

network security

ISSN 1353-4858 May 2020

www.networksecuritynewsletter.com

Featured in this issue: Optimising storage processes to reduce the risk of ransomware

t's evident that digital crime is just as lucrative for criminals as it is destructive for businesses. What's more, hackers are capitalising on advances in technology to launch more-advanced and destructive attacks that are even harder to detect.

The simple universal rules of this game are that all organisations face the danger of hackers destroying their data. The pressure is on for organisations to not only focus on the prevention, but also on the protection of the data itself, including its restoration. Businesses need a back-up solution that makes critical data easily recoverable in the event of a loss. In addition to reviewing security processes, organisations should re-examine their storage infrastructure to optimise their protection against attacks. Florian Malecki of StorageCraft provides a series of steps for businesses to take to ensure that they have the right elements in place to protect their data.

Full story on page 6...

Hybrid intrusion detection system using machine learning

Cloud-based architectures have reduced IT barriers and provided new capabilities of dynamic provisioning, monitoring and managing resources by providing immediate access to resources, enabling the easy scaling up of services.

However, sharing the same pool when requesting services involves the risk of data breaches, account compromises, injection vulnerabilities and distributed denial of service (DDoS) attacks. As a result, many customers rank cloud security as a major challenge that threatens their work and reduces their trust in cloud service providers. Amar Meryem and Bouabid EL Ouahidi propose an architecture that eradicates malicious behaviours by detecting known attacks using log files; blocks suspicious behaviours in real time; secures sensitive data; and establishes better adaptations of security measures by dynamically updating security rules.

Full story on page 8...

Chinese hackers attacking Covid-19 researchers, US warns

The US Department of Homeland Security's Cyber security and Infrastructure Security Agency (CISA) and the FBI have issued a joint statement accusing hackers based in the People's Republic of China (PRC) of attempting to steal research relating to Covid-19.

"These actors have been observed attempting to identify and illicitly obtain valuable intellectual property (IP) and public health data related to vaccines, treatments and testing from networks and personnel affiliated with Covid-19related research," the notice says. "The potential theft of this information jeopardises the delivery of secure, effective and efficient treatment options."

Continued on page 2...

Contents

NEWS	
Chinese hackers attacking Covid-19 researchers, US warns	1
InfinityBlack market taken down	2
New North Korean threats	3

FEATURES

Optimising storage processes to reduce the risk of ransomware

6

8

With the global cost of cybercrime estimated to be around \$6tr per year, it's evident that digital crime is just as lucrative for criminals as it is destructive for businesses. Hackers are capitalising on advances in technology to launch more-advanced and destructive attacks that are even harder to detect. Now more than ever, the pressure is on for organisations to focus on the protection of the data itself, including its restoration. Hackers need only one successful breach to compromise data. To mitigate this threat, businesses need a backup storage solution that makes critical data easily recoverable in the event of a loss. In addition to reviewing security processes, organisations should re-examine their storage infrastructure to optimise their protection against attacks. Florian Malecki of StorageCraft provides a series of steps that businesses should take to ensure they have the right elements in place to protect their data.

Hybrid intrusion detection system using machine learning

Cloud-based architectures have reduced IT barriers and provided new capabilities of dynamic provisioning, monitoring and managing resources by providing immediate access to resources, enabling the easy scaling-up of services. However, sharing the same pool when requesting services involves the risk of data breaches, account compromises, injection vulnerabilities and distributed denial of service (DDoS) attacks. As a result, many customers rank cloud security as a major challenge that threatens their work and reduces their trust in cloud service providers. In order to enhance cloud users' satisfaction and trust, a significant security improvement is required. Amar Meryem and Bouabid EL Ouahidi propose an architecture that eradicates malicious behaviours by detecting known attacks using log files; blocks suspicious behaviours in real time on behalf of recent architecture requests; secures sensitive data; and establishes better adaptations of security measures by dynamically updating security rules.

REGULARS

ThreatWatch	3
Report Analysis	4
News in brief	5
The Firewall	20
Events	20

ISSN 1353-4858/20 © 2020 Elsevier Ltd. All rights reserved

This publication and the individual contributions contained in it are protected under copyright by Elsevier Ltd, and the following terms and conditions apply to their use: Photocopying

Single photocopies of single articles may be made for personal use as allowed by national copyright laws. Permission of the publisher and payment of a fee is required for all other photocopying, including multiple or systematic copying, copying for advertising or promotional purposes, resale, and all forms of document delivery. Special rates are available for educational institutions that wish to make photocopies for non-profit educational classroom use.

Editorial Office:

Elsevier Ltd The Boulevard, Langford Lane, Kidlington, Oxford, OX5 1GB, United Kingdom Tel: +44 1865 843239 Web: www.networksecuritynewsletter.com

Publishing Director: Sarah Jenkins

Editor: Steve Mansfield-Devine E-mail: smd@contrarisk.com

Senior Editor: Sarah Gordon

Columnists: Ian Goslin, Karen Renaud, Dave Spence, Colin Tankard

International Editoral Advisory Board: Dario Forte, Edward Amoroso, AT&T Bell Laboratories; Fred Cohen, Fred Cohen & Associates; Jon David, The Fortress; Bill Hancock, Exodus Communications; Ken Lindup, Consultant at Cylink; Dennis Longley, Queensland University of Technology; Tim Myers, Novell; Tom Mulhall; Padget Petterson, Martin Marietta; Eugene Schultz, Hightower; Eugene Spafford, Purdue University; Winn Schwartau, InterPact

Production Support Manager: Lin Lucas E-mail: I.lucas@elsevier.com

Subscription Information

An annual subscription to Network Security includes 12 issues and online access for up to 5 users. Subscriptions run for 12 months, from the date payment is received.

More information: www.elsevier.com/journals/ institutional/network-security/1353-4858

Permissions may be sought directly from Elsevier Global Rights Department, PO Box 800, Oxford OX5 1DX, UK; phone: +44 1865 843830, fax: +44 1865 853333, email: permissions@elsevier.com. You may also contact Global Rights directly through Elsevier's home page (www.elsevier.com), selecting first 'Support & contact', then 'Copyright & permission'. In the USA, users may clear permissions and make payments through the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, USA; phone: +1 978 750 8400, fax: +1 978 750 4744, and in the UK through the Copyright Licensing Agency Rapid Clearance Service (CLARCS), 90 Tottenham Court Road, London W1P 0LP, UK; tel: +44 (0)20 7631 5555; fax: +44 (0)20 7631 5500. Other countries may have a local reprographic rights agency for payments.

Derivative Works

Subscribers may reproduce tables of contents or prepare lists of articles including abstracts for internal circulation within their institutions. Permission of the Publisher is required for resale or distribution outside the institution. Permission of the Publisher is required for all other derivative works, including compilations and translations.

Electronic Storage or Usage

Permission of the Publisher is required to store or use electronically any material contained in this publication, including any article or part of an article. Except as outlined above, no part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without prior written permission of the Publisher. Address permissions requests to: Elsevier Science Global Rights Department, at the mail, fax and email addresses noted above.

Notice

No responsibility is assumed by the Publisher for any injury and/ or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions or ideas contained in the material herein. Because of rapid advances in the medical sciences, in particular, independent verification of diagnoses and drug dosages should be made. Although all advertising material is expected to conform to ethical (medical) standards, inclusion in this publication does not constitute a guarantee or endorsement of the quality or value of such product or of the claims made of it by its manufacturer.

> 12987 Digitally Produced by Mayfield Press (Oxford) Limited

... Continued from front page

The hackers are affiliated with the Chinese Government, according to the statement. It adds: "Chinese government cyberthreat actors are actively exploiting trust relationships between information technology (IT) service providers – such as managed service providers and cloud service providers – and their customers", and one of the suggested mitigations for organisations is to, "ensure their providers have conducted a review to determine if there is a security concern or compromise and have implemented appropriate mitigation and detection tools for this cyber activity."

The warning came shortly after another joint statement between CISA and the UK's National Cyber Security Centre (NCSC), which claimed that advanced persistent threat (APT) groups are targeting healthcare and research organisations in the UK and US.

There's more information here: www.us-cert.gov/china.

The CISA/FBI warning did not provide specific details about the targets of attacks – it just stated that investigations are underway.

Although probably unrelated, there seems to have been a rise in sophisticated hacks against targets used for research. For example, around a dozen high-performance computing (HPC) facilities in Germany, the UK and Switzerland have had to shut down.

On 11 May 2020, Archer, the UK's National Supercomputing Service, found 'security exploitation' on its login nodes. It shut off all access and has invalidated all passwords and SSH keys. At the time of writing, the service was still offline and has warned users that all future access will require multi-factor authentication. There's more information here: www.archer.ac.uk/status/.

On the same day, the Baden-Württemberg High Performance Computing (bwHPC) project in Germany announced that five computing clusters across the country had to be shut down due to an unspecified security incident. A few days later, the Leibniz Supercomputing Centre and the Jülich Supercomputing Centre (JSC), both also in Germany, also went offline due to security issues. Shortly after, the Swiss Centre of Scientific Computations (CSCS) detected malicious activity on its systems and cut off access.

These HPC facilities are used for research in a wide range of disciplines, including chemistry, bioinformatics, physics and others.

Meanwhile, two construction firms in the UK that have been involved with the building of emergency hospitals to help deal with the pandemic have come under separate cyber attacks. Bam Construct, which worked on the Yorkshire and Humber hospital, seems to have suffered a straightforward ransomware attack which, the firm said, caused relatively little disruption.

A attack against Interserve, which worked on Birmingham's NHS Nightingale hospital and is also a major contractor to the Ministry of Defence, has had a major data breach. Around 100,000 records relating to former and current employees have been compromised. The data includes names, addresses, bank details, payroll information, next-of-kin details, personnel and disciplinary records. The firm has released little in the way of detail other than that the breach occurred "earlier this month".

InfinityBlack market taken down

Polish and Swiss law enforcement agencies have taken down the InfinityBlack underground market that was being used to trade hundreds of millions of stolen user credentials and sell hacking tools as well as committing other fraudulent activities.

The Polish National Police (Policja) arrested five people and seized equipment, hard drives and crypto-currency wallets worth around €100,000.

The operation, supported by Eurojust and Europol, also took down two other online platforms dealing in stolen credential databases with more than 170 million records.

"The hacking group created online platforms to sell user login credentials known as 'combos'," said Europol. "The group was efficiently organised into three

Threatwatch

RATicate targets industry

Sophos has spotted a hacking group that is using Nullsoft Scriptable Install System (NSIS) - an open-source tool used to create Windows installers - to sneak remote access trojans (RATs) and information-stealing malware into the networks of industrial organisations. Dubbed RATicate, the group is using the technique to install a variety of payloads, including LokiBot, FormBook, BetaBot, Agent Tesla and Netwire. The installers also drop a number of 'junk' files onto the system, presumably as an anti-analysis technique to disguise the malware's real actions. Sophos said it had identified six spamming campaigns being used to spread the malware, targeting industrial firms in Europe, the Middle East and South Korea. There's more information here: bit.ly/2WJlHcr.

Malicious incident report

Outpost24's Ghost Labs OffSec team has discovered multiple stored cross-site scripting vulnerabilities in the popular ServiceNow IT service management software. By taking advantage of this vulnerability, an attacker can create a malicious incident request that can then be sent out to users in the platform via a direct link to the request. If an attacker succeeds in exploiting this vulnerability, it would be possible to execute

defined teams. Developers created tools to test the quality of the stolen databases, while testers analysed the suitability of authorisation data. Project managers then distributed subscriptions against cryptocurrency payments. The hacking group's main source of revenue came from stealing loyalty scheme login credentials and selling them on to other, less technical criminal gangs. These gangs would then exchange the loyalty points for expensive electronic devices."

Using sophisticated scripts, a Swiss group was able to access a large number of accounts in the country. Europol put actual losses at €50,000, but said the group had access to funds of more than €650,000. Their activity came to light when some of the group – mainly minors and young adults – used the stolen data in shops in Switzerland in an attempt to cash out. Using shared intelligence sources and the services of Europol's network of cyber liaison officers (J-CAT), Swiss police then found links to a separate hacking group in Poland. custom JavaScript on the victim's browser. This means that the attacker could redirect the victim to a fake version of the ServiceNow site for social engineering purposes. It is believed that all customer installations are affected and possibly older ServiceNow installations as well. There's more information here: https://bit.ly/3eaS2Ph.

Kaiji targets Linux and IoT

A security researcher going by the name MalwareMustDie and the team at Intezer Labs have identified a new strain of malware that targets Linux servers and Internet of Things (IoT) devices. Dubbed Kaiji, the malware infects machines via root accounts, using brute-force SSH attacks, to use them as bots for distributed denial of service (DDoS) attacks. Unusually, the code was written using Go, rather than the more usual C/C++. This suggests some effort by the creators as most malware takes advantage of large libraries of exploit code, which isn't available in the Go language. The malware has been spotted in the wild, but seems to be making slow progress so far. There's more information here: https://bit.ly/3g4ruRt.

New features for Anubis

Anubis, an especially pernicious form of information-stealing malware for the Android platform,

New North Korean threats

The Lazarus Group – a state-backed hacking operation based in North Korea and also known as Hidden Cobra and Zinc – has evolved new forms of malware, primarily aimed at spying and network intrusion, including a new form of a remote access trojan (RAT) designed to work on Apple Macintosh machines.

The Dacls RAT was first seen by Qihoo 360 NetLab in December 2019 being used against Windows and Linux systems. Now it has been adapted for macOS and is being spread via a trojanised version of a two-factor authentication (2FA) app called MinaOTP, which is popular in China. The malware can operate as a fully covert RAT with capabilities including command execution, file management, traffic proxying and worm scanning. It communicates with its command and control (C&C) server using TLS encrypted connections.

There's an analysis of the malware at

may be getting clever new features soon, according to Alex Holden, founder of Hold Security. He has gained access to a new version of a web-based module used by Anubis operators to control infected devices. One of the new features is an eyeball icon that allows operators to halt malicious activities if the user is looking at the device – presumably to avoid detection. Anubis has been used extensively in the past for cyber espionage before pivoting to stealing banking credentials. The new features have not gone live yet, but Holden expects them to be available soon.

Android banking trojan

Cybereason's Nocturnus team has released details of new Android mobile malware targeting users of more than 200 financial apps, including banking, money transfer services and cryptocurrency wallets. EventBot is targeting users in countries across Europe and the US, and abuses Android's accessibility features to steal user data from financial applications, read user SMS messages and steal SMS messages to allow the malware to bypass two-factor authentication. The more than 200 financial apps impacted include: Paypal Business, Barclays, UniCredit, CapitalOne UK, HSBC UK, Santander UK, TransferWise, Coinbase and many more. There's more information here: https://bit.ly/2yjTTlu.

the Malwarebytes blog here: https://bit. ly/2TgHV3n.

Meanwhile, the US Department of Homeland Security's Cyber security and Infrastructure Security Agency (CISA) has issued details of three new malware variants also associated with the Lazarus Group.

CopperHedge is also a RAT - part of the Manuscrypt family used to target crypto-currency exchanges. It is capable of running arbitrary commands, performing system reconnaissance and exfiltrating data. TaintedScribe is a trojan that disguises itself as Microsoft's Narrator and has the capability to download, upload, delete and execute files; enable Windows CLI access; create and terminate processes; and perform target system enumeration. Finally, PebbleDash is a trojan that uses FakeTLS for session authentication and for network encoding and has similar capabilities.

There's more information, including indicators of compromise and mitigations, here: www.us-cert.gov/northkorea.

Report Analysis

AT&T Cybersecurity: The relationship between security maturity and business enablement

When you question cyber security practitioners about the specific dangers an organisation might face, or what countermeasures it should put in place, or how vulnerable it might be to a certain threat, an answer that crops up many times is, 'It depends on the organisation's security maturity'. That's undoubtedly true and also, often, unhelpful.

While there are plenty of frameworks and questionnaires out there that might give you some sort of metric, not everyone uses the same approach, and so comparing one organisation with another is tricky. A lot depends on perspective – about how much security matters to an organisation (especially when it comes to budgeting). And what does 'mature' mean anyway? After all, every organisation is vulnerable to something – it's just that some haven't found out what it is yet.

In wanting to evaluate how security maturity correlates with business success, AT&T Cybersecurity and its project partner the Enterprise Security Group (ESG) needed to define their own metrics for maturity. They based this around the US National Institute of Standards and Technology (NIST) cyber security framework (CSF), which breaks down security into five main functions – identification, protection, detection, response and recovery. The NIST CSF defines best practices that can create a secure environment that enables an organisation to innovate and to operate with maximum agility.

According to the report, the maturity model built by AT&T is based around questions covering processes, policies and controls in use, such as: how formalised is the organisation's cyber security programme? How frequently does it provide cyber security training? How diligently does it identify and prioritise threats? How is threat intelligence brought to bear? How extensively are data and assets segmented and encrypted? What technologies are used in event identification and resolution? How often is the organisation's security posture evaluated and revised over time? Based on the responses, organisations were grouped into three maturity levels: emerging (those that scored in the bottom 39%), following and leading (the top 20%).

It's as well to be cautious about such classifications. These are arbitrary divisions within AT&T's dataset. If the company happened to have chosen organisations that were all pretty bad at security, a fifth of them would still be classed as 'leaders'. It doesn't mean they are good at security, per se – just better than the others.

With that caveat in mind, how do these categories match up against business performance? The answer is, pretty much as you'd expect, but with a couple of interesting surprises.

Perhaps the least surprising conclusion is that 'leading' organisations – the ones with the greatest commitment to security – match the most successful businesses in the survey. "Leading organisations can weave strong cyber security into the business, IT and organisational culture," says the report. "This helps them be more aggressive with IT-driven business



initiatives, knowing they can count on a strong security foundation."

Leading organisations are also further advanced in implementing the five functions of the NIST CSF, especially when it comes to threat protection and incident response. They also know where they are weak in terms of security capabilities and are thus more likely to seek help from outside specialists.

No organisation is ever fully secure and those that topped AT&T's chart understand that better – that security is not a matter of installing a solution and thinking you're safe. It's an ongoing process. But one of the consequences of this is that leading organisations tend to get a better return on their security investments.

Of course, the amount they spend tends to be significantly higher than firms in the other two categories. Having said that, perhaps the most interesting and encouraging conclusion from this study is that security maturity is not dependent on the size of the organisation.

"The research shows that the median company size is identical across all three maturity levels," says the report. "The fact that there is no correlation between company size and maturity level indicates to us that doing cyber security well is less a function of resources and more a function of thoughtful consideration, planning and organisational culture."

The research also makes the point that security can and should be seen as an enabler – not just by IT people or even the board but by managers throughout the organisation. The link between security maturity and business success, "is likely anchored by trust, communication and collaboration between people – managers and staff from lines of business (LOB) and cyber security teams," the report claims. "Just over one quarter (26%) of respondents say that security is viewed as an enabler by line-of-business stakeholders. When this data is viewed through the maturity model, however, security teams are seen as 'enablers' by LOBs at 55% of leading organisations."

So how does an organisation put itself in the 'leading' category? According to AT&T, a couple of things are key. One is that proper alignment between security strategies and critical business assets is essential. This means putting your resources and money into protecting those things that matter most. Threat intelligence can play a key role here. Organisations also need to get their incident response act together, and that includes making sure that security event data is managed effectively and acted upon.

https://cyber security.att.com/resource-centre/ security-maturity-and-business-enablementsurvey-brief.



In brief

Top 10 threats

The US Government, via the Cyber security and Infrastructure Security Agency (CISA), has issued a list of the top 10 routinely exploited vulnerabilities. It lists CVE numbers for the most common attack vectors during the period 2016-2019, as well as looking at some of those cropping up frequently so far in 2020. In what will come as a surprise to no-one, Microsoft's Object Linking and Embedding (OLE) technology features heavily. This is used to enable embedded content between applications such as Word and Excel and is often subverted to create maliciously crafted documents that are sent as email attachments. OLE technology was also used in the top three threats associated with nation-state attack groups in China, Iran, North Korea and Russia. After OLE, the Apache Struts web framework was the next-most-exploited technology. Unfortunately, many organisations do not seem to be getting the message about software vulnerabilities. CISA notes that in December 2019, Chinese state-backed hackers were still exploiting a vulnerability (CVE-2012-0158) that the agency had assessed as being their most-used vector back in 2015. "This trend suggests that organisations have not yet widely implemented patches for this vulnerability and that Chinese state cyber actors may continue to incorporate dated flaws into their operational tradecraft as long as they remain effective," the agency says. In 2020, many attackers seem to be going after unpatched flaws in virtual private networks systems - notably in products from Citrix and Pulse Secure. There's more information here: www.us-cert.gov/ncas/ alerts/aa20-133a.

Criminals exploit RDP

With many employees now working from home due to the Covid-19 pandemic, cyber criminals are putting a lot of effort into attacking remote desktop protocol (RDP) applications. RDP is used to allow people to log into a computer and use it remotely, or to receive technical support at a distance. According to McAfee, the number of RDP ports exposed to the Internet grew from three million to 4.5 million in the period from January to March 2020. This has been matched by an increase in attacks against RDP ports and a boom in the sale of stolen credentials on hacker marketplaces. Kaspersky has also seen a massive uptick in RDP attacks, especially brute-force login attempts. "As far as we can tell, following the mass transition to home working, they logically concluded that the number of poorly configured RDP servers would increase, hence the rise in the number of attacks," the company says in a new report. For example, in the US, brute-force attacks against Internet-facing RDP servers have increased from 200,000 per day in early March to more than 1,200,000 during mid-April. The Kaspersky report is here: https://bit.ly/2LAIa4U.

Attack on UK power grid...

A company that supplies important services to the UK power grid has been hit with a cyber attack. The strike against Elexon, which monitors supplies to the grid and reconciles payments between energy companies, did not have an affect on the electrical supply, but it could have had an impact on the financial mechanisms in the market. Elexon hasn't revealed many details, other than that the attack affected internal systems and laptops, but it looks as though it might have been a ransomware infection that was largely contained before it could do too much damage. There's more information here: https://bit.ly/2Ze1ugk.

...and Israeli water

The Israeli National Cyber-Directorate (INCD) has advised organisations working in the energy and water supply sectors in the country to change passwords for all systems with connections to the Internet. With any systems where that's not possible, they should take them offline pending the implementation of tighter security. Similar alerts were issued by Israel's Computer Emergency Response Team (CERT) and by the government Water Authority. The latter told firms to change credentials, "with emphasis on operational systems and chlorine control devices in particular". The alerts were prompted by a report sent to the Government by cyber security firm ClearSky, which has been tracking an Islamist group's activities on social media. The group, Jerusalem Electronic Army, often posts screenshots from attacks it claims to have made against targets, including Israeli universities and government systems. However, there have been no confirmed reports of successful intrusion at water treatment and supply companies.

Microsoft takes down botnet

Microsoft's Digital Crimes Unit (DCU) has succeeded in taking down a 400,000-strong botnet that was being controlled by a single Internet of Things (IoT) device - an LED light control console. The botnet was being used for activities such as phishing campaigns, malware distribution, ransomware payload delivery and the launch of distributed denial of service (DDoS) attacks. As much as 1TB of malicious data was being sent out per week. Working with Taiwan's Ministry of Justice Investigation Bureau (MJIB), the DCU tracked down a single IP being used as a command and control server. This turned out to be a compromised Internet-enabled lighting console. By shutting down this device, the botnet has effectively been taken offline. There's more information here: https://bit.ly/2AGeu4f.

Trump held to ransom

The REevil (aka Sodinokibi) ransomware group, which appears to have stolen large

amounts of data from US law firm Grubman Shire Meiselas & Sacks (GSMLaw), is claiming that its haul includes significant amounts of "dirty laundry" on US President Donald Trump. Most of the stolen data relates to celebrities in the entertainment industries, and the group has already published samples relating to Lady Gaga and others. GSMLaw responded to the demand for a \$42m ransom, and the specific threat to Trump by saying that the FBI has classed this action as 'terrorism' and that "negotiating with or paying a ransom to terrorists is a violation of federal criminal law". This has prompted the REevil group to release 169 emails on the dark web which it claims is the "most harmless information". However, the emails appear to have nothing to do with Trump.

DDoS attacks respond to pandemic

The targets of distributed denial of service (DDoS) attacks have shifted somewhat to reflect life during a pandemic, according to research by Kaspersky. The most targeted resources in the first quarter of 2020 have been the websites of medical organisations, delivery services and gaming and educational platforms. Some of these attacks could be political or nation-state backed. For example, in mid-March attackers made an unsuccessful attempt to disable the website of the US Department of Health and Human Services (HHS), probably with the aim of preventing visitors from obtaining official data about the pandemic. At the same time, cyber actors were busy spreading misinformation on social networks and via text and email about the introduction of a nationwide guarantine in the US. An attack on the Paris-based group of hospitals Assistance Publique-Hôpitaux de Paris meant that remote hospital workers were unable to use software and email for some time. Other attacks were more likely to have been purely criminal, such as those against food delivery services Lieferando (Germany) and Thuisbezorgd (Netherlands). Online gaming platforms, which have been under heavy load during lockdowns, also came in for attack, including Battle.net, Eve Online and Wargaming. There's more information here: https://bit.ly/3bIw70b.

French thwart jackpotting scam

French authorities claimed to have taken down an international network involved in the 'jackpotting' of ATMs. This is where malware is inserted into the machines so that money mules can make them disburse all the cash they hold. The hacking group is believed to have been involved in 19 such incidents in France, which paid out \in 280,000. Now, Paris prosecutor Remy Heitz has announced that two suspects from the "Russian-speaking community" have been charged and are being held in custody.

Optimising storage processes to reduce the risk of ransomware



Modern-day hackers and cyber criminals are smarter than ever. With the global cost of cybercrime estimated to be around \$6tr per year, it's evident that digital crime is just as lucrative for criminals as it is destructive for businesses.¹ What's more, hackers are capitalising on advances in technology to launch more-advanced and destructive attacks that are even harder to detect.

The simple universal rules of this game are that all organisations face the danger of a hacker destroying their data, and security measures are only effective until (not if) they're breached. Now more than ever, the pressure is on for organisations to not only focus on the prevention, but also on the protection of the data itself, including its restoration. Hackers need only one successful breach to compromise data.

"Unstructured data is far less predictable and grows with unprecedented scale and speed. Most enterprise IT storage infrastructures are not designed with the scalability needed to handle this amount of data"

To mitigate this threat, businesses need a backup storage solution that makes critical data easily recoverable in the event of a loss. In addition to reviewing security processes, organisations should re-examine their storage infrastructure to optimise their protection against attacks. Following is a series of steps that businesses should take to ensure they have the right elements in place to protect their data.

Data is changing

Data is at the core of any digital transformation initiative. It provides the foundation for understanding where the business is positioned and for making informed decisions about the direction in which it is headed. As organisations grow, it is crucial that their technology systems develop in tandem, so they are able to best support, facilitate and protect business expansion. Many companies find that as they grow, the volume of data they send, receive and store grows exponentially, whether it's in the form of employee data, customer data or simply administrative content.

In addition to traditional data, the amount of unstructured data that companies are holding is exploding, with IDG Research predicting that 93% of all digital data will be unstructured by 2022.²

Finding correlations across multiple

disparate data sources

Unstructured data is freeform information that does not fit neatly into databases organised by fixed categories. It takes many different forms and includes, for example, audio and video files, PDFs, social media posts, machine-generated data such as medical 3D imaging, satellite imagery, or data borne from Internet of Things (IoT) devices. And with the rise in unstructured data, there are even more entry points for potential ransomware attacks.

For organisations that rely on legacy infrastructure systems to store and manage data, the influx of unstructured data poses a challenge. Structured data tends to grow predictably in a linear growth pattern and can be managed with traditional storage techniques. Unstructured data is far less predictable and grows with unprecedented scale

60%





and speed. Most enterprise IT storage infrastructures are not designed with the scalability needed to handle this amount of data, which can leave it incorrectly stored and, as a result, more vulnerable to a ransomware attack.

The importance of backup

Backing up all critical data and making it easy to recover in the event of an outage is one of the best lines of defence that a business can take against ransomware attacks. The most sophisticated attacks aim to encrypt both data files and Windows restore points, so this should be top of mind when installing a backup system.

These systems are under pressure from the amount of data being accumulated and stored. It's clear there is a problem, and a demand for appropriate backup and recovery strategies and systems.

A successful ransomware remediation and recovery plan should include:

- Immutable snapshots: To ensure that unstructured data can be recovered, companies should protect their information with continuous immutable snapshots. Data captured this way is 'frozen' and cannot be overwritten or deleted by ransomware attackers. This ensures that an organisation can revert to a secure set of data.
- Orchestration: A successful recovery process requires that business-critical data and applications are prioritised. Companies using cloud-based recovery should pre-determine the order in which their data and applications will be recovered. This 'orchestration' ensures minimal downtime, once data recovery begins.
- Immediate recovery: System downtime can prove inordinately costly for organisations. The speed of recovery following a ransomware attack is a crucial element of the remediation and recovery process. Storage solutions should provide the ability to recover virtual and physical infrastructures – and both structured and unstructured data – instantly.

Clear plan

Recent research has found that 68% of global businesses believe they have a clear plan in place and would be capable of quickly recovering from a ransomware attack.³ However, the research also found that nearly a quarter of the same companies do not test their recovery plans. And of those organisations that do test, nearly half only test their recovery plans once a year or less. This reveals a concerning disconnect between an organisation's confidence and its actual ability to recover from a ransomware attack, and underlines the importance of testing recovery plans.

In addition to performing regular backups, businesses should consider the following:

- **Update software:** Update all software according to a regular maintenance plan. If a workstation or server is too old to update, retire it. The few tasks it can perform do not outweigh the risk it creates for the other machines on the network.
- **Restrict usage:** Restrict administrator accounts to only a few people in the organisation. Create dedicated user (not admin) accounts on each work-station for each employee. End-users should not be logged into machines as administrators. The most destructive ransomware is designed to gain access to network areas only accessible via administrator accounts.
- Verify backups: Performing backups is just the first step because these will not be effective unless you know they work. And the only way to make sure they do is to verify backups by testing the data restore process. Occasionally the backup restores properly but does not include all critical files, so make sure your backup brings back all your data.
- **Training:** Employee training is often overlooked or not regularly updated for new employees. Do not assume employees are tech-savvy enough to recognise malware sent over email. Regular training takes time and resources, but apart from backup, it can have a major impact in deterring the spread of ransomware.

No backup?

For organisations without a current disaster recovery plan, a quality backup and restore solution must be implemented at once. Companies that have already fallen victim to ransomware and lost data due to a lack of appropriate security measures and/or backup must reassess their overall data protection policies and take the relevant prevention measures.

Organisations that have been compromised by ransomware must deploy their IT teams to assess the value of the data that has been encrypted. They need to make a decision as to whether it is worth hiring a security or ransomware expert to try to recover the data.

Some businesses might be tempted to pay the ransom, but that is not a good idea. Even if the ransom was to be paid, there is no guarantee of receiving the decryption keys and thieves often increase the ransom the longer they have to wait for it to be paid.

"Business leaders must realise that storage is not merely a container for enterprise data – it can also be designed to successfully mitigate the risks of ransomware"

In many ways, ransomware attacks are the perfect crime because the anonymity of the Internet makes it nearly impossible for authorities to track down the perpetrators. This means that cyber criminals often win and instead of being intercepted and stopped, they move forward to identify their next victims.

Looking to the future

Attacks will continue to happen and, in turn, will continue to evolve as companies learn to combat them. Businesses can no longer afford to sit back and hope that they will be the lucky ones that avoid an attack. Data is a highly sought-after asset, and its safeguarding must be of the utmost importance to businesses that wish to succeed in an increasingly threatening cyber landscape. Not only do thieves ask for ransom, the latest trends show that ransomware criminals now threaten to publicise the data they capture, adding a whole new dimension to their blackmailing strategy. Modern storage techniques are crucial in this fight against attacks. Business leaders must realise that storage is not merely a container for enterprise data – it can also be designed to successfully mitigate the risks of ransomware.

About the author

Florian Malecki is international product marketing senior director at StorageCraft.

He drives the development of the vendor's data protection and storage solutions. Prior to joining StorageCraft, Malecki worked in senior roles at SonicWall, Dell, Aventail, ClearSwift, Omgeo, Lucent Technologies and Air Products. He earned a master's degree in business, marketing and technology from the Ecole Supérieure Des Affaires et Des Technologies, France.

References

 Malecki, Florian. 'Best practices for preventing and recovering from a ransomware attack'. ITProPortal, 28 Jun 2019. Accessed May 2020. www. itproportal.com/features/best-practices-for-preventing-and-recoveringfrom-a-ransomware-attack/.

- '2016 Data & Analytics Research'. IDG, 7 May 2016. Accessed May 2020. www.idg.com/tools-formarketers/tech-2016-data-analyticsresearch/.
- StorageCraft research reveals need for a ransomware reality check'. StorageCraft, 11 Dec 2019. Accessed May 2020. www.storagecraft.com/ press-releases/storagecraft-researchreveals-need-ransomware-realitycheck.

Hybrid intrusion detection system using machine learning



Amar Meryem



Amar Meryem and Bouabid EL Ouahidi, Mohammed V University, Rabat

Recent technologies and innovations have encouraged users to adopt cloud-based architectures.^{1,2} This has reduced IT barriers and provided new capabilities of dynamic provisioning, monitoring and managing resources by providing immediate access to resources, enabling easy scaling up of services and implementation of new classes of existing applications. However, sharing the same pool when requesting services involves the risk of data breaches, account compromises, injection vulnerabilities, abusive use of features such as the use of trial periods and distributed denial of service (DDoS) attacks.^{3,4} As a result, many customers rank cloud security as a major challenge that threatens their work and reduces their trust in cloud service providers.

When a customer requests a cloud service, its connected machine generates log files. These flat documents contain chronological information about customer behaviour. They provide information about user and system identifiers, executed transactions and the time.⁵⁻⁸ Thus, log files are essential to provide feedback about activities, to provide an audit trail, to detect certain types of attack, to track insiders and to provide real-time alerts.

In order to enhance cloud users' satisfaction and trust, a significant security improvement is required. The proposed architecture eradicates malicious behaviours by detecting known attacks using log files, blocks suspicious behaviours in real time on behalf of recent architectures requests, secures sensitive data, and establishes better adaptations of security measures by dynamically updating security rules.

Detection techniques

Several techniques have been proposed to detect malware. Rule-based analysis identifies malicious behaviours based on static rules such as source code, serial numbers of developers, session fixation or attack signatures. It compares activities to known signatures of attacks and has two models. One is the negative security model (also known as a blacklist approach) that considers everything as normal except for blacklisted connections flagged as attacks. It is easy to implement and yields very few false positives. Nevertheless, predefined rules need to be constantly added to and adapted in line with new findings. The other approach – the positive security model – denies all connections by default and allows only whitelisted ones to pass. Firewalls are configured that way, where the whitelist is provided manually or during the learning phase.

Misuse detection helps predict anomalies in the short and medium term but considers each recorded event that is not on a list of known attack signatures as normal. However, it does not calculate the distance between a suspicious behaviour and a certain attack, nor does it predict new anomalies. Even worse, it is static and needs frequent updates.

Anomaly-based intrusion detection (ABID) - known as behavioural-based analysis when using machine learning algorithms - collects information at runtime, crafts a model and then matches each new comportment to the crafted model.^{9,10} This strategy is based on dynamic rules and usually builds a profile or a baseline of 'normal' traffic in a learning phase and considers every deviating behaviour as suspicious. Anomaly detection in the field of cyber security considers specific features in order to flag unknown behaviours as normal or outliers, such as the number of command executions, the number of times a specific server port was opened and how many times a user has requested root privileges.

"Anomaly detection in the field of cyber security considers specific features in order to flag unknown behaviours as normal or outliers"

Detecting anomalies by log file analysis was considered a manual misuse detection technique. In this paper, we see how it may be managed dynamically and considered as real-time log file processing.¹¹ We propose a hybrid machine learning architecture that analyses log events automatically to detect anomalies, and crafts real-time user profiles to decide whether behaviour is legitimate.

Related works

Jie Yang and Chen proposed a hybrid detection system combining both misuse detection and anomaly detection.¹² The misuse detection does a first discrimination of behaviours based on the packet protocol type, feature selection and anomaly detection phase discovery in newly discovered attacks. However, detection rates depend on the protocol type involved in the suspected connection. Experimentation was on the KDD dataset and results show that the Chisquare automatic interaction detector (CHAID) was at 95.93% detection rate and 183 items of data were falsely detected.

Norbik Bashash, Idris Bharanidharam and Abdul Manan anticipated a dynamic model for intrusion detection that includes neural network techniques and fuzzy logic along with data mining to process the network data.¹³ The proposed architecture mines both misuse and anomaly detection. For network intrusion detection, a Snort system first compares the input behaviour to historical mined data. If the input is suspicious, the detection phase is activated. Data mining is used with fuzzy logic known as 'fuzzy association rules'. The data mining discovers association rules from a large dataset and fuzzy logic upgrades the representation of the overlapping categories to create more abstract patterns. For host-based intrusion detection, the authors used self organising maps (SOM) to build user profiles.

Bambang Setiawan and Supeno Djanali used centroid-based classification for intrusion detection and provided composite indicators to evaluate the system.¹⁴ The evaluation was on three classification models – CANN, LSCANN and CASMN on NSL-KDD – and a selected group of individual indicators, namely accuracy, completeness, robustness and speed. LSCANN using feature selection gave the best accuracy, which varied between 97.61% and 99.74%, and a completeness (FAR) between 27.64% and 100%.

Zahoor Rehman, Sidra Nasim and Sung Wook Baik proposed a hybrid model to detect malware in Android Apps.¹⁵ This model is based on both signatures and heuristic analysis of the manifest.xml files. It compares, before each installation of a new Android package (APK), the constant strings of downloaded APK files to constant strings of malware applications and the manifest .xml files.

Yunaucheng Li, Rong Ma and Runhai Jiao used deep learning approaches in detecting malicious codes, based on autoencoders for feature extraction and a deep belief network (DBN) as a classifier.¹⁶ Results showed that the proposed model is more accurate than using only DBN. However, its major limitation is that the training dataset isn't updated. This process is different from ours. Our proposed solution updates the training dataset and the security system rules, merges different types of log files (audit, access, error, SSH error) to get all relevant features, then analyses them together to obtain network insight and to learn users' profiles.

"Using modular processing encourages high falsepositive rates, especially when there is a new normal behaviour that wasn't taken into consideration in the training phase"

Lekaning Djionang and Gilber Tindo proposed a modular architecture to improve intrusion detection systems.¹⁷ The first module is based on a technique that detects whether an unknown behaviour is legitimate or suspicious. If the input connection is suspicious, it is communicated to subsequent modules. Each module detects a unique type of attack. To do so, the system starts by normalising the NSL-KDD dataset and uses a two-class neural network. Then it identifies the best number of neurons to use in the hidden layer to have the best rate of recognition. Afterwards, each network is trained separately using a back-propagation algorithm. Results showed the importance of feature selection in increasing/decreasing the performance rate. However, using modular processing encourages high false-positive rates, especially when there is a new normal behaviour that wasn't taken into consideration in the training phase. We overcome these limits by using a hybrid system that updates the training dataset and calculates the distance between the normal and the unknown behaviour and compares it to a deviation measure when labelling it.18,19

Real-time big data processing

Real-time big data processing is required for the uninterrupted monitoring of

events, messages and processes on the network infrastructure.²⁰ In our case, log files are rapidly changing, inmemory datasets. The huge quantity of data that arrives at the pipeline is in different formats (SSH, audit, error and log files). However, these data contain detailed information about messages and events. Big data analysis enhances the analysis quality and reduces decision time responses.

Due to the accumulation of massive volumes and variety of data from different sources, several industries and institutions are moving towards real-time responses. To extract maximum value from moving data, cyber security engineers need to process data much faster and take timely actions, such as preventing the occurrence of suspect behaviours. The real-time processing improves correlation and pattern detection on a scale of millions of events and constantly moving data streams.

Among several big data anomaly detection techniques, we evaluate the six most-used ones.

K-nearest neighbour (KNN) is one of the most-used learning algorithms. It is a non-parametric learning algorithm whose purpose is to use a database of separated points and to predict the classification of a new sample.

KNN depends on the choice of distance parameter. It embraces the entire sampling set and incorporates the information in the set. The distance between each item in the sampling set must be processed for the purpose of classifying the points – the k closest passages in the sampling set are considered as the point in the far distance. One of the shortcomings is the similarity measure. The classification fails to calculate distances between points in a small subset. Plus, it requires memory computation.

A **Bayesian network** (BN) is a probabilistic model *P* on a finite directed acyclic graph (DAG). For each node *i* in the graph, there is a random variable X_i together with a conditional probability distribution $P(x_i|x_{\prod(i)})$ where $\prod(i)$ is the parent of *i* in the DAG. The joint probability distribution of the Bayesian network is the product of the conditional probability distributions.

$P(x_1,\ldots,x_n) = \prod_{i=1}^n P(x_i|x_{\prod(i)})$

Bayesian networks have been broadly utilised for grouping issues. With a BN grouped into a qualitative and a quantitative model, the qualitative piece of the system is spoken to by a co-ordinated non-cyclic graph, whose nodes denote the random factors in the problem domain whose edges systematise significance relations between the factors they interface.

Support vector machines (SVMs) have been used to construct a decision boundary, which has the most extreme edge between the typical data set and the source. The objective of the support vector machine algorithm is to find a hyperplane in N-dimensional space (where N is the number of features) that distinctly classifies the data points. That approach depends on a minimised classification threat as opposed to an optimal classification. It is helpful when the number of features, m, is high and the number of data points, n, is low (m >> n).

Hyperplanes are decision boundaries that help classify the data points. Data points falling on either side of the hyperplane can be attributed to different classes. Also, the dimension of the hyperplane depends upon the number of features. If the number of input features is two, then the hyperplane is just a line. Support vectors are data points that are closer to the hyperplane and influence the position and orientation of the hyperplane. Using these support vectors, we maximise the margin of the classifier. Deleting the support vectors will change the position of the hyperplane. These are the points that help us build our SVM.

Logistic regression was used in the biological sciences in the early 20th century. It predicts the probability of an outcome that can only have two values (binary output). The prediction is based on the use of one or several predictors (numerical and categorical). Logistic regression measures the relationship between the categorical dependent variable and one or more independent variables by estimating probabilities using a logistic function, which is the cumulative logistic distribution. Logistic regression is named for the function used at the core of the method, the logistic function.

The logistic function, also called the sigmoid function, was developed by statisticians to describe properties of population growth in ecology, rising quickly and maxing out at the carrying capacity of the environment. It's an S-shaped curve that can take any real-valued number and map it into a value between 0 and 1, but never exactly at those limits.

$$\sigma(t) = \frac{1}{1 - e^{-t}}$$

Random forest (RF) is a collaborative classifier used to improve accuracy. It consists of two stages – feature selection and classification. Random forest generates multiple decision trees from random subsets of data. One of the major advantages of random forest is that it yields low classification errors when compared with other traditional classifiers. Yet when working with large datasets and complex estimation procedures, RF consumes a lot of computational time.

K-means clustering is a simple and widely-used clustering algorithm. Given the value of k, it tries to build k clusters from samples in the dataset. Therefore, k is a hyperparameter of the model. In K-means, a cluster is a group of points, with a representative entity called a centroid. A centroid is also a point in the data space - the centre of all the points that make up the cluster. It is defined to be the arithmetic mean of the points. In general, when working with K-means, each data sample is represented in a d-dimensional numeric vector, for which it is easier to define an appropriate distance function.

Given k, K-means works as follows:

- 1. Randomly choose k data points (seeds) to be the initial centroids.
- 2. Assign each data point to the closest centroid.
- Re-compute (update) the centroids using the current cluster memberships.
- 4. If a convergence criterion is not met, go to step 2.

We can also terminate the algorithm when it reaches an iteration budget, which yields an approximate result.

Dynamic analysis and performance metrics

In cyber security, dynamic analysis relies on machine learning algorithms to detect malware. An IDS identifies the infection's patterns and controls the channels, such as identifying typically used websites, in simulating an attack and sorting out malicious executables from benign ones. The opposite of static analysis, which is very slow in detecting unknown attacks and requires few resources, dynamic analysis requires more resources for training the system at runtime but is much better at identifying new, unknown malware. As a result, all recent contributions focus on reducing the computational requirements and false positive key metrics.²¹

To visualise the performance of each adopted algorithm we plot a confusion matrix. It is a table with two rows and two columns that reports the number of true negatives (TN), false positives (FP), false negatives (FN), and true positives (TP).

		Predicted class			
		No	Yes		
Observed	No	TN	FP		
class	Yes	FN	TP		
Table 1: Confusion table.					

To evaluate the truthfulness of each new system and decide whether to trust it, there are several performance metrics for attack detection that we have used to evaluate our solution.

Accuracy (ACC) is the most common evaluation metric for the classification problem. It shows the overall effectiveness of a classifier. The accuracy is the ratio of the number of correctly predicted behaviours to the overall behaviours.

$$ACC = \frac{TP+TN}{TP+TN+FP+FN}$$

Detection rate (DR) is the ratio of correctly predicted behaviours to the total number of perfect predictions and wrongly, negatively labelled behaviours.

$$DR = \frac{TP}{TP+FN}$$

False alarm rate (FAR) is the ratio of wrongly predicted behaviours over the

total number of true negatives and false positives. It generally calculates the error of the detection system decision when it is caused by noise data.

$$FAR = \frac{FP}{TN+FP}$$

Precision shows the agreement of the data labels with predicted labels given by the classifier. It is the ratio of correctly predicted positive behaviours to the total number of positive observations. The higher the precision, the lower the false positive rate and the better the detection system.

$$Precision = \frac{TP}{TP + FP}$$

Recall shows the effectiveness of a classifier in detecting positive labels. Also called sensitivity, it is the ratio of correctly positive predicted behaviours to all actual class observations (TP and FN). The higher it is, the lower are the false negatives and true positives, which makes it insufficient to decide whether a system is trusted.

$$Recall = \frac{TP}{TP + FN}$$

F1-Score: This shows the relation between data positive labels and those given by the classifier. This score is more useful than accuracy when false positives and false negatives do not have the same costs.

$$F1 Score = \frac{2 \times Recall \times Precision}{Recall + Precision}$$

Receiver operating characteristics (**ROC**) is a popular performance metric that evaluates and compares the robustness of different classifiers.²² It is based on two basic evaluation measures: the FPR and FAR and sensitivity.

Proposed solution

Static analysis is fast and effective when it is based on predefined rules (Snort) and generates low false positive alarm rates.²³ However, various techniques can be used to dodge the system and enable it to detect polymorphic attacks.²⁴ Misuse detection can only detect attacks that are previously stored and may easily lack accuracy when attack signatures are outdated or when elected features are not resilient predictors, such as the Initial Sequence Number (ISN) used to define a trojan that alters a lot, or in the case of zero-day attacks. Anomaly detection can be utilised to acquire the signature information used by misuse-based IDS: it identifies attacks even when information is incomplete and no written rules are needed. However, it fails due to the high number of false positives, the definition of rules is difficult, the selection of the nature of the attack is impossible and it has a low detection rate.

In order to overcome the limitations of anomaly and misuse detection, this article proposes a hybrid approach that combines both signatures and 'normal' profile comparison in detecting attacks (Figure 1). Combining both techniques reduces the labelling time, enhances the detection and accuracy rates, reduces the number of false positives, enables the identification of new attack signatures and makes the updating of attack rules sets dynamic.

The adopted architecture does not only combine ABID and knowledge-





based intrusion detection (KBID) but also considers relevant factors in detecting intrusions based on their weight. It detects new attack signatures and dynamically updates the training data to reduce the number of false positives.

The main stages of our architecture (Figure 2) are:

- Step 1: Structuring log events.
- Step 2: Eradicating redundancy.
- Step 3: Classifying unlabelled behaviours.
- Step 4: Labelling user behaviours.
- Step 5: Identifying new signatures of attacks and updating security rules.

Structuring log events

Log files come from different layers and contain different information but are complementary and necessary for attack detection. We involve access logs to get features about network information (time, service, protocol, URL, transferred data, domain of accessed websites, status of the connection or attempt, etc). SSH

Network Security

logs give us insight into the number of failing attempts, the responses of the IDS and firewalls, the status of connection at flags (S0, S1, RSTO) granularity and so on. We can audit log records to find out about accessed resources (attempts to access sensitive files, password files, privileged documents and more) plus the destination that was targeted.

All these types of logs have different formats, are unstructured and have useless and repetitive information. These characteristics make the analysis a tedious task that takes significant time to train input datasets.

In this step, we start by moving all generated log files into a single cloud provider. It prevents insiders from deleting their traces and enhances cloud security auditing by uncovering complete information about users' manipulations. Generated log files amount to a size that makes their management, tracking and analysis better and easier when using a suitable predictive big data algorithm.

In this first step, the algorithm converts stored log data into a $41 \times N$ matrix named X_{log} where:

- matrix named X_{logs} where:
 41 is exactly the number of NSL-KDD features.²⁵
- *N* is the total number of centralised log files records.

Let $E = (e_1, ..., e_{41})$ be a set of NSL-KDD attributes, e_j the j^{th} attribute of the NSL-KDD data-set.

Let $X = \{x_1, ..., x_n\}$ be a set of $n \log records$. Each record is converted to a 41-dimensional vector, with x_i the i^{th} projected log row on E and x_{ij} the corresponding value of e_j in the $i^{th} \log$ row (see Figure 3).

Eradicating redundancy

Considering the increasing scale of log file volume and velocity, the resulting matrix from the previous step amounts to a similar size. Given the size of this matrix, traditional algorithms would be inadequate to compute frequencies and extract existing correlations between NSL-KDD attributes. In this step, we will be using a MapReduce program that takes in the X_{logs} matrix as input and calculates the number of occurrences



of each NSL-KDD feature. Figure 4 describes the MapReduce program.

In the map step, mappers attribute a unique identifier to each row of the matrix and initiate the occurrences to 1. It is important to note that each distinct row means a distinct behaviour present in the log file.

In the reduce step, reducers calculate the occurrences of each event row in the matrix. The output of the previous processing is a list of (X_i, v_i) couples, where C_i is the unique identifier of a distinct row and v_i is its occurrence. Then, each frequency is considered as a weight of the X_i user activity.

As a result, X_{logs} is reduced into a 41 × P matrix X where $P \le N$ and:

$$N = \sum_{i=0}^{p} v_i$$

Classifying unlabelled behaviours

At this stage, the log matrix is structured and reduced to a manageable size. In order to make the labelling easier, less time-consuming and more efficient in the upcoming step, this section is dedicated to clustering the N unknown structured logged events in X into k attack classes using K-means.^{26,27} Each of the classes has a centroid and the algorithm aims to minimise the distance between each point of a cluster and its centroid.

The basic idea behind partitioning X using K-means clustering is to define clusters such that the total intra-cluster variation (or the total within-cluster sum of square – WSS) is minimised.²⁸ The total WSS measures the compactness of the clustering and we want it to be as small as possible.

The elbow method looks at the total WSS as a function of the number of clusters. One should choose several clusters so that adding another cluster doesn't improve much over the total WSS.

In Figure 6, we see an elbow chart showing the SSE after running K-means clustering for k going from 1 to 10. For 352 events (first dataset), there is a clear elbow at k = 2, indicating that three is the best number of clusters. For a larger dataset of event logs (second dataset,





204,449 events) the elbow is three, which shows that yellow, grey and red classes may be joined into one class as plotted in Figure 5.

Labelling behaviours using NSL-KDD

The KDD'99 and ISCX datasets are the most commonly used datasets in

intrusion detection research. They reflect the traffic composition and have several characteristics that make them the appropriate datasets to use in our dynamic-based analysis. Being extensible, modifiable and reproducible, we combine in our approach both NSL-KDD features and log file insight into one architecture that maximises the accuracy of the derived analysis from

Figure 6: Elbow results on log events. Figure 6: Elbow $0.015 - \frac{x}{1} + \frac{x}{2} + \frac{x}{3} + \frac{x}{4} + \frac{x}{5} + \frac{x}{6} + \frac{x}{7} + \frac{x}{1} + \frac$	× × 8 9 5
events. The Elbow Method for 204449 even 0.10 - 0.08 -	s
0.10 -	
50.06 - 0.04 - 0.02 -	
	× × 8 9
% ACC DR FAR Precision Recall	F1
K-Nearest Neighbours (KNN)	
Vormal 98.89 98.80 0.91 99.09 98.80	98.95
Dos 99.93 85.29 44.12 65.91 85.29	74.36
12R 99.73 99.64 0.39 99.61 99.64	99.63
Victor 99.39 90.58 14.94 85.84 90.58 Victor 00.67 00.42 4.55 00.44 50.42	88.15
2000 9967 9967 9813 155 98711 9813	
Jaïve Baves (NB)	90.23
Value 50.15 1.55 50.144 50.15 Value 40.00 13405.22 0.45 60.87	0.90
Naïve Bayes (NB) 1.35 30.44 50.15 Vormal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51	0.90
Naïve Bayes (NB) 1.35 30.144 30.134 Vormal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 J2R 58.88 46.29 1.15 97.58 46.29	0.90
Naïve Bayes (NB) 55.13 1.33 30.44 56.13 Normal 47.95 60.87 13405.22 0.45 60.87 OoS 99.04 5.51 10.66 34.09 5.51 J2R 58.88 46.29 1.15 97.58 46.29 \2L 97.19 45.83 79.81 36.48 45.83	0.90 9.49 62.79 40.63
Naïve Bayes (NB) 1.33 30.44 30.13 Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35	0.90 9.49 62.79 40.63 56.04
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35	0.90 9.49 62.79 40.63 56.04
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) 92.08 4.65 95.20 92.08	0.90 9.49 62.79 40.63 56.04 93.61
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23	0.90 9.49 62.79 40.63 56.04 93.61 31.58
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74	0.90 9.49 62.79 40.63 56.04 93.61 31.58 96.65
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74 R2L 97.98 72.37 119.32 37.76 72.37	0.90 9.49 62.79 40.63 56.04 93.61 31.58 96.65 49.62
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74 R2L 97.98 72.37 119.32 37.76 72.37 Probes 97.13 86.19 17.44 83.17 86.19	0.90 9.49 62.79 40.63 56.04 93.61 31.58 96.65 49.62 84.65
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74 R2L 97.98 72.37 119.32 37.76 72.37 Probes 97.13 86.19 17.44 83.17 86.19 Support Vector Machine (SVM) Normal 95.27 92.76 92.76	0.90 9.49 62.79 40.63 56.04 93.61 31.58 96.65 49.62 84.65
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74 R2L 97.98 72.37 119.32 37.76 72.37 Probes 97.13 86.19 17.44 83.17 86.19 Support Vector Machine (SVM) Normal 95.27 92.76 1.29 98.63 92.76 DoS 99.85 0.00 0.00 0.00 0.00 0.00	93.61 9.49 62.79 40.63 56.04 93.61 31.58 96.65 49.62 84.65
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.98 72.37 119.32 37.76 72.37 Probes 97.13 86.19 17.44 83.17 86.19 Support Vector Machine (SVM) Normal 95.27 92.76 1.29 98.63 92.76 DoS 99.85 _ _ 0.00 _ _ I2R 98.35 98.96 3.73 96.36 92.76	0.90 9.49 62.79 40.63 56.04 93.61 31.58 96.65 49.62 84.65 95.60 0.00
Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74 R2L 97.98 72.37 119.32 37.76 72.37 Probes 97.13 86.19 17.44 83.17 86.19 Support Vector Machine (SVM) Normal 95.27 92.76 1.29 98.63 92.76 DoS 99.85 _ _ _ 0.00 _ U2R 98.35	93.61 9.49 9.49 62.79 40.63 56.04 93.61 31.58 96.65 96.65 49.62 84.65 95.60 0.00 97.64 38.16
Normal 47.95 60.87 13405.22 0.44 50.13 Normal 47.95 60.87 13405.22 0.45 60.87 DoS 99.04 5.51 10.66 34.09 5.51 U2R 58.88 46.29 1.15 97.58 46.29 R2L 97.19 45.83 79.81 36.48 45.83 Probes 86.30 40.35 3.64 91.73 40.35 Logistic Regression (LR) Normal 93.23 92.08 4.65 95.20 92.08 DoS 99.87 69.23 269.23 20.45 69.23 U2R 97.62 96.74 3.44 96.56 96.74 R2L 97.98 72.37 119.32 37.76 72.37 Probes 97.13 86.19 17.44 83.17 86.19 Support Vector Machine (SVM) Normal 95.27 92.76 1.29 98.63 92.76 DoS 99.85 _	93 93 93 93 93 93 93 94 95 95 95 95 95 95 95

Figure 7: Five classes of confusion matrices for KNN, RF, SVM and LR.

user traces, helps to build host profiles based on involved patterns using machine learning and reduces the false positive rate.

This step is divided into two main phases: training and evaluating.

Training and testing

The NSL-KDD dataset was used in a testing phase to see at each level that our attack signatures database is accurate in detecting attacks and to select which algorithm to use in labelling our logstructured matrix. In this step, we considered four machine learning algorithms

%	ACC	DR	FAR	Precision	Recall	F1	
K-Nearest Neighbours (KNN)							
Normal	98.77	98.64	1.00	99.00	98.75	98.87	
Attack	98.77	98.92	1.47	98.64	98.19	98.77	
Naïve Bayes (NB)							
Normal	87.24	85.64	8.89	90.59	85.64	88.05	
Attack	87.24	89.19	17.45	83.64	89.19	86.33	
Logistic Regression (LR)							
Normal	93.72	92.71	4.49	95.38	92.71	94.03	
Attack	93.72	94.87	8.34	91.92	94.87	93.37	
Support Vector Machine (SVM)							
Normal	95.19	92.95	1.74	98.16	92.95	95.49	
Attack	95.19	97.89	8.53	91.99	97.89	94.85	
Table 3: Comparison of performance metrics results for two-class detection.							

%	ACC	DR	FAR	Precision	Recall	F1	
K-Nearest Neighbours (KNN)							
Normal	96.69	94.89	0.57	99.40	94.89	97.09	
DoS	96.69	99.21	7.14	93.28	99.21	96.15	
Naïve Bayes (NB)							
Normal	45.36	53.19	304.26	14.88	53.19	23.26	
Attack	45.36	43.92	8.63	83.58	43.92	57.58	
Logistic Reg	Logistic Regression (LR)						
Normal	59.93	75.27	105.38	41.67	75.27	53.64	
Attack	59.93	53.11	11.00	82.84	53.11	64.72	
Support Vec	Support Vector Machine (SVM)						
Normal	79.14	88.89	35.56	71.43	88.89	79.21	
Attack	79.14	71.26	8.98	88.81	71.26	79.07	
Table 4: Performance metrics for audit log events							

for traffic classification and behaviour labelling: k-nearest neighbours, Naïve Bayes, support vector machine (SVM) and logistic regression (LR).²⁹ Classifications were made in two-class and five-class behaviours cases.

Figure 7 presents a confusion matrix for each applied algorithm and Table 2 shows the precise values of all traffic categories with KNN, Naïve Bayes, SVM and LR. For our experimental dataset, the classification model based on KNN produces higher precision for all five behavioural classes and is more stable.

Results showed (Table 2) that KNN is the most appropriate algorithm in classifying traffic network data, with 98.80% accuracy, 99.80% precision, 98.80% of recall and only a 0.9% false positive rate.

In order to identify the level at which the system maintains good accuracy, sensitivity and precision in detecting attacks when the class labels are reduced to only two (normal and attack categories), we made the same test on a two-class training and testing dataset.

Results in Table 3 showed that even with one attack class and normal behaviours, the KNN algorithm keeps a high accuracy and low false positive rate.

Evaluating and labelling log events

We extended the basic experimental setup adopted in many past researches and constructed a test data of audit logs. The test data is composed of 400 records of normal behaviours and 350 records of malicious behaviours. Results in Table 4 show that the number of false positives is reduced to 0.57% for normal behaviours and 7.14% for malicious behaviours. We may conclude that our system is accurate to 96.69% in labelling unknown behaviours.



Finally, we labelled our audit log files, containing 204,450 unknown behaviours. The system detected 29,333 DoS attacks, 40,420 probes, 27 R2L and 134,671 normal comportments.

Finding new attack signatures

Finding correlations between our training dataset features enables the system to discover new malicious behaviours and to identify new signatures when defining an attack. Furthermore, the rest of the unlabelled data from the formatted audit logs matrix can be labelled by substituting their highly correlated features. The new substituted matrix is communicated to the previous step to be analysed and labelled again, until the system converges.

In order to yield similarity scores between the different objects in our training data set, we have used several similarity metrics, each one depending on the prototype of the features.³⁰

The **Jaccard coefficient** is used when dealing with data objects having binary attributes. The objects are generally a disordered set of a collection of data. In our case, it is taken into consideration when seeing similarities between binary features such as logged_in, root_shell, is_host_login, is_guest_login.

$$J(X_i, X_j) = \frac{X_i \cap X_j}{X_i \cup X_j}$$

We have six binary features and from the results of the Jaccard similarity matrix we can presume that whether the user is a



guest logging in or not isn't important in defining a behaviour and may be deleted from our dataset, and that the location (land) of the connection, the status of the connection, the root shell privileges and the attempts to get the root user permissions are highly correlated to each other and can be reduced.

The **correlation metric** is investigated when seeking a relationship between two quantitative, continuous variables. It measures the strength of the association between the attack's features.

$$\rho_{X_i,X_j} = \frac{cov(X_i,X_j)}{\sigma_{X_i} \times \sigma_{X_j}}$$

 $cov(X_i, X_j)$ is the covariance measure. σ_{X_i} is the standard deviation measure.

From the correlation measures, we see that:

Rule 1: 'srv_count' is highly correlated to 'serror_rate', 'dst_host_srv_diff_ host_rate' and 'dst_host_serror_rate'. The number of connections to the same host and same service during a twosecond time window is more likely to activate the S0, S1, S2 or S3 flags. They are also proportional to connection to the same IP or the number of connections to different destinations (host) but the same port. **Rule 2:** 'serror_rate' is highly correlated to 'dst_host_srv_diff_host_rate' and 'dst_host_serror_rate'. The percentage of connections that were to the same destination and have activated one of the flags S0, S1, S2 or S3 during a twosecond time window is proportional to those that were to the same IP address and activated the same flags.

Rule 3: 'srv_serror_rate' is highly correlated to 'rerror_rate', 'dst_host_srv_serror_rate' and 'dst_host_rerror_rate'. The percentage of connections that have activated S0,S1, S2 or S3 flags among connections to the same service during a twosecond time window is proportional to the percentage of connections that have activated one of the same flags among the connections to the same port and related to the percentage of connections that have activated the REJ flag among connections to the same IP destination during a two-second time window.

Rule 4: 'rerror_rate' is highly correlated to 'dst_host_srv_serror_rate' and 'dst_host_ rerror_rate'. The percentage of connections that have activated the REJ flag among connections to the same service during a two-second time window is highly correlated to the percentage of connections that have activated one of the flags S0, S1, S2 or S3 among connections to the same port and related to the percentage of connections that have activated the REJ flag among connections to the same IP destination during a two-second time window.

Rule 5: 'dst_host_count' is highly correlated to 'dst_host_srv_count'. The number of connections to the same IP address during a two-second time window is highly related to the number of connections to the same port during the same period.

Rule 6: 'dst_host_srv_diff_host_rate' is highly correlated to 'dst_host_serror_rate'. The percentage of connections that were to the same destination and have activated one of the flags S0, S1, S2 or S3 is highly related to the percentage of connections to the same port but different machines during the same period.



Rule 7: 'dst_host_srv_serror_rate' is highly correlated to 'dst_host_rerror_ rate'. The percentage of connections that were to the same port and have activated one of the flags S0, S1, S2 or S3 is proportional to the percentage of connections that were to the same IP address and activated the REJ flag.

Cosine similarity is the usually used metric in the context of text mining while comparing documents, spams, log files and so on. In our case, it underlines the normalised product of each two features. When it is near to 1, it means that the two features have the same orientation and that one of the two should be more than enough to find the other one. Thus, the two attributes may be substituted.

$$cos(X_i, X_j) = \frac{X_i \times X_j}{\|X_i\| \times \|X_j\|}$$

From similarities metrics, we see that many features are correlated to each other and thus can be substituted to their highly correlated ones and drastically decrease the number of features following some rules:

Rule 1: Duration is highly correlated to protocol type, service and the percentage of connections to different hosts having the same service.

Rule 2: Protocol type is highly correlated to the percentage of connections to different hosts into the same port destination.

Rule 3: The number of connections to the same service in the same two-second window time is highly correlated to the number of connections to the same port outside the time window and to the percentage of connections that activated one of the flags S0, S1, S2 or S3 in the same period.

Rule 4: The number of connections to the same destination host that have activated the S0, S1, S2 or S3 flags in a two-second time window is mainly correlated to its number when it is to the same destination and overpasses the range time and the number of connections that were to different port destinations outside the time window.

Rule 5: The number of connections to the same service that have activated the S flags is proportional to the number of rejected connections to the same host destination



in the same time window and to the number of rejected connections in a different time window to the same host destination. It is also related to the number of connections to the same service that have activated the S flags outside the time window.

Rule 6: The percentage of rejected connections that were to the same host destination in the two seconds is highly related to its percentage when they are still to the same destination but outside the time window. It is also related to the number of connections to the same service that have activated the S flags outside the time window.

Rule 7: The number of connections that were to the same host destination is strongly related to the number of connections that were in the same service.

Rule 8: The percentage of connections that were to different destination machines, having the same service and outside the time window, is proportional to the percentage of connections that have activated the S flags when requesting the same destination host.

Rule 9: The percentage of connections outside the time window to the same service that have activated the S flags is proportional to the percentage of rejected connections to the same host machine that were outside the time window.

In the next step, the cloud auditor updates the tracking security system rules. It adds eligible signatures to the training dataset and evaluates constantly that high accuracy and precision are maintained in order to get very low false positive rates. Given that our model robustness is based on the AUC calculation too, which tends to be influenced by the imbalance of positive and negative data, we repeated our basic labelling experiment, adding normal and malicious audit log samples to the training dataset.

Those changes to the training/testing set showed obvious differences in the accuracy of the unknown behaviours identification, inferior results with the low update rate, and a continuous increasing of the system robustness through updates.

Finally, from AUC surfaces in Figure 12, results show that the more our training dataset is complete and accurate, the better are our prediction rates, accuracies and ROC surfaces.

Environment and results

The implementation was conducted on an Intel i5-7500 3.5GHz machine with 8GB of RAM. We used two separate datasets. The first consisted of the access logs of the Honeynet Project, composed of 204,449 log rows that we wanted to label. The second database was NSL-KDD, which is used as a knowledge database that validates the predictive system. Four-fifths of NSL-KDD was used to train the system and it took 1min 47s for the selected algorithm KNN. One-fifth of



the formatted access log and testing NSL-KDD dataset was used to validate the system. The labelling stage took less than 56 seconds to label 204,449 records. In streaming conditions, the labelling phase will take less than one second to flag one record or connection.

Our experiment has many objectives: enhancing the velocity of our system using big data techniques; selecting the best classifier over audit logs and NSL-KDD datasets; and finding a correlation between network features to finally detect attacks based on web application log files and our hybrid attack detection system. The NSL-KDD dataset was used as a training dataset and log files as a testing, evaluating and tracking database that helped us secure the system. It contained 61,685 signatures of normal behaviours, 42,706 rows of DoS attacks, 197 records of U2R attacks, 3,006 of R2L different signatures and 11,217 data rows of probes.

The audit logs used in our experiment are from the Honeynet Project, challenge scan 31 (SotM31). It traces 204,450 unknown behaviours that we aimed to label. As the fragments in audit logs were structured in the pre-processing step, the NSL-KDD and L matrices have four categories of features – basic (B) ones related to TCP connections, domain knowledge (D) containing both connection status and user attempts, information about a traffic network in a two-second window, and long-term attacks.

The experimental results demonstrated several standards. The first is the detection rate in Table 2. It exhibits that KNN is better than Naïve Bayes, random forest and support vector machine for both binary and multi-classes. As can be seen, KNN is 98.77% accurate in labelling behaviours, and has a very low false positive rate of around 1% to 1.47%. We also evaluated the robustness of the proposed architecture, applying a second test phase on the formatted audit logs matrix (L). The detection rate in Table 4 also shows that KNN is the best classifier, with an accuracy of 96.69% and an overall precision factor of 93.28% to over 99.4%. It also has a low score of false positives that varies from 0.57% to 7.14%.

The experiment also evaluated the importance of centralisation in increasing the network insight, selecting new attack signatures and finally updating the training dataset to enhance the system robustness.

In comparison with other published works – and especially Bambag and Supeno which also used hybrid machine learning (LSCANN) to detect intrusions and NSL-KDD as a validating database – their work achieved 99% accuracy and 27.64% FAR while ours surpasses it at 99.7% accuracy (in labelling the four above mentioned attacks) and 1.47% FAR rates. Furthermore, our proposed architecture not only detects intrusion at high effectiveness, but also selects new signatures of attacks based on feature correlations.

Conclusion

With the increase of connected devices in the cloud, the detection of anomalies in streaming data is becoming crucial. Our experimental results showed the successful adoption of our hybrid machine learning solution. It reduced error rates to 1.47%, enhanced the accuracy to 98.77% (for two-class labelling) and an average of 99.7% accuracy in labelling the four attacks. And it revealed its effectiveness in comparison with other machine learning algorithms that do not consider the importance of updating security rules (AUC = 73% when updating security rules).

Future work would focus more on how to predict the values of missing data and reduce the number of involved features in network traffic analysis using weight-of-evidence and informationvalue techniques. It could also evaluate outlier detection in defining the deviation probability at which an unknown behaviour is considered as a suspicious attempt. Studying these deviations is the motivation of our next search project.

About the authors

Amar Meryem is a business intelligence analyst and developer in the data preparation (ETL) and reporting fields. She is interested in the field of artificial intelligence and curious to find new approaches in detecting intrusions, attacks and malicious behaviours based on predictive models in cloud architectures. She is also a PhD student at Mohammed V University, Rabat in the Intelligent Processing and Security of Systems (IPSS) laboratory.

Bouabid EL Ouahidi has been a teacher and researcher at the Mohammed V University, Rabat since 1993 in the field of cyber security. He received his PhD degree from R&D of France Telecom and the University of Caen. He received a doctorate in state sciences in networks and distributed systems from both Institute Mines Telecom (Brest) and Mohammed V University.

References

- Ramachandra, G; Iftikhar, M; Khan, FA. 'A comprehensive survey on security in cloud computing'. Procedia Computer Science, pp.465-472, 2017.
- Lemoudden, M; Amar, M; El Ouahidi, B. 'A Binary-Based MapReduce Analysis for Cloud Logs'. Procedia Computer Science, vol.83, pp.1213-1218, 2016.
- Amar, M; Lemoudden, M; El Ouahidi, B. 'Log file's centralisation to improve cloud security'. 2nd International Conference on Cloud Computing Technologies and Applications

(CloudTech), IEEE, 2016, pp.178-183.

- Modi, C; Patel, D; Borisaniya, B; Patel, H; Patel, A; Rajarajan, M. 'A survey of intrusion detection techniques in cloud'. Journal of Network and Computer Applications, vol.36, no.11, pp.42-57, 2013.
- Lemoudden, M; El Ouahidi, B. <sup>(Managing cloud-generated logs using big data technologies^(A). International Conference on Wireless Networks and Mobile Communications (WINCOM), 2015, pp.1-7.
 </sup>
- Fageeri, SO; Ahmad, R. 'An efficient log file analysis algorithm using binary-based data structure'. Procedia-Social and Behavioral Sciences, vol.129, pp.518-526, 2014.
- Söderström, O; Moradian, E. 'Secure audit log management'. Procedia Computer Science, vol.22, pp.1249-1258, 2013.
- Muthurajkumar, S; Ganapathy, S; Vijayalakshmi, M; Kannan, A. 'Secured temporal log management techniques for cloud'. Procedia Computer Science, vol.46, pp.589-595, 2015.
- Sreeram, I; Vuppala, VPK. 'HTTP flood attack detection in application layer using machine learning metrics and bio inspired bat algorithm'. Applied computing and informatics, 2017.
- Salaken, SM; Khosravi, A; Nguyen, T; Nahavandi, S. 'Extreme learning machine based transfer learning algorithms: A survey'. Nerocomputing, vol.267, pp.516-524, 2017.
- Ahmad, S; Lavin, A; Purdy, S; Agha, Z. 'Unsupervised real-time anomaly detection for streaming data'. Neurocomputing, vol.262, pp.134-147, 2017.
- 12. Jie, Y; Xin, C; Xudong, X; Jianxiong, W. 'HIDS-DT: An Effective Hybrid Intrusion Detection System Based on Decision Tree'. IEEE Computer Society, no.1, International Conference

on Communications and Mobile Computing, 2010.

- Bashah, N; Idris, B; Abdul, MA. 'Hybrid intelligent intrusion detection system'. Engineering and Technology, World Academy of Science, vol.11, no.1, pp.23-26, 2005.
- Bambang, S; Supeno, D; Tohari, A; Nasrul, A. 'Assessing centroid-based classification models for intrusion detection system using composite indicators'. Procedia Computer Science, vol.161, no.1, Elsevier, pp.665-676, 2019.
- Rehman, Z-U; Khan, SN; Muhammad, K; Lee, JW; Lv, Z; Baik, SW; Shah, PA; Awan, K; Mehmood, I. 'Machine learning-assisted signature and heuristic-based detection of malwares in Android devices'. Computers & Electrical Engineering, vol.69, pp.828-841, 2018.
- 16. Li, Y; Ma, R; Jiao, R. 'A hybrid malicious code detection method based on deep learning'. International Journal of Security and Its Applications, vol.9, no.15, pp.205-216, 2015.
- Lekagning, BHD; Tindo, G. 'Vers une nouvelle architecture de detection d'intrusion reseaux a base de reseaux neuronaux'. 2016.
- Meryem, A; Samira, D; Bouabid, EO. 'Enhancing Cloud Security using advanced MapReduce k-means on log files'. Proceedings of the 2018 International Conference on Software Engineering and Information Management, pp.63-67, 2018.
- Ya, J; Liu, T; Li, Q; Shi, J; Zhang, H; Lv, P; Guo, L. 'Mining host behavior patterns from massive network and security logs'. Procedia Computer Science, vol.108, pp.38-47, 2017.
- Habeeb, RAA; Nasaruddin, F; Gani, A; Hashem, IAT; Ahmed, E; Imran, M. 'Real-time big data processing for anomaly detection: A survey'. International Journal of Information Management, 2018.

- Fan, Z; Liu, R. 'Investigation of machine learning based network traffic classification'. International Symposium on Wireless Communication Systems (ISWCS), pp.1-6, 2017.
- 22. Katzir, Z; Elovici, Y. 'Quantifying the resilience of machine learning classifiers used for cyber security'. Expert Systems with Applications, vol.92, pp.419-429, 2018.
- 23. Akash, G; Prachi, M. 'A hybrid intrusion detection system: A review'. IEEE, Conference on Intelligent Systems and Control, 2016.
- 24. Cathey, R; Ma, L; Goharian, N; Grossman, D. 'Misuse detection for information retrieval systems'. Proceedings of the twelfth international conference on Information and knowledge management, pp.183-190, 2003.
- 25. Dhanabal, L; Shantharajah, S. 'A study on NSL-KDD dataset for intrusion detection system based on classification algorithms'. International Journal of Advanced Research in Computer and Communication Engineering, vol.4, no.16, pp.446-452, 2015.
- 26. Capó, M; Pérez, A; Lozano, JA. 'An efficient approximation to the K-means clustering for massive data'. Knowledge-Based Systems, vol.117, pp.56-69, 2017.
- Jiang, Z; Li, T; Min, W; Qi, Z; Rao, Y. 'Fuzzy c-means clustering based on weights and gene expression programming'. Pattern Recognition Letters, vol.90, pp.1-7, 2017.
- 28. Meryem, A; Samira, D; Mouad, L; Bouabid, EO. 'A novel approach in detecting intrusions using NSL-KDD database and MapReduce programming'. Procedia Computer Science, vol.110, pp.230-235, 2017.
- Belavagi, MC; Muniyal, B.
 'Performance evaluation of supervised machine learning algorithms for intrusion detection'. Procedia Computer





A SUBSCRIPTION INCLUDES:

Online access for 5 users An archive of back issues

www.networksecuritynewsletter.com

The Firewall

Pandemic underpins need for SOAR



Colin Tankard, Digital Pathways

With the rapid change in working practices over the past few weeks, due to the Covid-19 outbreak, the need to get staff working from home has meant, in many cases, that the security measures normally in place when in the office have been ignored. This can be due to a lack of licences or simply that personal devices are being used that do not support corporate-level security systems.

The outcome will probably be loss of data, leading to subject access requests and the almost impossible task of finding data scattered across remote personal devices. Deciding if data has really been lost or if it is just a spoof ransom demand will require a smarter way of handling breach investigations, as time will be of the essence, given the Information Commissioner's Office rule of a 72-hour notification under the General Data Protection Regulation (GDPR).

Traditional approaches to the investigation of a breach would involve ploughing through firewall and server logs to identify any unusual behaviour. But if the data is outside the corporate network, this becomes difficult and so will take even longer to make decisions on data leakage.

What is required is a security orchestration, automation and response (SOAR) platform, which ingests the 'stolen' data and performs automated tests to see if the data is located in unusual places or has been modified. This process can reduce the total number of records under investigation to only those that match a change. These records can be further investigated to get a better picture of the breach. Imagine if you could get 10 million records down to 50,000 – how much easier would it be to work on the data?

The SOAR platform can be designed to interrogate many different sources using playbooks to access, for example, the dark web or third-party databases, such as mobile phone carriers.

One marker that can identify a change in data, and thus a possible fraudulent use, is a SIM swap. Here, a mobile number is transferred from one person to another. However, criminals can use this to steal an individual's money or gain access to corporate networks, as the use of two-factor authentication (2FA), using SMS texts, is increasing exponentially.

Someone stealing your phone number can result in severe problems, especially as most of our modern 2FA systems are based on SMS messages that can be intercepted using this technique. Criminals can hijack accounts, one by one, having a password reset sent to your phone. They can trick automated systems – such as banks – into thinking they're you. And they can use your hijacked number to break into your work email and documents.

Once a SIM swap has been identified, the platform checks other markers, such as the date of change and, if this falls within the expected timeline of the data theft, the number and its corresponding record can be pulled by the platform to enable further investigation, such as location of the SIM, where the SIM data was transferred to, if third parties were involved, etc. If these markers indicate a positive theft, the platform takes the results and feeds back to management to make decisions on the actions to follow.

In the current climate, speed of investigation is key. The only way to achieve this is to automate many of the processes and have access to multiple data sets to correlate against. This data can then be presented to investigators in a clear format, enabling better business decisions that are then fed back into the SOAR platform, to automate even further.

EVENTS CALENDAR

Due to the Covid-19 pandemic, many conferences are being cancelled, postponed or converted into virtual events. The events listed here were still planned to proceed at the time of publication.

1 October 2020 ArcticCon Anchorage, Alaska

https://arctic-con.com

6–8 October 2020 Critical Infrastructure Protection & Resilience Europe

Bucharest, Romania www.cipre-expo.com

7–8 October 2020 International Cybersecurity & Intelligence Conference (ICSIC)

Toronto, Canada https://www.icsicanada.org

8 October 2020 Florida Cyber Conference Orlando, Florida, US https://flcybercon.com

14–16 October 2020 International Conference on Digital Forensics & Cyber Crime (ICDF2C) Boston, US

http://icdf2c.org

20–21 October 2020 600Minutes Information and Cyber Security Stockholm, Sweden https://bit.ly/3bNwTcr

26–29 October 2020 International Conference on Internet of Things: Systems, Management and Security Milan, Italy http://emergingtechnet.org/IOTSMS2020/