

Case Report

Not peer-reviewed version

Instrumenting OpenCTI with a Capability for Attack Attribution Support

[Sami Ruohonen](#) ^{*}, [Alexey Kirichenko](#), Dmitriy Komashinskiy, Mariam Pogosova

Posted Date: 29 August 2023

doi: 10.20944/preprints202308.1936.v1

Keywords: cyberattack; technical cyberattack attribution; digital forensics; machine learning; cyber threat intelligence



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Case Report

Instrumenting OpenCTI with a Capability for Attack Attribution Support

Sami Ruohonen ^{1,*}, Alexey Kirichenko ^{1,2}, Dmitriy Komashinskiy ¹ and Mariam Pogosova ¹

¹ WithSecure Corporation, Tammasaarenkatu 7, 00180 Helsinki, Finland

² Independent consultant, Aallonhuippu 10 43, 02320 Espoo, Finland

* Correspondence: sami.ruohonen@withsecure.com

Abstract: In addition to identifying and prosecuting cyber attackers, attack attribution activities can provide valuable information guiding the defenders' security procedures and giving them greater confidence in incident response and remediation. However, technical analysis involved in cyberattack attribution requires high skills, experience, access to up-to-date Cyber Threat Intelligence, and significant investigators' effort. Attribution results are not always reliable, and skilful attackers often work hard to cover their traces and mislead or confuse investigators. In this article, we present a tool designed to support technical attack attribution and implemented as a machine learning model extending the OpenCTI platform. We also discuss the tool's performance in the investigation of a recent cyberattack.

Keywords: cyberattack; technical cyberattack attribution; digital forensics; machine learning; cyber threat intelligence

1. Introduction

Law Enforcement Agencies (LEAs), forensic institutes, national cybersecurity centres and Computer Emergency Response Teams (CERTs), and companies providing cybersecurity services routinely have to investigate cyberattacks on organisations and citizens. In many cases, a key question in such investigations is who is responsible for conducting a given cyberattack. This identification of the source of a cyberattack – which can be a nation state, a crime syndicate, other nefarious group, or even an individual cybercriminal – is often referred to as 'cyberattack attribution'. In this article, the focus is on *technical* attack attribution, which is based on the analysis of technical attack traces and Cyber Threat Intelligence (CTI). While it was pointed in [1] that "... questions of responsibility are rarely decided solely through a single technological tool or form of evidence ..." [1] (p. 382) and "... a legal approach, rather than a technological one, can solve the attribution problem." [1] (p. 376), technical attribution is nearly always an indispensable element of any attribution efforts, providing key facts and hypotheses.

Knowing the threat actor behind a cyberattack can be very important and valuable, though the attribution value and investigation priorities vary and depend significantly on the context. For internal cybersecurity teams, CERTs and commercial service providers, attribution efforts usually help understand the attacker's intentions, capabilities and level of sophistication, modi operandi, and expected behaviour, informing the defenders' security procedures from prevention to response and remediation and giving them greater confidence. For example, the understanding of the attacker's tactics, techniques, and procedures (TTPs) guides the defenders in what additional attack traces and artefacts they should look for and what vulnerabilities they have to prioritise for minimising the impact of the ongoing attack and the risk of future ones. In the context of cyberattacks driven by political, military or industrial competition reasons, the attribution (e.g., to a nation state) value can include a reliable view of the impact of sensitive information loss and can extend to driving foreign policy measures. Also, importantly for LEAs, the insights brought by attribution efforts can be instrumental in identifying and prosecuting attackers.

With all the potential benefits, technical analysis involved in cyberattack attribution requires high skills, experience, access to up-to-date CTI, and significant investigators' effort. Furthermore,



attribution results are not always reliable, and skilful attackers often work hard to cover their traces and mislead or confuse investigators. Recognising the challenges, the EU-funded CC-DRIVER [2] and CYBERSPACE [3] projects contributed to designing and developing a tool supporting cyberattack attribution. This article presents the tool and discusses the results of its application in the investigation of a recent cyberattack. We first briefly review several noteworthy challenges of technical attack attribution, the data used in attack analysis, the connections between attribution and other key questions that arise in digital forensics and cyber incident response activities, and the earlier work on applying machine learning to the attack attribution problem. We then explain the technical approach, present the tool, based on a machine learning model and implemented as an extension of the OpenCTI platform [4], and show its performance in the 'No Pineapple!' cyberattack investigation carried out by one of the CC-DRIVER and CYBERSPACE partners – WithSecure Corporation. The article is concluded by discussing the challenges and directions for future work.

2. Technical Attack Attribution

When running analysis in order to identify the source of a cyberattack, investigators face multiple problems. Cybercriminals and other perpetrators usually hide the origin of their attack network traffic by routing it via multiple links on the Internet, for instance, using proxy servers or onion-routing tools such as Tor [5] instead of directly connecting to the victim (which is essentially enabled by the structural design of the Internet), or by using compromised or stolen devices of other people (since identifying the source devices of an attack is not the same as identifying the people behind it). They also increasingly often rely on tools commonly available on the victim devices instead of using custom malware that can be fingerprinted and associated with their authors – the technique known as "living-off-the-land" [6]. Attribution activities are further complicated by the growing popularity of the "Crime-as-a-service" mode of cybercriminal operations (malware-as-a-service, ransomware-as-a-service, DDoS-as-a-service, bulletproof hosting, etc.), the use of malicious code which is open-sourced, shared or stolen from other attackers (and sometimes even from state security agencies and security researchers [7,8]), and the use of malicious infrastructure (such as command-and-control servers) and other TTPs previously attributed to other attackers. One should also note that CTI and other information crucial for attack attribution can be kept confidential by certain parties due to laws, contracts, and various – justified or unjustified – concerns.

Attack attribution is closely connected with several other questions asked typically by incident responders and investigators when trying to gain visibility into the threat actor's operations in the victim's cyber estate. Good examples of such questions are:

- When did the threat actor breach the victim's systems and networks?
- What level of privilege does the threat actor have at the moment?
- What assets has the threat actor touched and potentially compromised?
- What is the impact caused by the breach?

So, essentially any data collected in an incident response operation can be useful for attribution-related analysis, while the data revealing the threat actor's capabilities, objectives and behaviour is of particular value. This includes:

- Attacker's TTPs. The MITRE ATT&CK framework [9] is commonly used to structure and model this information.
- Indicators of Compromise (IOCs) and attacker's infrastructure, such as the file hashes of malicious payloads and IP-addresses which the attack traffic originates from or where the command-and-control (C2) servers are hosted.
- Malware analysis results (especially for victim-tailored malware with no public source code), which can provide high-value information. For instance, sometimes attackers make mistakes or leave traces in their malware code, and in other cases, they use evolving versions of the same malware for many years.
- Benign tools used by the attacker. These can be popular living-off-the-land binaries, such as Powershell and Windows Management Instrumentation (WMI), or other benign software found in the victim's estate and providing capabilities beneficial for the attacker.

- Exploited vulnerabilities, either previously unknown ones (zero-day vulnerabilities) or used earlier in other attacks. Exploitation techniques can be implemented in malware or by using appropriate benign tools, and we often see the same vulnerabilities used in multiple attacks conducted by the same threat actor.
- Attack metadata, such as the times when the attacker communicates with the victim's systems (which can hint at the attacker's geographical location) or information about the victim (as their operations, core business domains, location, etc. can reveal the attacker's objectives).

Given the nature of the attack attribution problem, an obvious approach is to look for similarities in the data collected from attacks and about attackers (presumably structured and stored in a convenient form). Identifying, ranking and aggregating such similarities in large volumes of highly heterogeneous data is, however, time-consuming for investigators and requires expertise and experience. So, the growing number and sophistication of cyberattacks prompt analysis automation, and machine learning techniques come here as a natural choice.

While machine learning has recently been very popular in attack detection and malware analysis methods, it seems very few reports are available on its applications to cyberattack attribution.

Han et al. [10] implemented WHAP, a web-hacking profiling system that uses a simple similarity measure for hacking cases, which is based on heuristically assigned similarity weights for selected features (such as IP addresses and domain names) and Case-Based Reasoning. While the use of feature vectors for representing website hacking cases and the defined similarity measure for those vectors are the only connections of the proposed approach to machine learning, conceptually it can be extended to attack attribution methods utilising similarity search and clustering based on "learning from data".

Noever et al. [11] presented a Random Forest classifier for attributing attack techniques (such as backdoor, man-in-the-middle, ransomware, DoS) to the types of threat actors (organised crime, nation-state, hacktivist, unknown). While this approach can be relevant, e.g., for policy discussions, it does not have the attribution of specific cyberattacks as its objective.

Noor et al. [12] presented a framework for attributing unstructured (natural language) CTI reports and documents. Since "low-level Indicators of Compromise (IOCs) are rarely re-used and can be easily modified and disguised resulting in a deceptive and biased cyber threat attribution" [12] (p. 227), the work focuses on common high-level attack patterns (i.e., TTPs) for mapping a CTI report to a threat actor. With the labels for high-level attack patterns taken from the MITRE ATT&CK taxonomy, Latent Semantic Analysis (LSA) is used to index CTI reports with relevant labels. Then a small set of CTI reports collected from publicly available datasets and marked with the threat actors behind the reported cyberattacks is used to train several machine learning models for attributing new reports. Although some of the models showed a very good performance in the tests, this is likely explained by the training and validation dataset's toy size. More generally, we think that fully focusing on high-level attack patterns and ignoring low-level indicators will result in poor real-world performance because: (i) many attackers use very similar TTPs (e.g., in ransomware attacks); (ii) high-level TTPs are easy to mimic in false flag operations; (iii) low-level indicators are actually re-used (mainly due to attacker's mistakes or time and cost pressures on their side) and very useful in such cases. We will further comment on this high-level vs. low-level balance issue in the "Discussion and Future Work" section.

The use of pattern recognition and anomaly detection methods for TTP and IOC extraction from raw log data was also proposed by Landauer et al. in [13], illustrated by system log data analysis.

3. STIX-based Attack Attribution Approach

A key technical objective defined in both CC-DRIVER and CYBERSPACE is to produce tools for following the threat landscape and actors (CTI management) and for investigating cyberattacks (digital forensics), and this toolkit serves a natural foundation for adding attack attribution capabilities. As fully automated reliable attack attribution is hardly feasible, we chose to build an attack attribution recommender, based on the Structured Threat Information Expression [14] (STIX

2) language and implemented as an OpenCTI extension, aiming to guide incident investigators and significantly reduce their efforts.

To facilitate the process of identifying threat actors responsible for cyberattacks, the problem was framed as follows: Design and implement a machine learning model that takes *a bundle of STIX 2 objects representing adversarial operations* as input and predicts "*the most probable*" threat actors behind the operations.

'Bundle' here is a STIX 2 term that refers to a collection of STIX 2 objects. While in principle any STIX 2 entities can be included in a bundle, we started with an important special case when a bundle is a set of incidents which were observed in a given (attacked) organisation in a given timeframe. In STIX 2, such 'incidents' represent information collected during attack investigation activities (conducted usually by law enforcement, CERTs or companies providing incident response services).

Threat actors are typically understood as identities representing an individual, a group or an organisation which operates in cyberspace with a malicious intent. We, however, chose to build our recommender model to predict *intrusion sets* (which can subsequently be mapped to identities, e.g., by LEAs) in order to provide greater flexibility. Cyberattacks are often leveraged by threat actors, e.g., a nation state or a cybercriminal group, as part of a coordinated campaign against a specific target and contain similar properties, behaviours and attributes in order to achieve multiple objectives over a significant period of time. Such an entire attack package is represented in STIX 2 as an intrusion set, and there are advantages in reasoning about attribution in terms of intrusion sets. For example, the threat actor behind a given attack may not be known but their multiple operations can be grouped together, in an intrusion set, and then a new attack can be attributed to that intrusion set. A threat actor can move from one intrusion set to another, changing their TTPs, or they can "utilise" multiple intrusion sets at the same time. Attribution relationships in STIX 2 are shown in Figure 1.

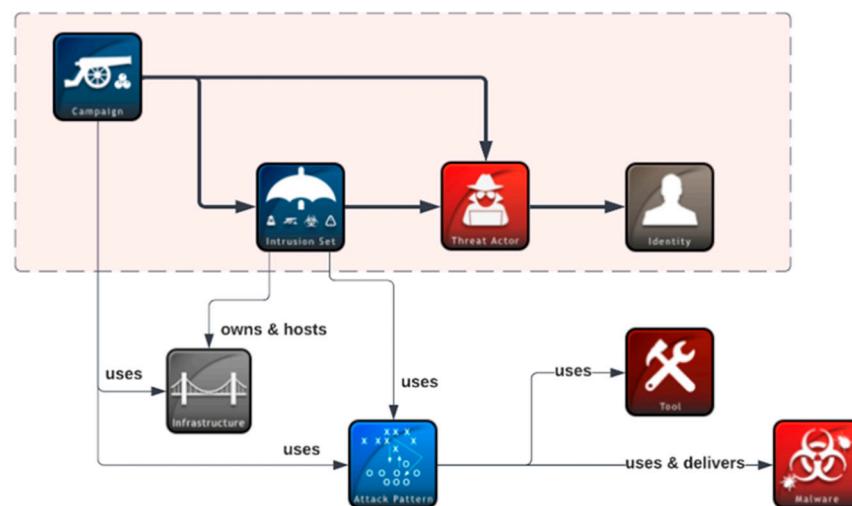


Figure 1. Attributed-to relationships in STIX 2 are shown with the arrows in bold. (the icons are taken from <https://github.com/MISP/intelligence-icons>).

Our preliminary investigations confirmed that obtaining sufficiently large incident datasets for training a good attribution support model would be challenging, particularly because such datasets are often considered highly confidential. So, in parallel to extending the data collection, we simplified our problem to predicting an intrusion set for a single incident (instead of a bundle), as shown in Figure 2.



Figure 2. The simplified version of the attribution problem. (the icons are taken from <https://github.com/MISP/intelligence-icons>).

For intrusion sets, a decent collection of over 350 entities (identified by their names, such as 'APT28' or 'Lazarus') was obtained from MITRE [15], AlienVault [16], Malpedia [17], and WithSecure. To compensate for the shortage of incident data, we chose to rely on the data augmentation approach, generating a number of synthetic incidents based on the data from the available intrusion sets. In producing incidents, specific rules designed together with cybersecurity experts, guided by their experience and observations, were followed:

- Incident data re-uses the elements (attributes) present in a specific intrusion set.
- The number of elements should be between 10 and 50 per incident, following a beta-binomial distribution with the median value around 15.
- From the attributes present in an intrusion set, the following STIX 2 objects are re-used: TTPs (up to 50%), tools (up to 20%), malware (up to 20%) and others (up to 10%). 'Others' here include indicators, locations and so on (all the entities that can be found in the intrusion set). The numbers in brackets indicate the upper bounds on the share of re-used attributes of a given type. However, if, for example, an intrusion set does not have 'tool' attributes at all, we will end up having zero tools added to synthetic incidents. Having the upper bounds, actual numbers of attributes of a given type are selected uniformly at random.
- To keep the synthetic dataset balanced, each non-empty intrusion set is used to generate the same number of incidents.

Using this approach, hundreds of thousands of synthetic incidents can be produced from the available intrusion sets, and those form the main body of a labeled dataset for supervised learning. It is split into training and testing sets, where the testing (validation) set has 20% of the data and the rest is used for training a model. CountVectorizer is applied as a one-hot encoder: the entity IDs and names seen in all the incidents are used as our features (names of malware families and tools, MITRE IDs of TTPs, etc.), and for a given incident the value of a specific feature is 1 if the corresponding entity is present in the incident and 0 otherwise.

After multiple rounds of experiments with several multi-class classification models (with the same collection of intrusion sets but new sets of synthetic incidents produced for each round), we selected for the recommender model the Bernoulli Naïve Bayes classifier, which showed good results with a low variation over the testing rounds. Of course, the good observed performance can be due to the synthetic nature of the data, so we are collecting more real-world incident data and planning further extensive modeling and validation experiments.

4. Attribution Results for 'No Pineapple!' Incident

The attack attribution recommender runs as an OpenCTI extension, and the OpenCTI platform, with a growing user community and a convenient framework for extending the platform's capabilities, has become a popular choice for storing, analysing and sharing both CTI and fresh digital forensics data from ongoing cyber incident investigations. STIX 2, the underlying OpenCTI data format, allows for a rich representation of incidents as collections of associated entities and observables (such as TTPs, malware, command-and-control infrastructure), combining high-level, abstracted views of attacks with relevant technical details. This data expression power explains why increasingly many incident response operations by WithSecure, a major European provider of cybersecurity services and solutions, rely on OpenCTI for data management and analysis, which

recently gave us an opportunity to validate the recommender as part of a real-world attack investigation engagement.

The attack, which was codenamed 'No Pineapple!' by the WithSecure's Threat Intelligence team due to one error message found in the malicious code, turned out to be part of a sophisticated campaign targeting public and private sector research organisations, the medical research and energy sectors as well as their supply chains. The WithSecure's engagement started when a threat hunt in a customer estate identified beaconing [18] to a Cobalt Strike C2 server. Since the C2 server IP was earlier listed as an IOC for the BianLian ransomware group and some other details also pointed in that direction, the initial (low confidence) assessment of the WithSecure's experts was that they were dealing with a potential ransomware incident. However, as more attacker tools, techniques and actions were collected from the customer environment, it became evident that the main objective of the attack was espionage, and a North Korean state-sponsored threat actor was behind it. Notably, the attacker took a serious effort of hiding their traces, clearing logs and deleting files, tools and other indicators of their presence [19].

The collected digital forensics data were added to OpenCTI as an incident object representing the details of a single attack against a single organisation. The object has quite a rich set of relationships, as can be seen in Figure 3.



Figure 3. Relationships of the 'No Pineapple!' incident as seen in OpenCTI.

We then applied the attack attribution recommender tool to the 'No Pineapple!' data and received the 'Lazarus' intrusion set associated with a North Korean state-sponsored threat actor on the top of the list. It should be noted that at the time of the recommender validation experiment, the Lazarus intrusion set was not updated with the 'No Pineapple!' investigation data but represented the state of knowledge prior to the investigation.

The top three results reported by the tool (with the respective model confidence values) were:

1. Lazarus Group: 0.996186486423268
2. Elephant Beetle: 0.003794891776652858
3. APT29: 0.000018620678059799746

So, the Lazarus group was suggested by the model as the most probable intrusion set for 'No Pineapple!' with an overwhelming confidence, and this was fully confirmed by the WithSecure's experts. Elephant Beetle, which is a financially motivated cybercrime group, was the second model's pick. While the model confidence for the Elephant Beetle intrusion set is low, we note that it shares a set of common attack techniques with Lazarus, including: blending in with the environment; deploying JSP web shells (JSP file browser, in particular); operating out of temp directories. It also exploits known vulnerabilities in public facing devices to gain initial access, although we do not know any vulnerabilities exploited by both Lazarus and Elephant Beetle. That is where the similarities end. Elephant Beetle is known to target different geographies, their operations have been financially motivated and they often target web services and their components.

5. Discussion and Future Work

The results obtained so far indicate that the approach of building machine learning models for attributing STIX 2 incidents to intrusion sets is promising and can bring significant value to incident investigators. The reliability of such models, especially when attackers actively work to counter attribution efforts through the use of false flags and other techniques, critically depends on the

availability of sufficiently rich incident data and on finding in the incident representation a suitable balance between high-level attack patterns and attributes and low-level indicators and other details. While STIX 2 is good for expressing TTPs, malware, tools, exploited vulnerabilities, targeted geography and sectors at certain level, more subtle details – such as malware code similarities, custom passwords, developer host information, attacker's email language, malicious domain registrar and registrant information – are not supported yet. For example, malware binaries and resource files are identified by their hash values, so even a very high similarity of two different files is of no use for our models, no matter how important it could be for attribution.

We see several ideas to explore for improving the attack attribution recommender:

- Acquiring more real-world incident data, preferably with attribution labels (but even unlabeled incidents can be useful), instead of heavily relying on synthetically generated incidents.
- If many organisations agree to combine their incident data, a high-quality attribution model can likely be trained, but incident data is usually highly sensitive. One way to address the data confidentiality issue is to train a model in a federated learning manner [20] on data of multiple organisations. In particular, joint efforts with the FATE project [21] working on collaborative confidentiality-preserving learning on CTI data can be considered.
- Use of inherited STIX 2 relationships (through the OpenCTI rule engine). At the moment, only the data of direct neighbours, i.e., first-level relationships, is used in the model for both incidents and intrusion sets. For example, a file associated with an incident may have another relationship with a custom directory where this file was located. If the same directory is associated with other files, this information may be valuable for attribution but is currently ignored.
- STIX 2 supports timestamps which can be used for building a timeline of attacker's actions. Because most of the incident data in our model training sets is produced synthetically from the intrusion sets, timestamps are currently ignored. Collecting timestamps whenever possible and including them in modeling should be explored for utilising attack timelines in attribution.
- Controlling the weights of features in the incident representation. For example, while many attackers use similar attack tactics or can easily imitate the tactics used by others, the presence of specific files can be a more reliable indicator for attribution. At the moment, the influence of specific features is learned implicitly when the model is trained. Combining the data-driven approach with expert-defined rules could be explored.

In conclusion, we would like to emphasize that even when large and clean training datasets are available, attack attribution models will make mistakes and can be deceived by skilful and determined attackers. Therefore, such models should primarily be used in the recommendation mode, with human experts verifying their output.

Funding: Parts of this research were supported by the CC-DRIVER project funding received from the European Union's Horizon 2020 research and innovation programme under grant agreement No 883543 and by the CYBERSPACE project funding received from the European Union's Internal Security Fund – Police (ISFP) programme under grant agreement No 101038738. The APC was funded by CYBERSPACE.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: The STIX 2 data used for producing the recommender model is aggregated from the following sources:

- MITRE database. The data is stored in GitHub [15] and publicly available. Third-party restrictions apply to the use of the data (used under license for the presented research).
- AlienVault. The data is publicly available in the AlienVault service [16]. Third-party restrictions apply to the use of the data (used under license for the presented research). A free account is required to access the data.
- Malpedia. The data is publicly available in the Malpedia website [17]. Third-party restrictions apply to the use of the data (used under license for the presented research). A free account is required to access the data.

All these data sources have open-source connectors for OpenCTI to facilitate data ingestion. The WithSecure data used for the presented research is not publicly available due to legal restrictions but constitutes only a small portion of the overall dataset.

Acknowledgments: The authors would like to thank their WithSecure colleagues from the Threat Intelligence and Incident Response teams for the help and valuable discussions. We also express our gratitude to the CC-DRIVER and CYBERSPACE project partners for their support.

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the presented research; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results.

References

1. Tran, Delbert. "The law of attribution: Rules for attribution the source of a cyber-attack." *Yale JL & Tech.* 20 (2018): 376-441. Available online: https://yjolt.org/sites/default/files/20_yale_j. l. tech. 376.pdf (accessed on 5 June 2023).
2. The EU-funded CC-DRIVER project: <https://www.ccdriver-h2020.com/> (accessed on 11 June 2023).
3. The EU-funded CYBERSPACE project: <https://cyberspaceproject.eu/> (accessed on 11 June 2023).
4. Open Cyber Threat Intelligence Platform – an open source platform for managing cyber threat intelligence knowledge and observables: <https://www.filigran.io/en/solutions/products/opentci/> (accessed on 11 June 2023).
5. Tor project: <https://www.torproject.org/> (accessed on 10 June 2023).
6. Living off the Land: Attackers Leverage Legitimate Tools for Malicious Ends. Available online: <https://symantec-enterprise-blogs.security.com/blogs/threat-intelligence/living-land-legitimate-tools-malicious> (accessed on 8 June 2023).
7. Stolen NSA hacking tools were used in the wild 14 months before Shadow Brokers leak. Available online: <https://arstechnica.com/information-technology/2019/05/stolen-nsa-hacking-tools-were-used-in-the-wild-14-months-before-shadow-brokers-leak/> (accessed on 8 June 2023).
8. Cobalt Strike, a penetration testing tool abused by criminals. Available online: <https://www.malwarebytes.com/blog/news/2021/06/cobalt-strike-a-penetration-testing-tool-popular-among-criminals> (accessed on 8 June 2023).
9. MITRE ATT&CK framework: <https://attack.mitre.org/> (accessed on 10 June 2023).
10. Han, Mee Lan, et al. "WHAP: Web-hacking profiling using case-based reasoning." 2016 IEEE Conference on Communications and Network Security (CNS). IEEE, 2016. DOI: 10.1109/CNS.2016.7860503
11. Noever, D. & Kinnaird, David. (2016). Identifying the Perpetrator: Attribution of Cyber-attacks based on the Integrated Crisis Early Warning System and the VERIS Community Database. 2016 International Conference on Social Computing, Behavioral-Cultural Modeling & Prediction and Behavior Representation in Modeling and Simulation. At: Washington DC. Available online: http://sbp-brims.org/2016/proceedings/CP_136.pdf (accessed on 5 June 2023).
12. Noor, Umara, et al. "A machine learning-based FinTech cyber threat attribution framework using high-level indicators of compromise." *Future Generation Computer Systems* 96 (2019): 227-242. DOI: 10.1016/j.future.2019.02.013
13. Landauer, Max, et al. "A framework for cyber threat intelligence extraction from raw log data." 2019 IEEE International Conference on Big Data (Big Data). IEEE, 2019. DOI: 10.1109/BigData47090.2019.9006328
14. STIX Version 2.1 OASIS Standard. Available online: <https://docs.oasis-open.org/cti/stix/v2.1/stix-v2.1.html> (accessed on 10 June 2023).
15. MITRE data: <https://raw.githubusercontent.com/mitre-attack/attack-stix-data/master/enterprise-attack/enterprise-attack.json> (accessed on 15 June 2023).
16. AlienVault, The World's First Truly Open Threat Intelligence Community: <https://otx.alienvault.com/> (accessed on 15 June 2023).
17. Malpedia: <https://malpedia.caad.fkie.fraunhofer.de/> (accessed on 15 June 2023).
18. Purple Team: About Beacons. Available online: <https://www.criticalinsight.com/resources/news/article/purple-team-about-beacons> (accessed on 12 June 2023).
19. No Pineapple! – DPRK Targeting of Medical Research and Technology Sector. Available online: <https://labs.withsecure.com/publications/no-pineapple-dprk-targeting-of-medical-research-and-technology-sector> (accessed on 12 June 2023).
20. Federated Learning: Collaborative Machine Learning without Centralized Training Data. Available online: <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html> (accessed on 12 June 2023).
21. FATE (Federated AI Technology Enabler) project: <https://github.com/FederatedAI/FATE> (accessed on 12 June 2023).

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s)

disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.