# Whitepaper
# Harnessing the Power of Metadata for Security

## Overview

As enterprises battle to keep pace with online traffic growth by throttling up network speeds, they are beginning to lose the war on cybercrime. Why? Because security tools are limited in how much traffic they can intelligently process. Further, these same security devices need to take on increasingly sophisticated functions to combat evermore advanced and persistent cyber attacks. As a consequence, many existing security applications will be rendered ineffective in the very near future.

Enter metadata, the new security super power.

Metadata is data about data or put another way, a kind of summary or high-level view of data. By providing security tools with summary takes of the packet data traversing networks, metadata can become a powerful weapon for enterprises looking to separate signals from noise, reduce time to threat detection, and improve overall security efficacy.

## The Challenge

### Too Much Data – too Little Compute

As enterprise networks continue to grow, so too, are the speeds at which the networks' resources are connected and share information. The security appliances, and devices which protect them, however are diminishing by an equal proportion in their ability to consume this additional data and analyze it for security purposes.

Threat complexity has required that security devices take on more complex analytics and this has strained already scarce compute on appliances that could barely match 10Gb speeds let alone 40Gb and 100Gb. Consider this real world data from a leading intrusion prevention system which shows top inspection speeds at 5Gb.[1]
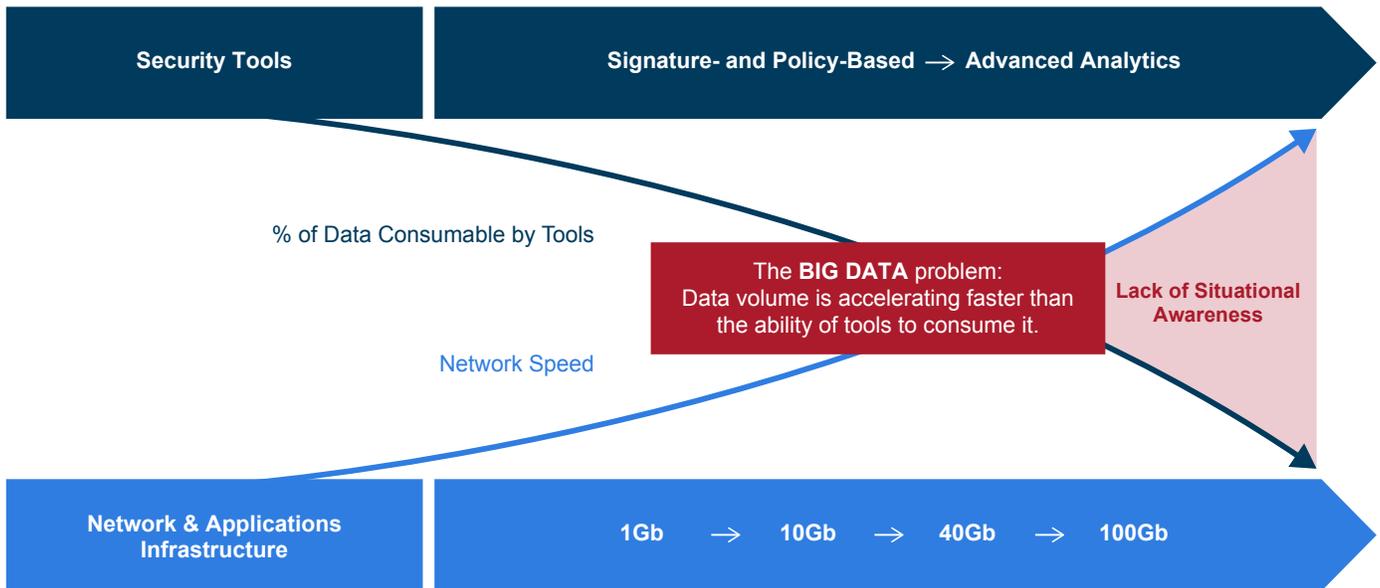


Figure 1: Mapping network speeds to speed of security infrastructure

[1] http://www.cisco.com/c/en/us/products/collateral/security/ips-4500-series-sensors/white_paper_c11-716084.html

# 6.7ns

Time to process a single Ethernet frame on a 100Gbs link with minimum size packets

**Load the packet into memory**

**Extract relevant application information**

**Examine application data and match against signatures**

**Decide what action to take**

**Get ready to load the next packet in**

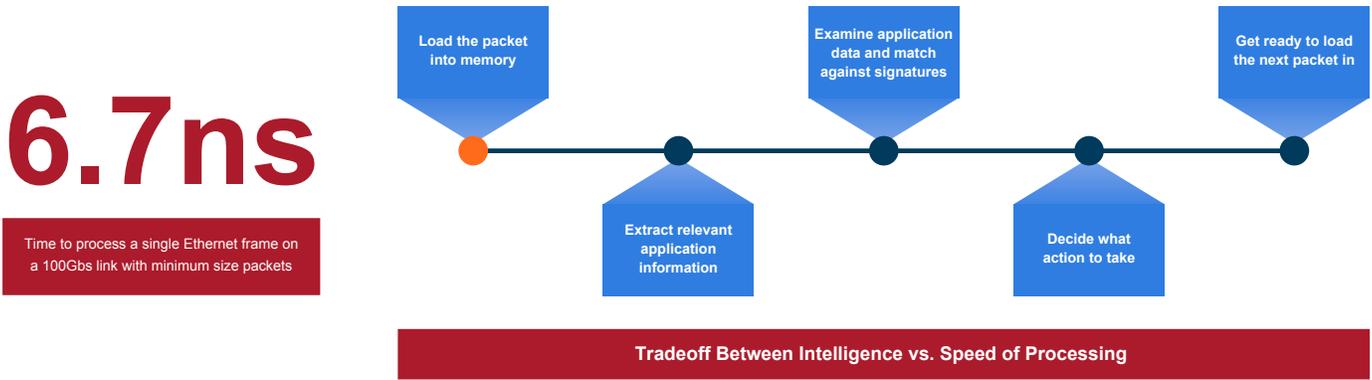**Tradeoff Between Intelligence vs. Speed of Processing**

*Figure 2: Speed of data versus volume of security processing needed in a 100Gb network*

## Prevention is Morphing to Fast Detection and Prediction

The first intrusion prevention systems were signature-based and provided protection for known threats by sequential matching of packets against repositories of attack signatures. Today, security technologies are more advanced and aim at predicting threats before they are known. This zero day detection may use sophisticated modeling based on machine learning, support vector machines, flock analysis, and virtual machine execution, to spot anomalous patterns and bad actors. As network speeds continue to increase, however, the mean time between packets is reduced to 6.7 nanoseconds for 64 byte packets in a 100Gb link. What this means for a security device is that all of the processing associated with examining this packet must happen in this nanosecond time envelope and start anew for the next packet to arrive.

Figure 2 makes it clear that the steps associated with security processing, including extracting the packet, loading it into memory, determining the type of application it is, performing the relevant checks for protocol conformity and anomalous pattern matching, performing signature lookups against known patterns, extracting attachments, computing hashes, etc., are computationally intensive unto themselves. This is not an exhaustive list of all of the security inspection, enforcement, and recording tasks that need to be performed.

To increase processing speed, some security appliances can offload certain operations on user-configurable integrated circuits called field-programmable gate arrays (FPGAs) and custom ASICs. But these approaches too have limits that are easily exceeded when multiple advanced functions are turned on and security effectiveness can suffer.[2]

## Everyone Can be a Hacker

Early malware writers, often motivated by fame, did not make much effort to hide themselves. That is no longer the case today. Threats are stealthier, and attackers go to great lengths to avoid detection because longer dwell times mean attackers can steal more valuable information and reap larger rewards for their efforts. Methods to avoid discovery include disabling antivirus software on infected machines, hiding in encrypted tunnels or non-standard protocols, and erasing their activity on target servers.

In fact, evading detection can be as simple as modifying a known attack to create a zero day or one for which there is no signature. The diagram below demonstrates this. By simply appending a null to the end of the malware, an attacker changes the hash value without affecting the behavior in any form.
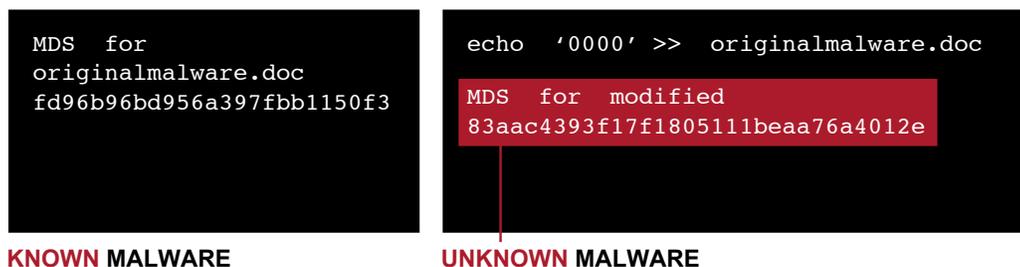
```
MDS for
originalmalware.doc
fd96b96bd956a397fbb1150f3
```

**KNOWN MALWARE**

```
echo '0000' >> originalmalware.doc

MDS for modified
83aac4393f17f1805111beaa76a4012e
```

**UNKNOWN MALWARE**

*Figure 3: Changing the signature of a known malware to an unknown malware[3]*

---

[2]https://www.nsslabs.com/company/news/press-releases/nss-labs-publishes-first-test-of-next-generation-intrusion-prevention-system-products/
[3]https://www.checkpoint.com/resources/2015securityreport/

Figure 3 highlights a frightening new reality in the democratization of hacking. Anyone wishing to do harm needs little skill and few resources to pull off a successful breach. Malware packages and ways to distribute them are readily sold on the black market and cheap to obtain. As a result, the volume and frequency of attacks is growing and the time to successful breach and data theft or exfiltration is getting shorter. This is putting new pressure on security architectures to shorten the time to detection and response.

### How Threats Persist

While modifying malware to get past defenses maybe the core way end-user devices are compromised, the process for successful data theft is a little more involved. Most advanced persistent threats (APTs) follow a Kill Chain (introduced by Lockheed Martin and popularized by Mandiant), which is essentially a six-step sequence that has a specific set of activities associated with each step.

(1) **Reconnaissance.** Attackers target the organization and its employees researching information available on social media and public channels to understand which individuals are strategic and have broad access privileges and how they may be vulnerable to phishing scams or drive by malware.

(2) **Zero day.** Attackers weaponize the information they gather in reconnaissance by succeeding in surreptitiously installing malicious code on the target's device. This code is a toehold or hook into the network of the organization being targeted.

(3) **Back door.** Once on the target device, attackers then reach out on the Internet and make a connection to servers that will serve as the command and control infrastructure. Command and Control (also called C&C or C2) servers can download more malware onto infected devices and help provide additional hacking tools inside the target network.

(4) **Lateral movement.** Using newly downloaded tools from the C&C communication, hackers can move deeper in the network by hoping onto available resources or escalating their privileges through compromised authentication infrastructure.

(5) **Data Gathering.** Attackers begin collecting information about hosts in the network, network topologies, trusted relationships, Windows domain structures, file system locations and much more. During this phase, attackers identify where valuable data is stored.

(6) **Exfiltration**. In this final phase, the valuable data is extracted from the network and sent to deposit locations throughout the web. When this stage completes successfully the network breach has resulted in data loss and theft and is now a security event.

Most attacks follow this kill chain with very minor modifications. Even attacks that look different will follow this pattern. This creates an opportunity to mitigate or stop data theft and loss by identifying attackers and their stage in the kill chain. The steps attackers take at each stage can betray their intentions if organizations are in a position to monitor and analyze their movement inside networks. This sort of identification focuses on security analytics that can apply advanced techniques for pattern matching, a new and rapidly evolving area of cybersecurity called behavioral analytics.

### Building Models of Normal and Bad Behavior

Identifying anomalous behavior like a bad actor moving laterally inside a network requires a couple key steps: 1) building an understanding of what is normal behavior inside networks so that deviations can be flagged and 2) understanding the context for the behavior to accurately identify the deviation as appropriate or malintended.
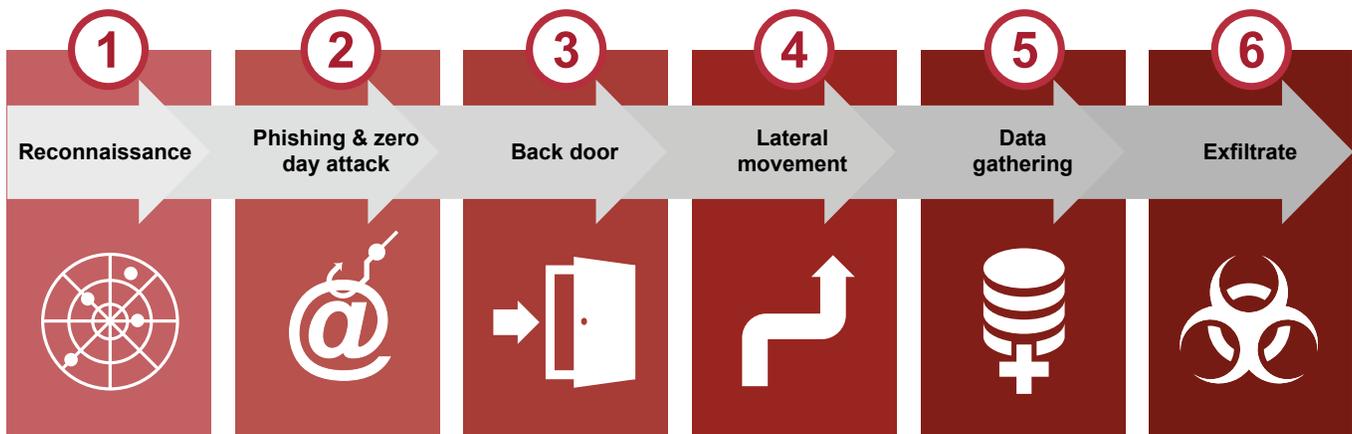


| 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|
| Reconnaissance | Phishing & zero day attack | Back door | Lateral movement | Data gathering | Exfiltrate |

*Figure 4: Six-step sequence of specific activities for an Advanced Persistent Threat*

What is considered normal behavior can be very specific to an organization depending on the type of business it is engaged in and how its network is architected, all of which create somewhat unique traffic patterns. To establish an organization's normal means building a model through observation of data flows and connectivity. Similarly, models of bad or malicious behavior can be built using known data sets like those available through malware databases and threat feeds. These models need to be continuously informed by context for what is normal (ex. traffic spikes during end of month promotions are typical and recurring) and intent (ex. what seems like too many unauthorized access attempts is actually a forgotten password issue not a DDoS attack). Once these models and a system for updating them is established, they are continuously triangulated with respect to one another so anomalous and truly malicious behavior can be identified.

This seems straightforward enough as a plan but it is quite difficult to implement. Establishing context for network use is challenging because there are simply too many sources of information, all belonging to different departments (e.g., groups for endpoints, network routing and switching, applications, virtualized systems, security devices like firewalls, intrusion prevention and data loss prevention systems) with differing access rights and formats (e.g., syslog, json, and text files). And even if it is possible to gather all of this information, there is simply too much of it to store and analyze easily.

Many organizations receive millions (if not billions) of events from devices. To make matters worse, some of these sources of information are not available at all. Auditing functions can be computationally intensive and punishing to security devices whose primary function is to protect, so by default granular logging may be off.
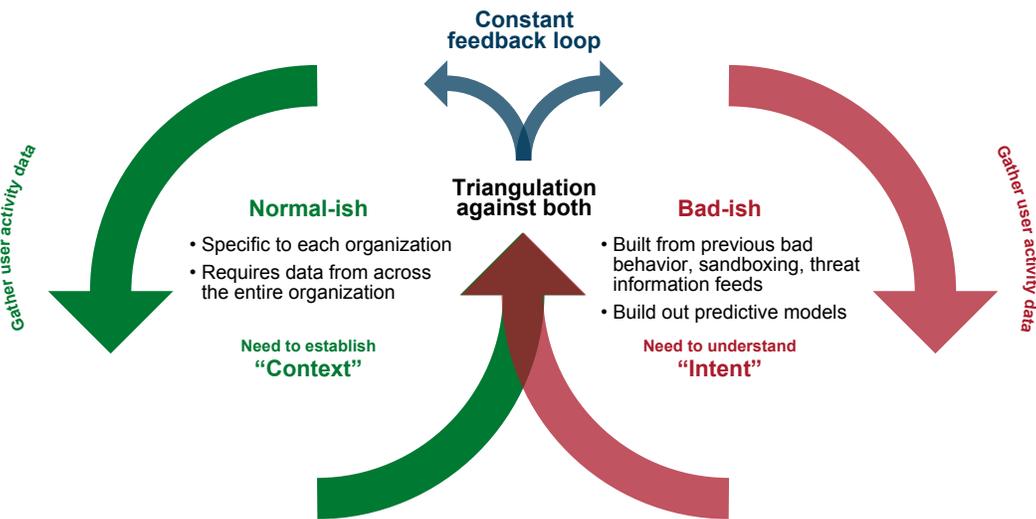


*Figure 5: Determining normal versus bad behavior in an organization*

**Sources**
- Endpoints and servers
- Applications
- Switches, routers
- Network appliances

**Challenges**
- Different departments
- Different access rights
- Different formats
- Agent requirements

**Consequences**
- Massive inefficiencies
- Too much data
- Less control
- Performance impact

**Slows Down Analysis, Slows Down Response, Slows Down the Feedback Cycle**

*Figure 6: Challenges associated with getting relevant security data in organizations*

## The Solution
### The Case for Network Metadata

Deriving context for network use by going to each of the individual source classes for this information is quite difficult to impossible. Fortunately, there is an alternate route to retrieving this information and that is by leveraging the network. The network is the massive web of interconnectivity that spans the physical, virtual, and cloud environments and connects users, devices, and applications. The network is the single best source of traffic and the context by which it is traveling. So by harvesting the information directly from the network, organizations can obtain the context needed to refine behavioral models without the herculean effort required in touching every single network endpoint and source.

Along with traffic or packets, flowing within the network is summary information about that traffic, like where it is going and where it has been. This summary information or metadata can provide valuable clues to lingering threats inside networks. Behavioral and security analytics using metadata gives organizations an approximation of the location of hot spots or areas of suspected threat activity. Rather than searching the entire network, security analysts can focus on the identified trouble spots and conduct a more thorough investigation by using traffic or packet analysis. Consider the example of a DNS request made by a laptop to a suspicious server. Investigators can examine all connectivity from that device to ascertain if the endpoint is infected, whether it has forwarded malicious content and what other devices might be implicated.

Building models of normal and bad behavior require refinement by context. This context is easily retrieved from metadata, which is readily available inside networks and does not require touching individual endpoints. Security and behavioral analytics' approaches that use metadata as a first step in the research can create approximations of where threats reside making the in depth investigation more focused and the time to threat discovery shorter.

## Gigamon's Metadata Engine
### Overview

Gigamon is a company that pioneered the area of network visibility. Gigamon's GigaSECURE® is a Security Delivery Platform (SDP) that provides pervasive reach across the infrastructure, spanning cloud, on-premise data centers/locations, and remote sites in order to bring the network to the security and performance management devices that require access to it in order to conduct their functions. Security appliances specifically, simply connect into the GigaSECURE platform to receive a high-fidelity stream of relevant traffic from anywhere in the network infrastructure.

One key pillar of the GigaSECURE Security Delivery Platform is the ability to generate summary information about packets from network traffic and put it in the NetFlow format. This functionality provides security and behavioral analytics products with valuable unsampled information about traffic without impacting the performance of infrastructure. Recently, Gigamon transformed NetFlow generation into a full-blown Metadata Engine that sits inside the GigaSECURE SDP and serves as the single source of network truth for all kinds of information about applications, users, and devices, and sends relevant information to the security tools connected to the SDP.
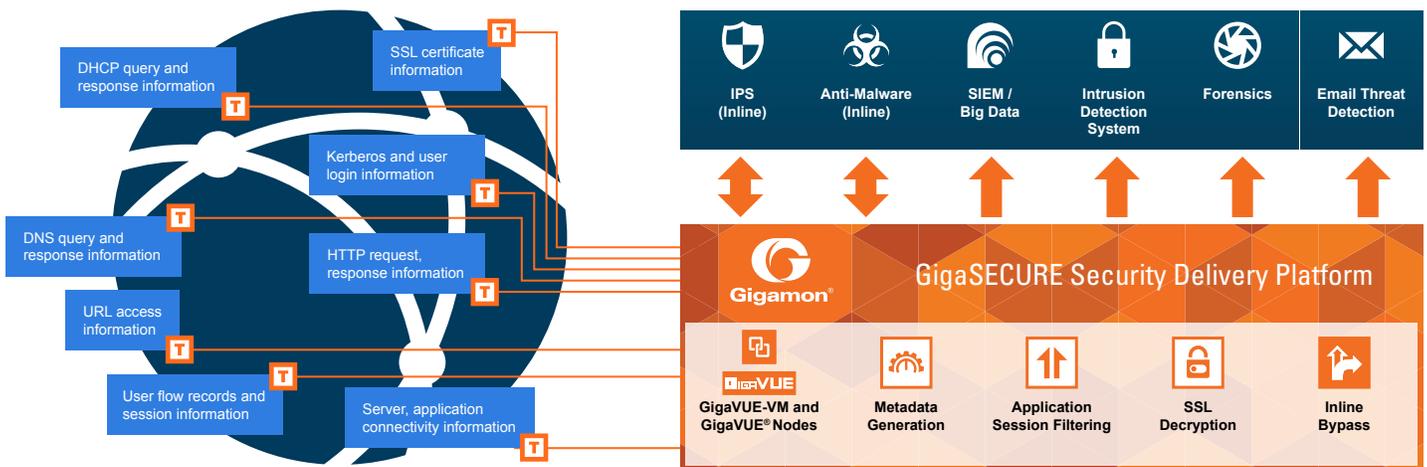


*Figure 7: The GigaSECURE Metadata Engine*

The GigaSECURE Metadata Engine provides powerful infrastructure by which to collect network metadata elements (inclusive of NetFlow) that are critical for security operations, It offers a single source of valuable data for both network and security operations.

## Examples of Gigamon-provided Metadata and Uses

### HTTP URLs and Response Codes

HTTP URLs and response codes are a valuable source of threat information and Gigamon can generate both.

Web security gateways for instance can look at URLs provided by Gigamon's GigaSECURE platform to determine a host of security vulnerabilities (ex. whether the URL is a match on blacklists, is a known C&C server, or is consistent with SQL injection). If security analysis identifies the presence of a malicious site lookup in the network traffic, Gigamon users can search for that location everywhere in the network metadata to see which endpoints maybe have accessed it.

HTTP response codes are similarly important for determining anomalous behavior. Response codes are divided into 5 categories:

- 100-199 = informational
- 200-299 = success related
- 300-399 = redirection
- 400-499 = client requests
- 500- 599 = server related

For the casual Web user, an encounter with a 404 error code means trying a couple more times and moving on after unsuccessfully accessing the page. However, amass all of these codes over time and, through the power of analytics, a picture of suspicious and nefarious activity can come into view.

For instance, an analysis of 2xx codes that denote successful access to resources requiring special authorization might surface an atypical or unauthorized user getting through. A flurry of unusual HTTP 2xx codes in a short time span from servers way beyond what is considered normal for the network could mean an attacker has found a way to send malicious requests that are being answered by servers. 3xx codes need to be looked at closely for redirections that are anomalous and may lead to sensitive URLs. Too many redirects could indicate a compromise of internal servers. Likewise, too many 4xx codes in a short time span could signal that an infected machine is searching to make contact with a command-and-control server.

### Domain Name Service (DNS) Discovery

DNS information can also provide extremely valuable security clues. Take an example of an infected device trying to reach its command-and-control center (e.g., www.evil.com). The first thing it does is initiate a DNS request to resolve the IP address for the domain www.evil.com. Providing the metadata for the DNS request for this domain from the original DNS request and response (also called authoritative DNS exchange) can help a security analytics device put the pieces of this multipart request together and quickly flag that this action is consistent with C&C communications and a potential threat.

Another field that is of interest in the DNS request is "canonical name". It gives a list of other domain names that are being served by the same IP address. In this example, another domain (e.g., www.veryevil.com) could be using the same server and the security analysis will surface this destination as a source of threat as well.

Further, the time to live (TTL) value returned in the DNS server is also useful. Normal servers typically have a time to live value corresponding to one day. However, most malicious servers have very low TTL values because as they are discovered and blocked they change their DNS associations.
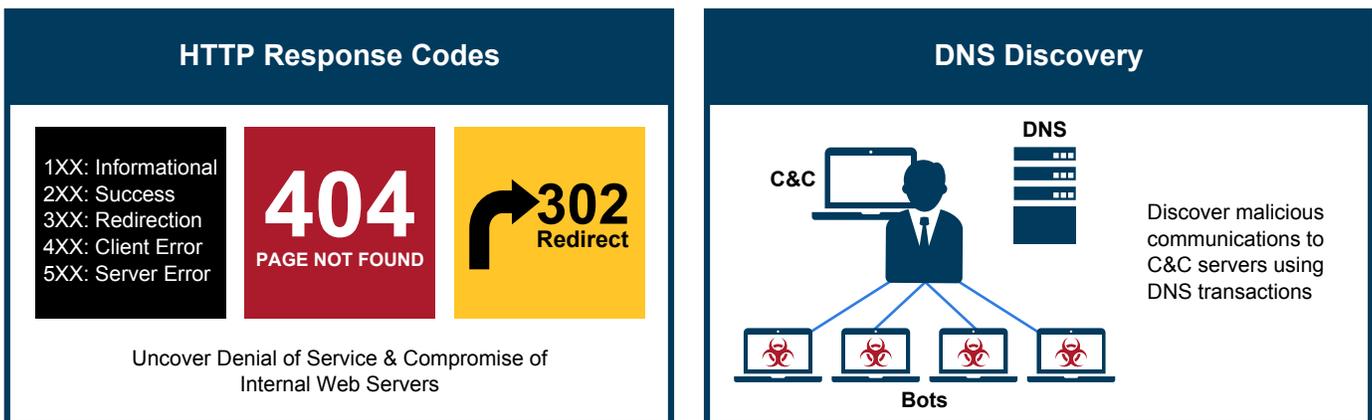


*Figure 8: Metadata extensions with HTTP response codes and DNS discovery*

**Certificate Anomalies**

Secure Socket Layer (SSL) certificates are another source of valuable security information. Certificates are typically used by servers to authenticate themselves to clients prior to starting the encrypted SSL session (often shown as a lock in most browsers). All certificates have an Issuer Field that lists the name of the Certificate Authority (CA) granting the certificate. There are only a few reputable certificate authorities that sign the majority of these certificates. Anything signed by a CA that is unfamiliar should raise a concern and warrant investigation.

Similarly, self-signed certificates or certificates that have expired should be treated with suspicion. The subject field in the certificate indicates the entity to whom the certificate has been issued. A website that offers a certificate with a subject field that does not match the name of that same website should prompt concern with a security stakeholder.

**The Metadata Engine and Application Session Filtering**

Application Session Filtering (ASF) is another key functional pillar of the GigaSECURE Security Delivery Platform. ASF provides a powerful filtering engine that identifies applications based on signatures or patterns that appear across any part of the packet payload. ASF can extract entire sessions corresponding to a specific application.

The combination of network metadata with Application Session Filtering creates a powerful method for incident analysis and control. While network metadata can be used to provide valuable clues to detect anomalous behavior, ASF can be used for the deeper analysis or drill down into suspicious traffic streams. As an example, if metadata is used to uncover a malicious URL corresponding to a command-and-control server, ASF can be dynamically turned on to look for the command-and-control communication corresponding to that server across all network Web traffic. This search will effectively identify infected endpoints that may be communicating or connecting to the malicious server, helping to stem or stop data exfiltration.

## Conclusion

Security will increasingly rely on the modeling of good and bad behavior made accurate through tuning for context and intent. Security analytics triangulating these models to uncover anomalous behavior and identify data breaches are already gaining broad adoption in the form of SIEM, NetFlow and user behavioral analytics solutions. As traffic speeds grow, ensuring that the models and analytics approaches are informed by not only Big Data, but the right data, will be key to making threat detection pragmatic and effective in modern high-bandwidth networks.

Gigamon's GigaSECURE Security Delivery Platform provides metadata that is rich in context and helps make security analytics faster at approximating the location of breaches as a first step. Then the same platform enables more in-depth security analysis of traffic that is focused on the areas of concern, ultimately resulting in faster time to detection, response, and mitigation.

4068-01 06/16