This we can do using maximum likelihood estimation; we consider  $\Pi_n f_{\alpha\beta}(t_n)$  as a function of  $\alpha$  and  $\beta$  (since the  $t_n$ 's are now fixed), and find the maximum value. The values of  $\alpha$  and  $\beta$  leading to this maximal value are then the maximum likelihood estimates of the parameters, and are denoted  $\hat{\alpha}$  and  $\hat{\beta}$ , respectively.

Typically, rather than maximize  $\Pi_n f_{\alpha\beta}(t_n)$  directly, it is more convenient to maximize its logarithm. This leads us to the log likelihood

$$\begin{aligned} l(\alpha, \beta) &= \sum \log f_{\alpha\beta}(t_n) \\ &= \sum \varphi_{\alpha\beta}(t_n) - N \log I_{\alpha\beta} \\ &= \sum (\alpha_i c_i + \beta_i s_i) - N \log I_{\alpha\beta}, \end{aligned}$$

where  $c_i$  and  $s_i$  are the statistics

$$\begin{aligned} c_i &= \sum \cos \frac{2\pi t_n}{k_i} \\ s_i &= \sum \sin \frac{2\pi t_n}{k_i}. \end{aligned}$$



In particular,  $(N, (c_i), (s_i))$  constitutes a sufficient statistic. Thus as new data arrive, the number of values that must be stored for estimation remains constant.

For each actor, the sufficient statistics allow us to fit a curve representing the traffic rate as a function of time, with multiple time scales taken into account simultaneously. Once the density function is determined, we can use it to estimate the expected traffic rate at a given point in time. Suppose observed times span from  $t_0$  to  $t_1$ . Our estimated rate can be found using the criteria

$$r(t) = \gamma f_{\hat{\alpha}\hat{\beta}}(t)$$

$$\int r(t)dt = N.$$

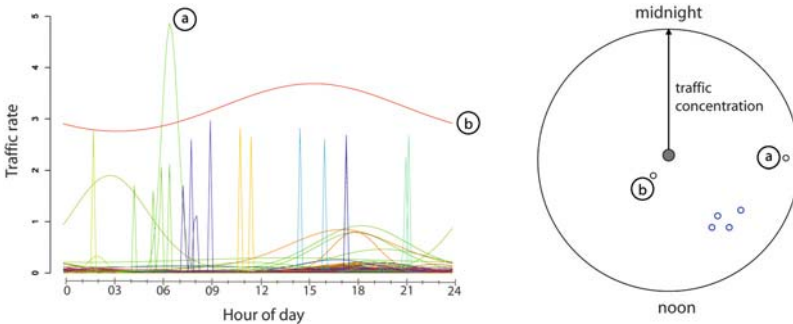
Here  $\gamma$  is a normalizing constant. We then have

$$N = \int r(t)dt = \gamma \int f_{\hat{\alpha}\hat{\beta}}(t)dt,$$

so the estimated traffic rate for each actor is

$$r(t) = \gamma f_{\hat{\alpha}\hat{\beta}}(t) = \frac{N f_{\hat{\alpha}\hat{\beta}}(t)}{\int f_{\hat{\alpha}\hat{\beta}}(t)dt}.$$

The NUANCE statistical methods result in evolving behavioral models for each actor on a network. Figure 1 shows a sample of behavioral profiles using this modeling approach. Each colored curve represents a unique actor. Some actors demonstrate characteristically “bursty” behavior, such as the profile labeled “A”, engaging in short sessions of traffic at various times throughout the day. Others exhibit continuously high levels of traffic, such as the profile labeled “B”; this is characteristic of a server, or, if the IP address is external, of a search engine crawler. A low-and-slow port scan would appear “bursty” in this model, because the traffic would appear as short periods of activity separated by long periods of quiet.



**Fig. 1** (left): Diurnal behavior profiles (traffic over time) for a selection of actors; each curve represents a distinct actor, whose activity over all ports is aggregated (port-specific models can be constructed). Y-axis scale is packets per second. (right): Visual transformation of profiles to “clock” view allows cohorts of similar actors (blue) to be detected

In the special periodicity case  $k = (2\pi)$ , an actor's traffic function reduces to the density for a von Mises distribution, which is a continuous distribution describing a set of points situated on a circle. This von Mises plot (at the right of Fig. 1) represents temporal information, with time moving clockwise around the plot, and provides an even higher-level summary of the nature of each actor's behavior. Actors who typically engage in brief bursts of activity (e.g., A) plot nearer to the perimeter of the circle, while those who are more continuously active over the course of time period being examined (e.g., B) plot closer to the center. The position of the actor around the "clock" indicates its typical peak activity time (e.g., 6:00 am for actor A). The von Mises plot in Fig. 1 shows 24 h of activity, but shorter or longer time periods can be selected. The von Mises distribution dramatically reduces the storage and processor time required to generate behavioral models (only three values are required to express each actor's diurnal traffic patterns), allowing models to be constructed on the fly for large segments of the IP address space. Moreover, actor models can be clustered with K-means. These clusters represent behavioral cliques, and contain actors who, regardless of organization or geographic location, exhibit similar behavioral profiles. This clustering approach is useful for detecting distributed attacks that show the same "modus operandi" but would not otherwise be associated into the same actor group. It is also possible to detect actors who change behavioral clique over time, which is often an indicator for insider threat.

We have currently modeled up to 100,000 unique actors successfully in real-time on commodity hardware, and are working to increase by an order of magnitude the number of models the system can process. One limitation to the number of models that can be accurately constructed is that the curve fitter cycles independently over the actors, updating fitted curves as it goes. As new transactions arrive, we walk the model list and tune each model in sequence. As a result, as the number of behavioral models to be stored increases, the time to complete the walk increases and the accuracy of the resulting models decreases. This limitation can be mitigated by parallelizing the curve fitter.

### 3.3.2 Dynamic Histograms

To detect off-normal behavioral conditions, we compare the activity models described above with empirical traffic rates. When observed traffic varies from the predicted behavioral model by a user-specified threshold, the difference can be visualized for the analyst and appropriate contextual information retrieved. A natural approach for representing observed traffic rates is to use a histogram of transaction (e.g., packet or session) counts. However, in storing these histograms there is a trade-off between high resolution and long memory. As with the modeling component, we aim to maintain constant space (i.e., not infinitely increasing the storage and processing requirements as the amount of data increases), so we require the number of bins in the histogram to be fixed.

To accommodate this trade-off, we maintain high temporal resolution for recent data and drop the usual implicit assumption that bins are adjacent. Each bin in a

NUANCE histogram represents 1 min of activity from each actor for whom a model is also being built. Each minute a new bin is introduced and an older bin chosen at random is dropped. Thus the number of chances a bin has had to be dropped is proportional to its age, so older bins are progressively farther apart. The result is a time-decayed histogram, representing more recent activity on a finer time granularity than older activity. Our approach allows observed data to age in this way with very little overhead.

### ***3.4 Building Context***

NUANCE also introduces a new method for associating contextual content with network behaviors. This process involves gathering and filtering text data to construct a vocabulary that describes each actor and group being modeled, a monitoring component that collects real-time content from sources of the analyst's choosing, and classification routines that determine, on the basis of their vocabularies, to which actors incoming content is relevant.

#### **3.4.1 Vocabulary Construction**

Using the same metadata by which IP addresses are assigned to groups (organization name, geographic area, and so on), NUANCE actively constructs a text vocabulary to describe each IP and group. The vocabulary is generated by performing an automated web search around those metadata terms; the result of this harvest is a corpus of training documents representative of the actor (alternatively, training documents can come from the content of traffic in which that group is involved). These documents are merged into a single term list, processed for stopwords (terms to exclude) and major terms (frequent, statistically descriptive terms), and converted into a vector that represents topics characteristic of the group. The value of the vocabulary construction approach is that we can make associations between text documents such as security notices and group activity on the basis of more than just keyword matching. That is, we do not look for the term “Philippines” in a news article and automatically associate that article with our traffic model for the Philippines. Instead, the web harvest builds a more complete representation of the topics relevant to the Philippines (and can be tailored to favor security-related web sites over general information sites). The NUANCE vocabulary allows notices of exploits to be associated with actors that might show evidence of using that exploit, even if the notice does not mention that actor at all (as is usually the case, when we do not know from what IP addresses to expect an attack). There is a significant bootstrapping challenge in generating these vocabularies, however, as one must be constructed – and updated – for each IP address and actor group being modeled.

### 3.4.2 Content Analysis

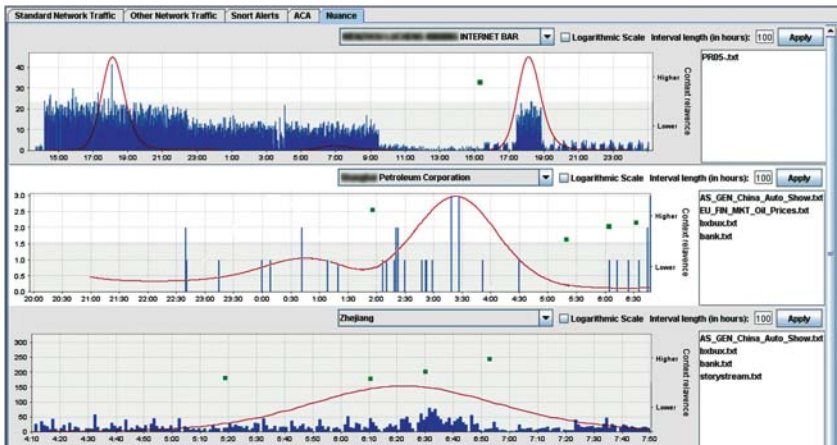
Once baseline vocabularies have been constructed for each group NUANCE has also created a behavioral model for, the user can specify a set of text feeds to monitor for computer security and geopolitical events of interest. Currently, NUANCE collects text reports every hour from a variety of security and news web sites. These sources include US-CERT, SANS Internet Storm Center, PacketStorm advisories and exploits, online bug trackers, BBC News, New York Times, and a selection of user-specifiable sites such as message boards. Typically, NUANCE ingests about 300 text reports a day, although this is a minimum and can grow as the user adds new sources to the NUANCE “reading list”. The time of publication for each text document is extracted and later used to correlate the document with network events. Just as each document in the training set for a group vocabulary is processed into a term vector, each incoming context document is processed into a similar vector.

### 3.4.3 Calculating Contextual Relevance

Once vectors for both characteristic group topics and incoming textual events have been created, they can be compared to determine to which groups context events are relevant. Measuring the relevance of a feed relative to a group is based on a cosine similarity metric. Each time a new contextual feed item is received, it is compared against each of the group vocabularies currently in the system. The smaller the angle between each feed item and the group vocabulary, the more similar or relevant the feed is to the group. We currently use a pre-defined similarity score threshold of 80% to determine whether a context item should be associated with a group. It is possible for a given context item to “hit” on multiple groups, meaning that there are potentially multiple behavioral events to which the context alert is relevant. It is also likely that for many incoming context documents, there will be no relevant actor; the goal is simply to create a broad harvest of current security topics so that if a relevant actor appears on the network, the analyst can be notified.

## 3.5 Visualizing Behavior in Context

Figure 2 shows NUANCE in operational mode. The main application allows the analyst to vertically tile and sort multiple behavioral models (three are shown in this view), facilitating comparison of activities across groups. Selections from an analyst’s “favorite” behavior models are made through a drop-down list at the top of each pane. When the system monitors large numbers of actors, analysts can call up models through a search interface. NUANCE currently defines actors automatically, aggregating transactions by IP address or group of addresses representing a geographic or organizational unit. Analysts do not have to define these actors manually, but we anticipate creating a visual interface that allows them to do so. Models



**Fig. 2** NUANCE behavioral models for three actor groups. Histograms (blue) show current activity in 1-min intervals. Red curves represent expected behavior for each group. Relevant contextual information is automatically attached to the group at the time it is received (green dots). Selecting a context item in the chart or the on the list to the right of each plot opens it

can also be “pushed” onto the view based on pre-defined anomaly thresholds (i.e., observed behavior differing from expected behavior by more than a specified percentage) coupled with availability of context (i.e., groups with large anomalies plus highly relevant context).

Each NUANCE plot is essentially a timeline showing a zoomable temporal window on the *x*-axis and traffic level over that window on the *y*-axis (here, shown as packets per minute over all ports, although views could be segregated by port). The top plot in Fig. 2 shows approximately the last day and a half of traffic from a particular internet café. Histogram data shows per-minute observed traffic levels, while the red periodic curve represents the current expected behavior pattern for that group. In this example, NUANCE has been processing real-time data for about 3 days, so the models are beginning to tune themselves but have not yet reached full precision (we have found that approximately one week’s worth of data is required to maximize the modeling approach’s ability to represent observed traffic). For traffic from the internet café, the model correctly predicted the two peaks in traffic that occurred over this period, although the model did not predict the generally high level of activity over the first half of the plot. This difference can trigger a flag which adds the group to the analyst’s watch list. At the same time, incoming context feeds that met the scoring threshold for relevance to this group are displayed as green dots on the plot. The right *y*-axis describes the approximate score of each context item, with more relevant items plotting toward the top of the chart. In the case of this internet café, we can see that one contextual item has fused to the group; selecting that item (either the dot on the chart or the item title in the context list to the right of the chart) opens a reader where the analyst can view the context item. This particular context feed warned of a SQL injection attack, and it fused to this internet café because

the vocabulary for that group indicated that similar exploits had been identified as sourcing from the same location. The analyst, without knowing that traffic from this source should be of concern, now has both an understanding of how the traffic has departed from recent historical levels and, through the context item, an explanation of what that traffic may involve. Drilling down into the histogram can expose the underlying transactional information.

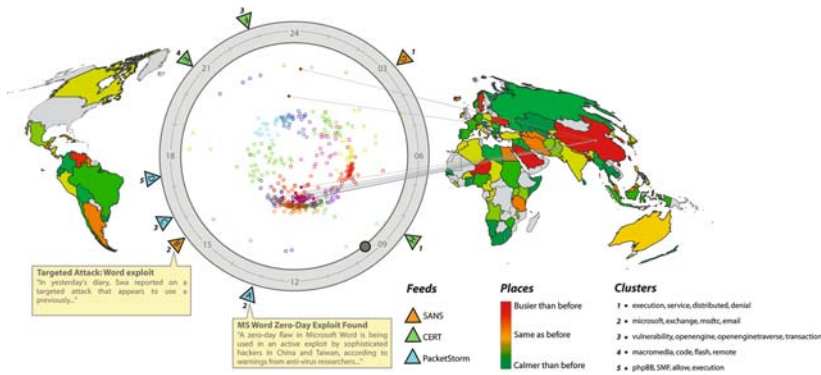
The behavioral models in the lower two charts (for traffic from a petroleum company in the middle plot, and from Zhejiang province in China for the lower plot) generally track their group's observed activity, and will continue to refine at each update increment. Here too, context elements help explain changes in the observed activity rate for each. In the middle chart, a news alert about rapidly rising petroleum prices in the region in which this company does business (leftmost green dot, corresponding to the "EU\_FIN\_MKT\_Oil.Prices" story) precedes a rise, roughly 90 min later, in this organization's activity. Just as news items can help explain observed changes in stock prices, our analysts found that contextual items could often help them determine why activity changes might be taking place, leading to a better ability to triage responses.

### 3.5.1 Situational Awareness Dashboard

Summative evaluation from our analyst team on the NUANCE application suggests that the ability to visually relate context and actor behavior fits best in an *exploration* phase of analysis. That is, analysts can drill down through these charts to uncover the original transaction-level data that went into the histogram and modeling routines and can "play" with the data by comparing actor charts, looking up models for actors of interest, and scrolling forward and back in time to uncover subtleties in traffic.

However, analysts in our workshops wanted to complement this exploratory interface with a similar context-driven visual tool that would serve them in their *monitoring* activities. For many of our analysts, monitoring preceded exploration; they might use a suite of monitoring tools to alert them of suspicious activity over the course of the day, and if an event merited further investigation they would turn to an analysis tool to explore it further. As a result, we developed a prototype monitoring interface that reduces the NUANCE contextualization and behavior modeling approach to a simpler "one-look" view.

Figure 3 shows the enhanced NUANCE view designed to support real-time monitoring. At the center of the display is a von Mises circle depicting the past 24 h of activity (shorter and longer time periods can be selected), containing a point for each actor. Through the method described in Sect. 4.1, each point corresponds to one of the time charts in Fig. 2; the two displays can be linked such that selection of an actor in this view will bring up the relevant histogram and activity model. The location of each actor around the plot perimeter indicates its expected time of appearance, while its location along the radius indicates the nature of its activity, from "bursty" to continuous. The actor models in Fig. 3 have been clustered using K-means; each cluster is represented by a unique color in the central circle.



**Fig. 3** NUANCE situational awareness tool. Von Mises circle at center shows a point for each actor observed over the last 24 h; colors indicate behavioral cliques. A geographic region (China) is selected, highlighting actors from that location

Surrounding the time wheel is a map that shows the analyst the level of “chatter” about geographic locations in the context feeds being monitored. As more context items mention a location, its country is colored a deeper red. In this example, China has been selected, and the actors representing traffic from China have been highlighted in the circle. Context items relevant to this geographic group then appear as flags around the perimeter of the time wheel at the time they were published. Context items are also clustered in real-time, using standard text-clustering techniques (Hetzler et al., 1998); each flag is numbered with the context cluster to which it belongs. For instance, the two highlighted context items referring to an MS Word exploit correspond to cluster 2, which contains the “Microsoft, exchange, msdtc, email” topics.

The dashboard view simplifies user interaction by restricting the amount of detail shown about each actor. By reducing each behavioral profile to a single point, more actors can be shown at once. In a single view, it is also possible to ascertain both the current geopolitical and security landscape. We are currently exploring visual interactions that help users link high-level concepts derived from the dashboard with the specific activities (as in Fig. 2) that manifest these changes.

4 Future Work

NUANCE behavioral modeling is suited equally well to modeling the activities of external hosts or machines internal to a network. In the internal case, NUANCE models can offer analysts an evolving picture of the expected state of machines on their network. One challenge in modeling external hosts is the large number of potential actors of interest; we do not currently have an adequate mechanism for enabling analysts to navigate through the models that exist and choose which



to display. An additional visual component that helps in this regard (for instance, clustering related models to provide hierarchical navigation) may be called for.

We are currently working on extending NUANCE to support predictive analysis. Once NUANCE has linked a context item with a behavioral profile, it is possible to use machine learning techniques to reinforce the association between the context topics and resulting behaviors. Then, when particular collections of context topics are observed in the future, analysts can be presented with representations of the kinds of behaviors to expect. Associating behavioral profiles with historical security events helps detect those activities that are not in themselves anomalous but that match the longer-term patterns of actors who eventually did perform malicious acts. We are also developing techniques to visualize higher-order threat models that link behaviors observed across distributed actors. This work helps overcome the difficulty in detecting multi-stage attacks emanating from spoofed sources. Finally, to address the challenges inherent in processing streaming network data in real-time, we have begun implementing NUANCE on high-performance computing infrastructures. Distributing the model generation and text classification work across a cluster can improve the accuracy of the models (by reducing the time between updates) and speed the generation of training vocabularies for contextualization as new actors are observed.

## 5 Conclusions

Based on requirements gathered during a series of formative sessions with practicing analysts, we have developed a new behavioral modeling and contextualization approach that helps users visualize the associations between changes in network activity and explanatory external events. Our efforts to build automated context gathering and classification methods into visualization tools give analysts an improved ability to situate threats in the real world, a practice in which they already engage but for which current tools offer little support.

The NUANCE method's behavioral modeling technique is capable of representing the periodic nature of network activity over multiple time scales, even when periods are not evenly spaced. NUANCE can also create a unique behavioral model for each actor observed on a network, allowing it to present a detailed picture of network activity. Assessing the deviation of each actor's current activity from historical precedent offers specific warning of suspicious activity.

Behavioral models are one technique for accommodating surprise, by helping the analyst understand whether what he or she is seeing fits into historical patterns. However, these models alone will never produce perfect forecasts, especially when dramatic events force a change in behavior (for instance, Hurricane Katrina fundamentally changed the nature of traffic to and from certain regions of the US). Therefore, we incorporate a text content collection and fusion technique that helps analysts discover possible reasons why traffic is behaving as it is. Ultimately, by fusing heterogeneous information sources – including those, such as news feeds,



that are part of the analyst's toolkit but typically ignored in security analysis tools – we can improve the analyst's ability to detect and respond to threats.

## References

- Barsalou L, Yeh W, et al. (1993) Concepts and meaning. Chicago Linguistics Society: Papers from the Parasession on Conceptual Representations. K Beals, G Cooke, D Kathman, et al. Chicago, University of Chicago 29:23–61.
- Conti G, Abdullah K, et al. (2006) Countering security information overload through alert and packet visualization. *IEEE Computer Graphics and Applications* (March/April):30–40.
- Cordella LP, Finizio I, et al. (2005) Using behavior knowledge space and temporal information for detecting intrusions in computer networks. *Pattern Recognition and Image Analysis, Pt 2, Proceedings* 3687:94–102.
- D'Amico A and Larkin M (2001) Methods of visualizing temporal patterns in and mission impact of computer security breaches. *DARPA Information Survivability Conference and Exposition II*, IEEE Computer Society.
- Erbacher R, Walker K, et al. (2002) Intrusion and misuse detection in large-scale systems. *Computer Graphics and Applications* 22(1):38–48.
- Foresti S, Agutter J, et al. (2006) Visual correlation of network alerts. *IEEE Computer Graphics and Applications* 26(2):48–59.
- Goodall J, Lutters W, et al. (2004) I know my network: Collaboration and expertise in intrusion detection. *Proceedings of the ACM Conference on Computer-Supported Cooperative Work, CHI Letters*, ACM Press, New York, pp. 342–345.
- Hertzog P (2006) Visualizations to improve reactivity towards security incidents inside corporate networks. *ACM Conference on Computer and Communications Security, Third International Workshop on Visualization for Computer Security*, Alexandria, VA, ACM Press, New York, pp. 95–102.
- Hetzler B, Harris W, et al. (1998) Visualizing the full spectrum of document relationships. *International Structures and Relations in Knowledge Organization, Proceedings of the Fifth International ISKO Conference*, Wurzburg, ERGON Verlag, pp. 168–175.
- Kafadar K and Wegman EJ (2006) Visualizing “typical” and “exotic” Internet traffic data. *Computational Statistics & Data Analysis* 50:3721–3743.
- Ko C, Frincke DA, et al. (1993) Analysis of an algorithm for distributed recognition and accountability. *Proceedings of the First ACM Conference on Computer and Communications Security*, Fairfax, Virginia, United States, ACM Press, New York.
- Krasser S, Conti G, et al. (2005) Real-time and forensic network data analysis using animated and coordinated visualization. *IEEE Workshop of Information Assurance and Security*, West Point, NY.
- Lakkaraju K, Yurcik W, et al. (2004) NVisionIP: Netflow visualizations of system state for security situational awareness. *ACM Conference on Computer and Communications Security, Workshop on Visualization and Data Mining for Computer Security*, Washington, DC, ACM Press, New York, pp. 65–72.
- Leckie T and Yasinsac A (2004) Metadata for anomaly-based security protocol attack deduction. *IEEE Transactions on Knowledge and Data Engineering* 16(9):1157–1168.
- Paulson L (2004) Researchers develop network-security visualization tools. *Computer* 37(4): 17–18.
- Porras P, Fong M, et al. (2002) A mission-impact-based approach to INFOSEC alarm correlation. *Fifth International Symposium on Recent Advances in Intrusion Detection (RAID 2002)*, Zurich, Switzerland, Springer, Berlin Heidelberg New York, pp. 95–114.
- Randazzo M, Keeney M, et al. (2004) Insider threat study: Illicit cyber activity in the banking and finance sector, U.S. Secret Service and CERT Coordination Center.

- Stolfo S, Hershkop S, et al. (2003) Behavior profiling of email. First NSF/NIJ Symposium on Intelligence and Security Informatics, Tucson, AZ, Springer, Berlin Heidelberg New York, pp. 74–90.
- Thompson R, Rantanen E, et al. (2006) Network intrusion detection cognitive task analysis: Textual and visual tool usage and recommendations. Proceedings of the Human Factors and Ergonomics Society (HFES 06), pp. 669–673.
- Wegman EJ and Marchette D (2003) On some techniques for streaming data: A case study of Internet packet headers. *Journal of Computational and Graphical Statistics* 12:893–914.