

Adaptive Alert Throttling for Intrusion Detection Systems

Gianni Tedesco (Corresponding Author)
43 Runswick Terrace
Bradford
BD5 8LR
United Kingdom
gianni@scaramanga.co.uk

Uwe Aickelin
School of Computer Science
University of Nottingham
Nottingham
NG8 1BB
United Kingdom
+44 115 9514215
uwe.aickelin@nottingham.ac.uk

ABSTRACT

Each time an intrusion detection system raises an alert it must make some attempt to communicate the information to an operator. This communication channel can easily become the target of a denial of service attack because, like all communication channels, it has a fixed capacity. If this channel can become overwhelmed with bogus data, an attacker can quickly achieve complete neutralisation of intrusion detection capability. Although these types of attack are very hard to stop completely, our aim is to present techniques that improve alert throughput and capacity to such an extent that the resources required to successfully mount the attack become prohibitive.

Keywords

Intrusion Detection System, Denial of Service Attacks, Token Bucket Filter

1. INTRODUCTION

As global awareness of information security issues has increased, so has the proliferation of intrusion detection technology. It is widely agreed that network intrusion detection systems (NIDSs or IDSs) are quickly becoming a crucial part of the Internet security infrastructure. In March 2001, there was a media furore [3] when the FBI Internet crime division issued a warning concerning the then unreleased Stick [5] program which “essentially disarms intrusion detection systems.” With attackers now focusing their energy on network intrusion detection systems, the importance of protecting them becomes ever greater.

IDSs come in many different forms [7]. Our work is based on network IDSs (NIDSs), which detect attacks by directly analysing network traffic in real-time. Techniques used by NIDSs still have a lot of room to improve. Ning [2001] identifies a number of problems associated with current NIDSs.

Examples of current NIDS research include methods of identifying new or varied attacks [12], ‘honeypots’ [1] and anomaly detection [10]. Our area of research focusses on preventing denial of service attacks.

There are, in fact, numerous possible types of denial of service attack against an IDS [13], but we will focus on one particular attack type which is fundamental to all systems: (automated) alert flooding. We will also focus primarily on NIDSs using the signature matching model described below. However, most of the techniques presented will be generally applicable to any type of IDS.

The pattern matching [19] model is currently the most commonly used methodology for detecting intrusion attempts. In this model the NIDS is configured with a database of known attack patterns (also called signatures). An example of a signature is shown in Figure 1. This signature alerts on traffic generated by the well-known “BackOrifice” trojan horse program and detects any incoming packets destined to user datagram protocol (UDP) port 31337, containing a specific sequence of bytes anywhere within its payload.

Figure 1: A Sample Rule as used e.g. by SNORT

```
alert udp $EXTERNAL_NET any -> $HOME_NET 31337
(msg:"BACKDOOR BackOrifice access"; content:
"|ce63 d1d2 16e7 13cf 39a5 a586|";)
```

Alert flooding attacks are achieved by transmitting packets that simulate intrusion attempts and which the IDS will recognise as true attacks. Taking the example signature in Figure 1, an attacker must craft a UDP packet, set the destination port to 31337, include the sequence of bytes given in the signature and flood the target network with these packets.

The possible ramifications of this type of attack against an IDS are threefold:

1. Sensor storage becomes full, preventing further logging.
2. Sensor exceeds maximum alert throughput, causing alerts to be lost, or the sensor to cease functioning.

3. The analyst becomes deluged with false information and becomes unable to distinguish real attacks from the false ones.

The alert flooding technique has been automated, and hence popularised, by tools such as Stick [5] and Snot [17] which read in signatures directly from the freely available Snort [15] IDS. Each packet sent could also have crucial fields such as source and destination address modulated by adding random data into them. This random noise makes it difficult to block the attack using a simple packet filter or firewall.

Alert floods can also be exacerbated by the poor alerting performance of IDS systems in general. A quick examination of the Snort [15] system reveals that, in its preferred output mode (called “unified”), Snort flushes its buffers needlessly in at least two places. This causes a reduction in the effectiveness of the buffering and on UNIX like systems results in added system call overhead for every logged alert.

Performance in this area can be understandably overlooked by the IDS system designer. After all, good engineering practise tells us to optimise for the common case, and, in the world of intrusion detection, an alert is not usually the common case. In fact, on a high-speed network it should be a very rare event indeed.

The next section of this paper will show that previous approaches have focused on removing shortcuts to perform the attack in order to raise the cost to the attacker. In section 3 we define the problem more carefully and examine the flaws in existing solutions. We then further the line of thought, suggesting how an adaptive throttling technique can increase the cost to the attacker as the magnitude of the attack increases. We will also provide some experimental data, and some calculations of the effectiveness of our technique. We finish by presenting a summary and some concluding remarks.

2. CURRENT STATE OF THE ART

The Snort [15] team addressed the problems of wide spread proliferation of automated alert flooding tools like Stick [5] and Snot [17] in their 1.8 release. Their solution was to implement a Transmission Control Protocol (TCP) state tracking system which they called “stream4”.

TCP is the transport protocol of the majority of Internet traffic. It is a connection oriented two-way stream protocol (similar to a telephone call) and has a connection initialisation procedure known as the “three-way handshake”. A TCP message (called a segment) cannot exist outside an established connection in much the same way that a telephone conversation cannot exist without first dialling a number and waiting for the other side to pick up the receiver.

By keeping track of TCP connection states, stream4 is able to ignore any segments which are not part of such a conversation. In order to make the IDS raise an alert the attacker is now forced to transmit at least three segments, rather than just one. More importantly, because the three-way handshake requires two hosts to be communicating, the external attacker must find a host on the monitored network

willing to participate. This might be prevented by a firewall blocking connections.

Currently most systems keep track of TCP states. This is mainly to protect against desynchronisation attacks as described by Ptacek and Newsham [13], but there is also the additional benefit of making sure that there is no such shortcut in carrying out an alert flooding attack. Further to performing TCP state tracking, it is also possible to track any application layer state, enabling us to remove shortcuts even for protocols running over stateless transports such as UDP.

Looking at it pragmatically, we can assert that with the stateless implementation the attacker had a shortcut to make the attack cheaper, equivalent to simply picking up a telephone receiver and starting to talk. Naturally, the solution is to remove that shortcut. This forces the attacker to find a valid number, dial it, and wait for an answer.

While this approach is a positive step, it cannot cover all cases: For example, some signatures must ignore state information as some exploits can exist as a single packet (i.e. statelessly); or because in other cases, they work over inherently stateless protocols. In the following sections, we will go on to show ways to reduce the alerting throughput during an alert flood attack. In particular, we will present a method that works regardless of signature configuration.

3. ADAPTIVE DATA REDUCTION

Our idea is to create a technique for reducing alert throughput during an alert flood attack by detecting that a flood is occurring and then adapting to the threat by making the IDS more terse, even if that means discarding or “dropping” alerts. We show that the technique scales up such that it effectively nullifies the threat of an alert flood attack. A similar approach to this has been successfully deployed in order to drastically slow the spread of worms and viruses across the Internet [24]. We then show how some simple compression techniques can be used to ensure that we can still log important information about the attack.

In section 1 and 2 we have described that the problem is fundamentally that of resource exhaustion. For example our human IDS operator is a resource, and one which cannot cope with having to examine many thousands of bogus alerts at the rate at which a sustained attack can produce them.

There are two approaches to solving this type of problem: one is to increase the amount of resources you have, the other is to reduce the amount of resources required. While it is conceivable that one could scale the sensor hardware to be fully able to cope with alert floods at a given rate for a given length of time it seems rather more complex to scale the human operator.

Unfortunately, the processing capacity of human computer operators scales linearly at absolute best [4]. However, for complex analytical tasks requiring coordination and communication, adding more people is so inefficient that it can even make a task take longer. Besides which, even if humans could scale to the task linearly it would seem an incredibly expensive proposition for an organisation to multiply the

number of operators even by a factor of ten. It would seem that the more beneficial approach would be that of data reduction.

3.1 Token Bucket Filter

Perhaps the simplest way to reduce data output while maintaining the same intrusion detection capability is to make minor modifications to the signatures to make sure that the IDS is as terse as possible. Such modifications are often used to reduce the number of false positive alerts generated. In fact generally speaking, signatures are usually a subtle compromise between allowing false negative and false positive alerts.

One way to make the IDS less verbose is to fine-tune signatures to examine only those packets destined for the relevant hosts. Let us consider BIND, DNS server software infamous for its security vulnerabilities. In this situation, the signatures may be modified to only look for BIND exploits if the destination address on the packet matches a pre-defined list of DNS servers. Of course, the operator may actually be interested to know that someone is attempting a BIND exploit on a workstation or a web server. That is to say, this approach tips the false alarm compromise towards the false negative side. The interested reader is pointed to 'Target-based IDS', a new approach towards this direction [18].

Perhaps a more beneficial approach would be adaptive. In the case where there is no need to conserve alert throughput we log all suspicious activity. Once a certain rate threshold has been met, indicating an alert flood attack, we could switch to a more terse behaviour or even drop alerts exceeding the threshold all together. This can be achieved very efficiently by using a token bucket filter [23].

A token bucket filter is an algorithm for controlling the rate of flow of data. Token bucket filters have traditionally been used in a number of applications where rate limiting has been needed. Some good examples are:

1. Network bandwidth management systems [8].
2. Flood protection in network chat / text conferencing systems such as Internet Relay Chat.
3. Flow control in network transport protocols [14].
4. Flood protection for programs that log externally generated events such as UNIX syslog.

A token bucket filter has two parameters, bucket size, and token rate [23]. Tokens are generated at the token rate and stored in a buffer called the "bucket". If the bucket becomes full, the extra tokens are just discarded. Each alert that arrives must have a token to pass through the filter. Any alert that does not have a token is called "over-limit" and does not pass the filter. If the alert rate is less than the token-rate then credit is allowed to accumulate in the bucket. This stored credit allows for the alert-rate to temporarily exceed the token rate (or "burst").

The token bucket filter could be applied per signature, per attack type, globally, or even in complex hierarchies as in HTB3 [6].

One NIDS that supports token bucket filtering is Firestorm [20] which is a high performance NIDS that is similar to Snort but with additional benefits such as protocol anomaly detection, intelligent TCP stream reassembly and full application layer decodes. In Firestorm, it is possible to set per signature and per "generator" (signature type) thresholds. The filtering is readily configurable as an extension to the Snort signature format which it supports natively. Firestorm completely drops all over-limit alerts.

We can perform a simple test with the Firestorm [20] system running off-line against a tcpdump [11] capture file containing a stick attack of one million packets at a rate of around 7,500 packets per second. This is a realistic packet rate for a high bandwidth network and similar rates have been observed in the wild, for example a Shmoo Group defcon capture [2] shows a rate of around 7,343 packets per second for an ICMP flood. Our stick generated data can be found on the Firestorm website [21].

The Firestorm system has been chosen as it already supports token bucket filtering and supports the same signatures and intrusion detection techniques as the popular Snort system. This makes it roughly representative of NIDS systems that are currently deployed. The software is also, with available source code that allows us to readily add instrumentation and measurement code.

In the both tests, we have a full signature database loaded containing around 1,600 signatures, with the network data read directly from the hard disk. The test machine was a 2.6GHz AMD Athlon XP running Linux 2.6. The results shown are an average of three iterations for both runs to factor out any random fluctuations such as disk seek latency.

The first run (#1) is a control run. The second run (#2) is identical except for the addition of token bucket filtering. There is a filter for each rule which is set to 1 alert per second and a burst of 10 alerts. Each protocol (IP, TCP, UDP and ICMP) is set with a total rate of no more than 10 alerts per second with a burst of 50 alerts. These parameters are rather arbitrary but based on our experience are reasonable to differentiate a real alert flood from a small sequence of genuine attacks.

As we can see in Table 1, the amount of data logged was reduced by several orders of magnitude and the run time decreased disproportionately to the CPU time. While the run time was reduced by around 40%, the CPU time only reduced by around 20%. This indicates that with the token bucket filter enabled, the Firestorm process is not wasting as much time waiting for I/O completion.

From these results it is clear that we can effectively boost performance and capacity, allowing the sensor to carry on working during an alert flood rather than becoming overwhelmed and possibly exhausting the storage on the sensor. Even if the attack contained twice as many packets in the same space of time, it would not double the amount of data logged as the token rate is fixed.

3.2 Alert Compression

Table 1: Token Bucket Filter Alert Throughput

#	Data Size (Kilo Bytes)	Alerts	CPU Time (seconds)	Elapsed Time (seconds)
1	130,001	523,256	5.568	7.151
2	942.367	4381	4.490	4.501

Dropping alerts entirely may seem a less than ideal way of dealing with an alert flooding situation. The IDS operator may be uncomfortable with not having all the evidence with which to make judgements. However, if we assume that alert flood attacks are mainly repetitive, we may exploit this similarity to achieve effective compression of over-limit alerts instead of dropping them.

If we assume that the attack repeatedly transmits the same packet then we can use a “run length encoding” (RLE) to represent it. RLE is a simple compression technique which replaces recurring sequences of symbols (called runs) with a single symbol and a run count N . To decompress, one simply copies the symbol into the output stream N times. This is an approach familiar to UNIX users who have ever tried to flood the syslog program and seen its “last message repeated N times” warning.

To implement RLE compression in our case all that is required is to store a reference to the first over-limit alert against the token bucket data structure and increment a counter for all further over-limit alerts. When there is enough credit in the token bucket to permit new alerts, we flush out the alert and the counter to permanent storage.

The operator can then see that alert A was repeated X times. The only data lost is the timing of the $(X - 1)$ alerts which constitute the run. However, the time can be encoded along with the alert as delta values. For example if the IDS timestamps are kept in millisecond resolution we could use a 16 bit unsigned integer to represent the time deltas. That would allow us to keep totally accurate timestamps provided that no two alerts in the run are more than 65 seconds apart.

An advantage of combining delta compression with RLE is the ability to present the operator with one composite alert for each run while still providing the ability to see the complete unfiltered data or “drill down” if that is required. This can be implemented by only displaying the first alert in each run and providing the operator with an option to see the whole run, which could be decompressed on the fly.

Using our previous experiment as a basis, we can calculate the storage requirements for this type of compression. The Firestorm alert log format has enough reserved space for an RLE run counter which means that RLE compression has exactly the same storage requirements as the previous “token bucket only” result. We also provide a comparison with popular compression programs gzip [9] and bzip2 [16]. In this case, we simply ran the gzip/bzip2 commands with default parameters and took the resulting file size.

As Table 2 shows, the RLE combined with delta compression provides a very impressive compression ratio (almost an order of magnitude greater than existing techniques) and does this without loss of fidelity when all the alerts are identical.

It is also worth noting that the gzip and bzip2 compression took around 11 and 62 seconds to complete respectively. This assures us that our technique does not present any significant computational overhead in the way that general purpose compressors do.

To extend this technique we can further generalise our assumption that alerts are identical. In reality, an attacker can randomly modulate fields in the packet structures, notably source and destination addresses (cf. Stick [5].) In this instance, we can use the delta compression technique for those extra fields or indeed for *all* fields in our alert as opposed to just the timestamp.

If faced with modulation within the actual packet payload one can choose a trade-off between losing all payload data for alerts in the run or saving all data, possibly by using an algorithm like rsync [22] as a generic delta compressor.

A real world system could be configured and fine-tuned allowing the operator to make the trade-off between efficiency and data fidelity based on operational parameters such as hardware and network capacity and the perceived threat of alert flooding. The operator may decide, for example, that the IP type-of-service (TOS) field provides so little information that it is not worth logging.

4. SUMMARY AND CONCLUSIONS

Alert flooding is a problem that will probably always exist with intrusion detection systems and one that cannot be eliminated entirely. However, we have shown that it is possible to drastically reduce the effects by recognising an attack and responding proportionately.

We detect the alert flooding attack using a token bucket filter and react by either dropping or compressing the alerts which make up the flood. We showed that delta compression of relevant fields can provide massive compression ratios in this situation, meaning that we can minimise the amount of useful data that is lost in reacting to the threat.

As an added benefit of the techniques presented, we are also able to group the many related alerts that make up the attack together, greatly simplifying analysis.

However, more investigation is needed to produce optimal token bucket filter parameters as the performance of the system rests on this variable. An ideal scheme would probably be hierarchical, providing greater flexibility in configuration.

A method for efficiently compressing full packet payload without loss of fidelity would also be useful. This should not present a serious problem as the rsync [22] algorithm provides a good delta compression solution. There is also a substantial body of knowledge on generic data compression algorithms. Gzip [9] would appear to be a good choice.

Table 2: Comparison of Compression Techniques

Algorithm	Data Size (Kilo Bytes)	Compression Ratio
Uncompressed	130,001	1
RLE Only	924.367	137.952
RLE With Timestamp Delta	2970.199	43.768
Gzip	24263.5	5.357
Bzip2	16634.16	7.815

Acknowledgements

We would like to thank Louisa Parry, Matthew Hall and John Leach for input with writing the article.

5. REFERENCES

- [1] *Hogwash*. <http://hogwash.sourceforge.net/>.
- [2] *CCTF Defcon Data*. <http://www.shmoo.com/cctf/>, 2001.
- [3] *ZDNet UK News*. 2001.
- [4] F. P. Brooks. *The Mythical Man-Month*. Addison Wesley, 1995.
- [5] G. Coretez. *Fun with Packets: Designing a Stick*. 2002.
- [6] M. Devera. *Hierarchical token bucket theory*. <http://luxik.cdi.cz/devik/qos/htb/manual/theory.htm>, 2002.
- [7] Escamilla. *Intrusion Detection*. Wiley, 1998.
- [8] R. R. G. Woodruff and P. Richards. *A congestion control framework for high-speed integrated packetized transport*. 1988.
- [9] J.-L. Gailly and M. Adler. *The gzip compression algorithm*. <http://www.gzip.org/>, 2003.
- [10] J. Hoagland and S. Staniford. *Statistical Packet Anomaly Detection*. 2001.
- [11] L. V. Jacobson, Craig, and S. McCanne.
- [12] D. R. Ning and Y. Cui. *Correlating Alerts Using Prerequisites of Intrusions*. 2001.
- [13] T. H. Ptacek and N. N. Newsham. *Insertion, Evasion and Denial of Service: Eluding Network Intrusion Detection*. January 1998.
- [14] M. K. R. Wade and P. Dew. *Study of a Transport Protocol Employing Bottleneck Probing and Token Bucket Flow Control*. 2000.
- [15] M. Roesch. *Snort - Lightweight Intrusion Detection for Networks*. <http://www.snort.org/>, 1999.
- [16] J. Seward. *bzip2 compressor*. <http://sources.redhat.com/bzip2/>, 2002.
- [17] Sniph. *Snort*. <http://www.stolenshoes.net/sniph/index.html>, 2001.
- [18] J. Snyder. *Taking Aim*. http://infosecuritymag.techtarget.com/ss/0,295796,sid6_iss306_art540,00.html, January 2004.
- [19] C. Systems. *The Science of Intrusion Detection System Attack Identification*. <http://www.cisco.com/warp/public/cc/pd/sqsw/sqidsz/prodlit/ids2002>.
- [20] G. Tedesco. *Firestorm Network Intrusion Detection System*. 2004.
- [21] G. Tedesco. *Stick Dataset*. 2004.
- [22] A. Tridgell. *Efficient Algorithms for Sorting and Synchronization*. April 2000.
- [23] J. Turner. *New directions in communications (or which way to the information age?)*, volume 24. 1986.
- [24] J. Twycross and M. M. Williamson. *Implementing and testing a virus throttle*. 2003.

Biographies

Gianni Tedesco has worked in the security industry for 4 years, primarily in the West Yorkshire area. He is the author of the Firestorm Network Intrusion Detection System which is a high-performance signature based IDS. He has also worked on other projects within the security arena such as the Linux Netfilter firewalling framework.

Uwe Aickelin is a lecturer in Computer Science at the University of Nottingham. Previously, he worked as a lecturer in Computer Science at the University of Bradford and in Mathematical Sciences at the University of the West of England. His main research interests are Artificial Immune Systems, Intrusion Detection, Scheduling and Evolutionary Computation. Uwe holds a PhD and a Masters degree from the University of Wales, Swansea and an undergraduate degree from the University of Mannheim.