

NetVis: A network traffic visualization tool

SUBMISSION ID
AFFILIATION ID

May 24, 2013

Abstract

Computer network traffic visualizations attempt to deliver improved understanding of traffic on a network to an observer. Many existing tools opt for graph or plot-based visualizations to detect patterns or outliers in the data, but still largely provide a segmented view of any data feed. In this paper, we present a novel network traffic visualization framework that makes use of a variety of complementary visualizations to obtain better situational awareness. Our proposed solution is to look at different..

1 Introduction

Analysing and monitoring network traffic is essential for ensuring network security. Since the amounts of data involved in this task are extraordinarily large, it can be hard to make effective use of them. Visualizations offer a solution: they give a compact representation of data to a human analyst who can use them to detect patterns in the network. An effective visualization of network traffic will make it easy to identify outliers while making clear the general patterns of network usage, and it will facilitate detection of intrusions and of malicious activity.

We present here an application which provides a framework for analysing network activity in an interactive and holistic [**better word?**] manner. Time series data in the form of packet captures are processed in real-time, and simultaneously rendered in multiple interconnected visualizations. The user can switch between the available representations of the data and dynamically influence both the visualization's layout and the amount and type of the data displayed. The software is designed to make it easy to spot irregular activity and then to investigate it from multiple angles.

The framework developed here aims to provide an environment in which situational awareness can more effectively be obtained. It is therefore built to be easily extendable to fit the demands of the specific situation. New visualizations will in a natural way complement the existing ones: The underlying data processing engine provides a standard, or *normalised*, data basis which is shared by all displays.

The purpose of the tools presented in this paper is to assist a human observer to make sense of what is happening in a given network. They emphasize exceptions, show comparisons, and try to answer a wide array of question in a concise and simple fashion. We want the user to generate good hypotheses in response to the visualized information. To achieve this, many aids are given to the user in order to avoid cluttering displays with irrelevant data. In particular, we have implemented a powerful filtering system which can be dynamically adjusted in response to changing situations.

The application's main purpose is to monitor the current activity in the network. However, the same architecture can be used to forensically analyse recorded data. The system will simulate the data records as if they represented a live network. This exploits the important dimension of time, and makes it easier to understand what is actually happening.

2 Related Work

JH

3 NetVis Architecture

Add information on CSV, Data, non-existence of analysis of packet content, and add info on data flow through the application.

Normalisers: A normaliser is a class responsible with taking a packet and returning a normalised value of a certain attribute of that packet and the other way around. Each normalising class is able to create a temporary filter in the application that will filter its corresponding attribute on a certain range. Since the normalising class is in control of its filter it creates a zooming effect based on the range of the filter: the lower bound is normalised to 0, the upper bound to 1.

Flow of Data through the application: The data feeder creates Packet objects based on the traffic in analyses and feeds it to the Data Controller. We implemented a CSV data feeder that reads CSV files (/ * maybe here a little more about the structure of the CSV's? */) created with Wireshark but the simple interface of the data feeder makes it easy to extend the application for reading live data.

The Data Controller uses the listener pattern and feeds new packets as they arrive to its listeners (visualizations, analysis panel). If there are any active filters then the data controller will apply those filters to the data and supply only what is relevant.

When a filter is applied or changed a reset signal is sent from the data controller informing the visualizations that all the data has changed. Along with this signal a list of all the past packets (with filters applied) is provided.

Filtering: **High-level: What do filters do?** The application supports two types of filters - filters that the user explicitly defines, and filters applied on-the-fly from within the

visualisations. Both types are applied to all visualisations and information displayed, and can be adjusted at any time without losing prior data. The user has access to a variety of different filter controls. First, there is a menu for filtering by transit protocol. By default, all protocols are selected and therefore included. Protocols are sorted into menus by protocol family, appearing in multiple places where appropriate. Second, there is a control to select the range of ports the application uses, which defaults to the maximum port range. The source and destination ports can be set separately, enabling a user to view all data entering/exiting a port as desired. Next, there are IP and MAC address filters. These work on a blacklist/whitelist system, allowing a user to only view packets to or from a particular set of addresses, or ignore packets going to or from a different set. This enables a user to, for instance, ignore traffic from sources they know are irrelevant or focus only on an address that is causing concern. The second type of filter will be dealt with in more detail in sections 4.2 and 4.3.

Modularity: **[Make this better.]** The application was designed to be very extensible. Adding a new visualisation is simply a matter of adding it to the application's list of those available, which will cause it to be included and kept up-to-date as packets come in and are filtered. The same is true for filters and normalisers. **[wrong place:]** Adding a filter would result in its controls automatically being included in the right panel and all packets would then be filtered according to the criteria it specifies. Adding a normaliser would cause it to be integrated with the Spinning Cube, Dataflow and Attribute Distribution visualisations. This modularity extends even so far as data input. If a class were written to accept packets from a different source, it would be trivial to switch the application to use this class.

4 NetVis Visualizations

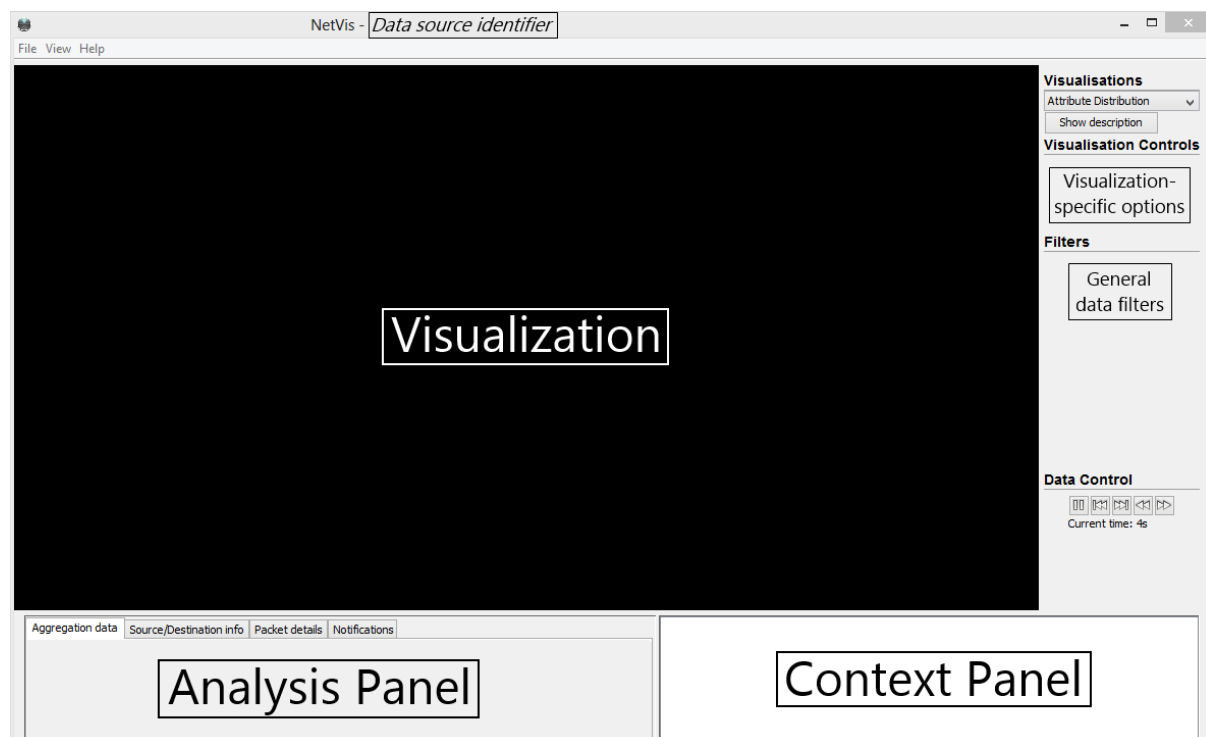
The user interface of the application focuses on the displayed visualization, which takes up the majority of the window, and can be maximised to fill the available space. Other panels occupy the right and bottom sections of the window, and are as follows.

Choice of visualization is presented to the user by a list in the top of the right panel, followed underneath by options specific to that visualisation, and then by general filters to refine the processed data. If a recorded data source is active, time controls are displayed on the right panel to allow the user to fine-tune the speed at which data is processed and visualised, by pausing, doubling and halving the current data rate. The right panel is ultimately aimed at providing the user with options to adapt the shown data.

The function of the bottom panel is to show fine and aggregate details of the processed data. To fit the large amount of data available into such a limited space, a tabbed pane is used on the left to categorise distinct types of data. The right section of this panel contains a 'Context Pane', the purpose of which is to show further detail of selected data on demand. The left 'Analysis' panel and the right 'Context' panel are split in two by a split pane, so the user can modify the proportion of each they wish to see.

Window tools, an exclusive-mode full screen option and the facility to select a new data

source are all accessed by the main menu bar. Messages to the user come in the form of warning dialogues for user or program errors, notifications in the bottom ‘Analysis’ panel for program notifications, and text messages in the bottom ‘Context’ panel for help messages and suggestions. A diagram of the application layout is included below.



4.1 Attribute Distribution

Displays a graph of the distribution of some packet attribute using its normalised value. It improves upon existing distribution visualizations by offering a dual representation. There is a logarithmic line graph layer useful for detecting spikes in the data and a layer of circles that indicates the actual distribution (no logarithm applied) through their area. Everything is color coded to give an extra indicative of the volume of traffic (red means bigger than blue). All these ideas have a single purpose: make the analyst understand what is going on. // maybe you can change this into something that makes more sense (combine different ideas into one visualisation...) In this visualisation you can select a range for an attribute and the visualisation will apply a removable data filter that will only allow data in that range to get to the visualisations. [make nicer]

4.2 Dataflow

The Data Flow visualization applies the idea of parallel coordinates proposed by Inselberg1985plane to network traffic. In a similar fashion to the work of Keith Fligg and Genevieve Max, the visualization shows the ‘flow’ of each packet, representing each as

a line through the parallel coordinates. This provides an informative view of the whole traffic in the network. It visualises the distribution of various packet attributes while also giving an intuitive insight into relationships and correlations between distinct aspects of the packets.

Lines representing packets are coloured based on their value in some coordinate. They fade out based on how old the packet is, thereby placing focus on newer developments in the network.

A common criticism of parallel coordinate plots is that they make it hard to interpret data that is uniformly distributed between coordinates or that is very concentrated around few values (applied2004). Our implementation addresses these concerns using animation. Lines representing new packets are randomly perturbed and move subtly. This makes it easier to spot whether a colourful line represents one or more packets. The animation allows the user to distinguish between single packets and concentrated groups of packets with the same characteristics.

For example, if multiple requests to a server came from a single MAC address through the same ports and using the same protocol, in a non-animated visualization they would all be represented by the same line. A user would erroneously interpret this as a single request. Adding randomness makes the pack of requests more visible without significantly impacting the accuracy of the data representation.

Hovering over one coordinate displays a scale for the attribute it represents. Clicking transfers you to the distribution visualisation showing the distribution of that particular attribute.

4.3 Spinning Cube

The Spinning Cube is an implementation of an existing visualization tool known as the spinning cube of potential doom¹. It displays packets as dots inside a spinning cube. Their 3D position and color is based on 3 packet attributes. (eg. x-axis \Rightarrow Source Port / y-axis \Rightarrow Destination Port / z-axis \Rightarrow Source IP). However it improves upon the Spinning Cube of Potential Doom by making the three axis customisable and adding randomness to the position of the points in the cube. The latter uses the same principle of randomness as in the Data Flow visualization: a group of points is undistinguishable from a single point if they have the same position so it is easier to spot the density of a spot in the cube by adding a random offset to the points. **[make nicer]**

4.4 Traffic Volume

The Traffic Volume visualization is inspired by the ‘FlowScan’ graph (www.caida.org/tools/utilities/flowsan/). It is realised as a stacked bar chart which displays the volume of data arriving in each time interval. Each column is segmented into segments with heights proportional to the total number of packets transmitted for each protocol.

¹See www.kismetwireless.net/doomcube/

In addition, the column segments are colour-coded and can be cross-referenced with a protocol key underneath the visualization itself.

To help distinguish between different protocols, colours are selected from a colour palette which provides colours from a qualitative colour scheme. The objective of this visualisation is to allow users to see at a glance if a particular protocol is being exploited in the network. The increase in both column height and colour proportion should draw attention to any protocol which becomes overburdened.

A control panel is provided in the right panel for the purpose of adjusting the number of time intervals displayed at once on the x-axis. The y-axis scales automatically to fit the relevant data.

4.5 Heat Map

Heatmap visualization uses a hexagonal map to display a set of communicating nodes. The activity of a specific machine in the network is represented by the color of the corresponding hexagon. As the packets arrive nodes representing machines that send them "heat up" and as the time passes "heat down". Hexagonal grid is being filled in a compact manner - node by node.

This visualization helps a user find the machines most active in the network. It gives a quick overview of the number of the clients and how many of them are currently active. Presenting most commonly used protocol of every client, it lets users do some basic filtration choices easily. It provides a method of finding out about hotspots in the network.

4.6 Activity Groups

Acting on the hexagonal grid, this visualization puts to use two important visual factors: size and proximity.

As the traffic in the network arises the Activity visualization groups together all the machines sending packets to a one specific other device. As more and more machines communicate with a specific address, more and more nodes appear around the hexagon representing the that destination.

Color of the communicating nodes represents how much data they send. More commonly used destinations will then have more orange and red nodes, and destinations used by a lot of machines will have more nodes around them. Therefore, both a size of the segment and its color indicate important information.

The placement of the nodes is automatic and it always assures that there is enough space around a center and that no two nodes overlap.

This visualization helps detecting the machines that perform a server role in the network, as well as it allows a user to quickly guess what type of service is the device providing.

Each node specifies the most commonly used protocol, thus revealing the possible role of the server node.

5 Discussion

The main determinants of traffic in a network are the network's topology, the flow of the traffic, and the type of the traffic. The graphical visualizations presented here ensure that the available data can be interpreted from all three of these perspectives. Moreover, the underlying idea of our design is that better insight into the activity in the network can be gained by combining these perspectives.

→: Each visualisation was chosen to fill a specific purpose. The Heat Map visualisation fills the need to gain a very quick impression of overall network load and size. The Activity Groups visualisation provides the same information on a server-by-server basis, providing more data at the expense of simplicity. The Dataflow visualisation provides the ability to see what packets currently coming over the network 'look' like, and the Spinning Cube and Attribute Distribution visualisations allow the user specify which aspects of the packets are shown.

As such, no two visualisations fill the same role, each providing something that is not found in any of the others. However, to get a complete understanding of the network no one visualisation will suffice - it is the combination of these complementary visualisations that makes the application unique.

[combine these two, or decide between them, or rewrite]

5.1 User Workflow

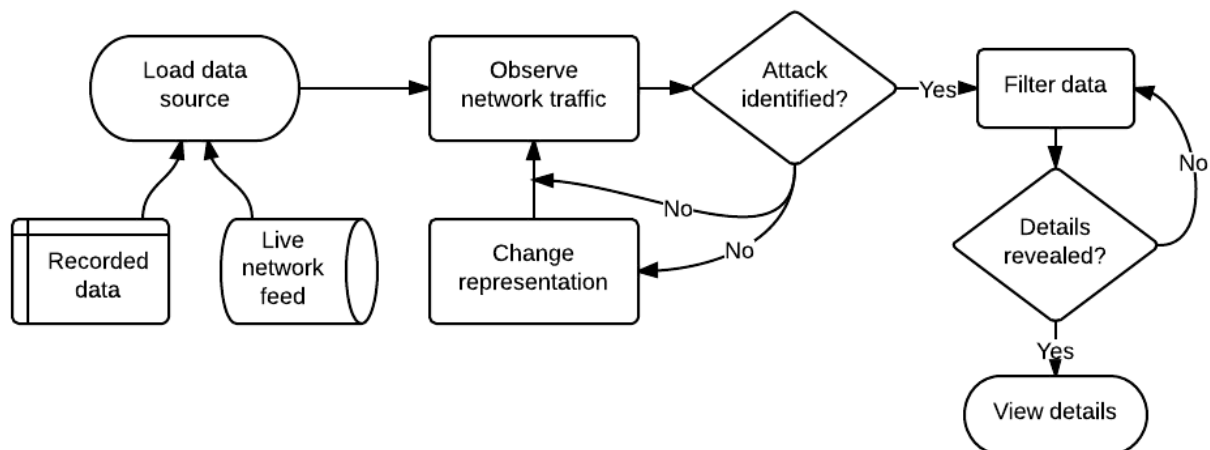
From the perspective of an analyst examining a network, a typical workflow would begin with considering all activity from a broad overview, using any number of visualizations. Once a network attack has been identified, it should be easy for the analyst to "zoom in" to the data of interest, by applying successive filters to the source, and changing the representation to suit the particular attack. Finally, when the exceptional behaviour has been singled out, the analyst may explore details of the intrusion on demand.

This workflow illustrates the visual design framework suggested by Ben Shneiderman, University of Maryland, in 1996;

Overview first, zoom and filter, then details-on-demand.

To assist with the process of singling out suspicious activity, certain visualizations allow filtering by mouse interaction. Where appropriate, the analyst may click or drag their mouse cursor over sections of the active visualization to create a new filter, or switch to a different representation of the selected data. Navigation between complementing visualizations in this fashion is highly intuitive and allows the analyst to tweak the active filters quickly and effectively.

The following is a graphical representation of the typical workflow expected of an analyst studying a network using NetVis.



5.2 Advantages

Existing network visualization solutions are typically single-purpose programs which present given data in one specific way (see for instance the catalogue in applied2004). The analyst is expected to choose in advance which kind of visualization will provide the most insight, and is then limited to what the specific application is able to show. Existing applications also require different input file formats which makes it complicated to use them simultaneously (for a discussion of these problems see).

The program presented here takes a single data input stream and simultaneously visualizes it in different presentations. Since these visualizations are built upon a common base, they are able to complement and inform each other. This improves the user's ability to investigate anomalies.

Network visualizations have to process and display a large amount of data. Often, this can cause a visualization to be cluttered and lose its effectiveness. An obvious solution is to provide filters which allow the analyst to focus on phenomena of interest. The filters we provide apply application-wide and can be defined in an intuitive manner from within the visualizations themselves in a “click-and-zoom” fashion. This is a significant improvement over existing workflows.

The framework developed here can be used both for real-time monitoring, but also for historical analysis, which makes it applicable for a wide variety of use cases. Furthermore, the application can be used as a learning tool: by observing archived data captures, new users can familiarise themselves with the general patterns of network usage, and can also see how suspicious patterns show up in a visualization.

[The following will probably not be used.]

→ The complementary nature of the visualisations means that virtually all information a user could want about the network is available somewhere. It also makes it very easy to

spot a point of interest in one visualisation, then follow it through the others to see more information.

The dynamic filtering makes focusing on a point of interest extremely easy, allowing a user to remove any unnecessary or distracting data. As the filtering is applied retroactively, it allows the user to see past events from different perspectives in order to determine the cause.

As the application works in real-time, the user can see the current state of the network. Hence, they could see a developing network attack and take preventative measures, e.g. blocking the IP of a client attempting a brute-force SSH attack.

5.3 Limitations

The tools in our framework support only the analysis of time series data of packet captures. A multi-layer approach might help inform the user more concretely about the state of the network. In particular, information from intrusion detection systems, routers, and firewalls could be incorporated. The application in its current state does not analyse the content of packets.

5.4 Future Work

Our main avenue of future work is a more powerful data preprocessing system. Our current visualizations focus on giving the user a complete view on the network while allowing interactive filtering. However, if statistical inference is applied to the data, and historical observations are incorporated, this can help make visualizations more effective in pointing out suspicious patterns and anomalies in the data.

Further future developments will also give advanced tools for assisting users in investigating phenomena and better observing specific forms of network attack. For this, a library of applicable data filtering configurations should be available, and an adaptive data selection procedure should be offered.

6 Conclusion

In this paper we have presented..