

iGen: Toward Automatic Generation and
Analysis of Indicators of Compromise (IOCs)
using Convolutional Neural Network

by

Anupam Panwar

A Thesis Presented in Partial Fulfillment
of the Requirements for the Degree
Master of Science

Approved April 2017 by the
Graduate Supervisory Committee:

Gail-Joon Ahn, Chair
Adam Doupé
Ziming Zhao

ARIZONA STATE UNIVERSITY

May 2017

ABSTRACT

Field of cyber threats is evolving rapidly and every day multitude of new information about malware and Advanced Persistent Threats (APTs) is generated in the form of malware reports, blog articles, forum posts, etc. However, current Threat Intelligence (TI) systems have several limitations. First, most of the TI systems examine and interpret data manually with the help of analysts. Second, some of them generate Indicators of Compromise (IOCs) directly using regular expressions without understanding the contextual meaning of those IOCs from the data sources which allows the tools to include lot of false positives. Third, lot of TI systems consider either one or two data sources for the generation of IOCs, and misses some of the most valuable IOCs from other data sources.

To overcome these limitations, we propose iGen, a novel approach to fully automate the process of IOC generation and analysis. Proposed approach is based on the idea that our model can understand English texts like human beings, and extract the IOCs from the different data sources intelligently. Identification of the IOCs is done on the basis of the syntax and semantics of the sentence as well as context words (e.g., “attacked”, “suspicious”) present in the sentence which helps the approach work on any kind of data source. Our proposed technique, first removes the words with no contextual meaning like stop words and punctuations etc. Then using the rest of the words in the sentence and output label (IOC or non-IOC sentence), our model intelligently learn to classify sentences into IOC and non-IOC sentences. Once IOC sentences are identified using this learned Convolutional Neural Network (CNN) based approach, next step is to identify the IOC tokens (like domains, IP, URL) in the sentences. This CNN based classification model helps in removing false positives (like IPs which are not malicious). Afterwards, IOCs extracted from different data sources are correlated to find the links between thousands of apparently unrelated

attack instances, particularly infrastructures shared between them. Our approach fully automates the process of IOC generation from gathering data from different sources to creating rules (e.g. OpenIOC, snort rules, STIX rules) for deployment on the security infrastructure.

iGen has collected around 400K IOCs till now with a precision of 95%, better than any state-of-the-art method.

To my loving family and friends for their patience and support. I want to give special note of appreciation to my late mother Madhu Panwar for her sacrifice and hard work which she put in to teach me the importance of hardwork and being good human. She always wanted to see me succeed. She sacrificed a lot of her personal interest and ambitions to provide all the necessary resources needed for my success.

Mom, this is for you! My twin brother Anurag motivated me to do research and taught me never give up attitude. My father Dr. H R Panwar always supported my decision. My elder brother Anshuman worked as catalyst for my success. Finally, I am forever indebted to my family for their understanding, endless patience and encouragement when it was most required.

ACKNOWLEDGMENTS

My journey through my Master of Science (Computer Science) degree has been an exciting one which has played a major role in my career and life, especially in enhancing my knowledge and experience in this field of study. Having had the opportunity to work with Dr. Gail-Joon Ahn for two years has been an insightful experience, and I am utmost grateful to him for having given me the opportunity to work on cutting edge research projects at Center for Cybersecurity and Digital Forensics (CDF). Dr. Ahn has been a constant source of motivation through my graduate career at Arizona State University, and has contributed significantly towards my abilities in reasoning, approaching and solving research problems impacting the society as a whole. I owe him a lot for having given directions to my thoughts and unlimited encouragement. Without his insight and patient guidance this work would not have been possible. I would like to thank Dr. Adam Doupé and Dr. Ziming Zhao for serving on my committee and providing their valuable feedback on my thesis.

My experience at the CDF has provided me with great opportunities in my professional career, and has enabled me to connect with brilliant and innovative individuals who have always been there to provide valuable inputs during my research work. I thank my fellow labmates in Threat Intelligence Analytics (TIA) Project: Zhibo (Eric) Sun, Faris Bugra Kokulu, Tejas Khairnar, and Ajay Modi, for the stimulating discussions, for the sleepless nights we were working together before deadlines, and for all the fun we have had in the last two years.

TABLE OF CONTENTS

	Page
LIST OF TABLES	vii
LIST OF FIGURES	viii
CHAPTER	
1. INTRODUCTION	1
2. BACKGROUND	8
2.1. Threat Intelligence	8
2.1.1. Consumption of Threat Intelligence	9
2.2. Convolutional Neural Network Architecture	10
2.2.1. Convolution	10
2.2.2. Convolution Neural Network	11
3. DATA COLLECTION	12
4. SYSTEM DESIGN	14
4.0.1. Architecture Overview	15
4.1. Data Acquisition Module	16
4.1.1. Blog Scrapper (BS)	16
4.1.2. APT Report Collector (ARC)	16
4.1.3. Cuckoo Malware Analysis System (CMAS)	17
4.2. IOC Extraction Module	17
4.2.1. Relevant Content Picker (RCP)	18
4.2.2. Sentence Classification (SC)	19
4.2.3. IOC Extractor (IE)	22
4.3. Rule Generation System	22
4.3.1. IOC Generator (IG)	23
4.3.2. IOC Formatter (IF)	23

CHAPTER	Page
5. IMPLEMENTATION	24
5.1. CNN based Sentence Classification	24
6. EVALUATION	28
6.1. Experiments	28
6.2. Results	30
6.2.1. Loss Function and Accuracy	30
6.2.2. Method Comparison	31
7. DISCUSSION AND FUTURE WORK	33
7.0.1. Limitations	33
8. RELATED WORK	35
8.1. Knowledge discovery	35
8.2. Convolutional Neural Network	36
8.3. Sharing Threat Intelligence	36
9. CONCLUSION	38
REFERENCES	39
APPENDIX	
A. LIST OF LINKS TO WEBSITES	44

LIST OF TABLES

Table	Page
1. Summary of Dataset	12
2. Sample Regular Expressions	19
3. Change in Accuracy with Embedding Size.	28
4. Change in Accuracy with Number of Filters.	29
5. Comparative Analysis with Other Methods	32

LIST OF FIGURES

Figure	Page
1. IOC Lifecycle	2
2. iGen Architecture Overview.	14
3. Architecture of iGen CNN based Sentence Classification Model.	25
4. Change in Loss Function over Steps.	30
5. Change in Accuracy over Steps.	31

Chapter 1

INTRODUCTION

According to Verizon's 2016 Data Breach Investigations Report [20], over 100,000 security incidents were reported in 2016 across 82 countries, which is a 25% increase over the prior year. Because the number of security threats and breaches is rapidly increasing over time, every organization is attempting to protect their systems and their data. The threat landscape is always progressing, and the information security risk is increasing because of the organization's dependence on computing systems. This constantly shifting and constantly increasing number of threats results in tremendous pressure on organizations to manage threats. Though abundant information is available in the form of unstructured data, it is very difficult and time-consuming to mine meaningful information based on which preemptive measures can be established. This attracts more and more researchers towards Threat Intelligence (TI) as it helps to understand threats using the deluge of data and provides actionable insights.

Threat intelligence (TI) is proof-based knowledge, which includes reasoning, context, mechanism, indicators, implications, and actionable advice, about an existing or evolving cyber-attack that can be used to create preventive measures in advance [51]. Attackers consistently exploit systems and networks to steal sensitive information, to take control of the target system, or for ransom (using ransomware) [57]. TI allows an organization to expand its visibility into the fast-growing threat landscape, can allow early identification of an attack, and successfully prevent the attack.

Indicators of Compromise (IOCs) are forensic artifacts that are used as signs that a system has been compromised by an attacker or that a system has been infected with a particular piece of malware [27]. The intent of accumulating IOCs related to

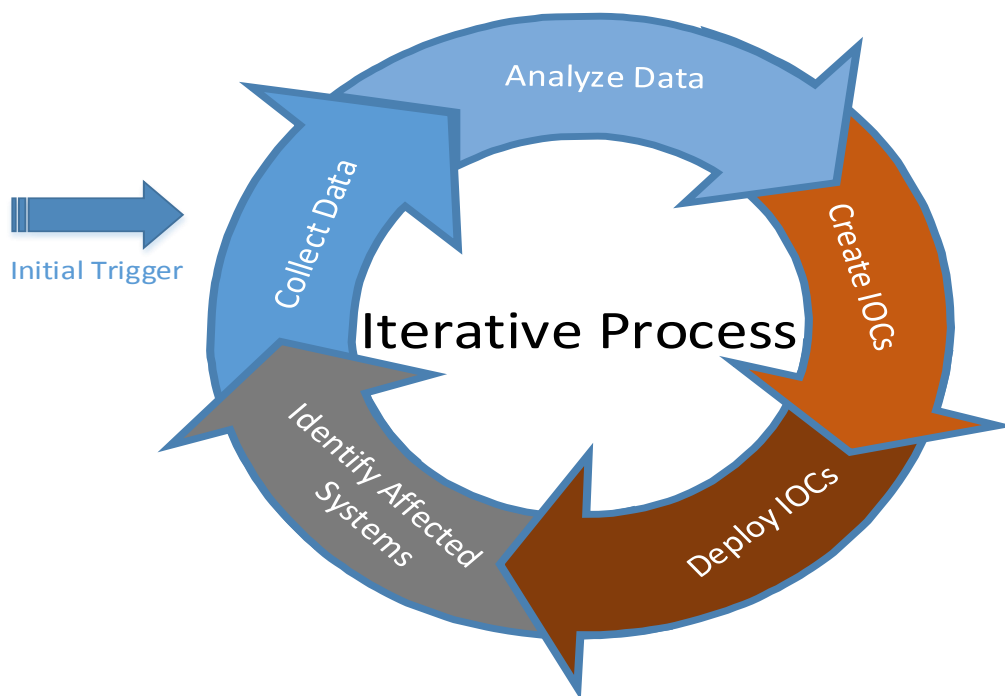


Figure 1: IOC Lifecycle. The process starts from the collection of data from internal (e.g., malware analysis reports, network logs, etc.) and external (e.g., security blogs, security white-papers, etc.) sources. Data analysis identifies the IOCs using classification techniques, log analysis, forensics analysis, false positive identification, etc. All IOCs are converted into security formats such as OpenIOC [18], snort [13], etc., during the IOC creation step. Then, these IOCs are deployed on security tools such as Intrusion Detection Systems (IDS), Host based Intrusion Detection Systems (HIDS), Security Information and Event Management (SIEM), and other investigative systems [49]. During the 5th stage, suspected systems are identified. Then, log data and forensic images of suspected systems are provided to the data collection step for the creation of new IOCs.

malware or an attack is to be able to state, with a relatively high degree of confidence, whether or not such artifacts are present in a given environment. The goal of TI

is to help security professional provide data-backed reasoning of why an artifact is identified as an IOC or not. Concretely, IOCs are composed of some combination of IP addresses, hostnames, filenames, processes, services, Windows registry entries, or hashes [23].

Figure 1 explains the typical life cycle for the use of IOCs in an organization. This life cycle consists of five steps: data collection, data analysis, IOC creation, deployment of IOCs on security infrastructure, and identification of affected systems using deployed IOCs. During data collection step, unstructured data is collected from different data sources. Data is analyzed manually or automatically during the data analysis step. Based on the analysis, IOC rules are created from the IOCs. Then, IOC rules are deployed on security infrastructure like Network Intrusion Detection Systems (NIDS) etc. As a final step, log data of these compromised systems is collected and used for the creation new IOCs, which feeds back into the IOC life cycle in a cyclical way.

Several standards are commonly used to represent IOCs for expressing cyber-threat intelligence information such as: OpenIOC [18], Structured Threat Information eXpression (STIX) [14], Cyber Observable eXpression (CybOX) [6], Trusted Automated eXchange of Indicator Information (TAXII) [19], snort rules [13], suricata rules [15], YARA rules [21], Malware Attribute Enumeration and Characterization (MAEC) [8], and Common Attack Pattern Enumeration and Classification (CAPEC) [4].

IOCs can help an organization's security personnel to attain full automation: Given a set of IOCs for a particular security event, security tools scan through an environment or infrastructure to identify the existence of any IOC on the systems in question. IOCs complement and augment existing solutions, such as Intrusion Detection Sys-

tems (IDS) or Security Information and Event Managements (SIEMs), by providing an additional and important set of information that can be used to decide whether a particular artifact being examined is malicious or not. IOCs provide a rapid route to detect new or zero-day attacks for which virus signatures or detection rules have yet to be developed for existing security tools. Thus, timely generation of IOCs is important. Also, IOC formats document a threat in a consistent fashion, thus it becomes easier for organizations to share this threat information.

However, effective collection and sharing of threat information is a challenge. Most of the current threat intelligence systems have following limitations:

- Threat data received from different sources, such as malware reports, APT white-papers, etc., is examined and interpreted by security analysts *manually*, a procedure that is timely and error-prone, which limits the practical usefulness of this data in an organizations' security infrastructure.
- Most of the traditional approaches consider either *one or two* data sources for the generation of Indicators of Compromises (IOCs) and miss valuable IOCs from alternative data sources, such as blog articles about recent attacks.
- Current state-of-the-art lacks support of *diverse IOC formats* (e.g., STIX, OpenIOC, or snort rules) which further reduces the effectiveness of these tools.
- Most importantly, many TI tools generate IOCs directly using regular expressions and white-listings from security reports or blog articles [1] without understanding the *contextual meaning* of those IOCs from the data sources which allows the tools to include lot of false positives during the process of IOCs creation.

In this project we present iGen, which is a system that tackles all of these limitation by intelligently collecting threat information from publicly available security resources (e.g., security whitepapers, blog articles etc.) and sharing the threat information in the form of IOCs (e.g., virus signatures, IPs, domains, or MD5 hashes). iGen automates the entire process of collecting publicly available data (data collection) to IOC generation. iGen supports a flexible system to collect data from diverse data sources. Currently, iGen collects data from 20 security blogs, an APT reports database [2], and a malware analysis system. The flexible nature of the system means that new data sources can be easily added. Note that the input to iGen is in unstructured English text, therefore iGen must understand the semantics of the sentences present in the input data to categorize IOC-alike strings into IOC or non-IOCs. Also, we added a practical improvement to iGen that allows transforming these IOCs into various threat information sharing standards.

First module in iGen is data acquisition module which collects key observations from the data sources like APT reports, malware analysis reports and blog articles. This module collects these reports in real-time and keep iGen updated with new IOCs from recent malware and APTs. In parallel, malware analysis system looks for new malware on web and generate malware analysis report for those samples. All the different kinds of reports accumulated from different data sources are then presented to machine learning based IOC extractor. IOC extractor module decides the parsing strategy based on the type of report. Since, malware analysis reports are based on real malware and available in machine readable format (e.g. JSON file). In this case, IOC extraction was done using a parser which accrue IOCs based on the JSON schema. But security reports or blog articles are relatively complex which tend to describe IOCs in intricate manner in English text using some context terms like “download”, “malware” etc. Google’s tensorflow [22] based Convolutional Neural Netowrk [45]

model was used to classification all the sentences in these reports into IOC or non-IOC sentences [25, 67, 52, 29]. Once sentence classification is done, next step is to identify IOCs intelligently. More explicitly, after classifying the sentence in the article into IOC and non-IOC sentences. iGen utilizes a set of regular expressions (regex) to locate IOC tokens (e.g., IP, hostname, domain, md5 etc.) within the sentence. Then, iGen extract these IOC tokens and convert them to machine-readable formats (IOC formats) which can be used directly by security tools.

Simple extraction of IOCs from publicly available security resources is straightforward: regular expressions can extract MD5s, URLs, filenames, filepaths, etc. However, this naïve approach will result in a significant amount of false positives: these security resources will include MD5s, URLs, and filenames that are *relevant to the security incident but are not malicious, and thus not an IOC*. Therefore, we need an approach that can automatically filter these potential IOCs *based on the context*. While prior work has shown that filtering potential IOCs can be done with manual creation of features [47], iGen leverages a Convolutional Neural Network, which selects features automatically. This enables iGen to be more effective and to use diverse data sets.

The main contributions of this paper are the following:

- We present the design of a novel system for intelligently generating IOCs, with high confidence, from different data sources using Convolutional Neural Networks (CNNs) (Chapter 4 and 5).
- A prototype implementation of our design as a tool called iGen (Chapter 5).

- An evaluation of iGen with other state-of-the-art methods. iGen identified around 400,000 IOCs with a precision and recall of 95% and 99% respectively (Chapter 6).

Structure of document This thesis document is grouped into the following sections:

- Chapter 2 discusses the basics of Threat Intelligence (TI) as well as different formats for sharing threat intelligence. This chapter also describes the Convolutional Neural Networks (CNN) which is a core of iGen system.
- Chapter 3 discusses the data collection from different sources.
- Chapter 4 discusses the system design, the architecture, and the components of iGen system.
- Chapter 5 discusses the implementation details of iGen system.
- Chapter 6 presents our findings and comparative analysis of the results.
- Chapter 7 continues the discussion of the results; the lessons learned, limitations, and improvements that are possible on current iGen system.
- Chapter 8 explores related work in the area.
- Chapter 9 concludes this thesis, with ideas to develop the research in this area.

Chapter 2

BACKGROUND

In this section, we provide the background necessary for understanding the design of iGen. In Section 2.1, we describe Threat intelligence and how it is consumed by security tools. Then, in Section 2.2, we describe Convolutional Neural Network (CNN) and how it can be used in iGen.

2.1 Threat Intelligence

Given the increasing number and rapidly evolving nature of current threats, quick sharing and exchange of relevant threat information is the key to swiftly detecting, understanding, and responding to cyber-attacks. Threat intelligence is two things: (1) a process of collection of knowledge that defines security threats, which empowers an organization to determine its responses at the strategic, operational, and tactical levels, and (2) information that has been analyzed to discover actionable insights. Actionable threat intelligence is insight that an organization can act on—it enables informed decision-making that results in improved outcomes.

Threat intelligence consists of two words: “Threat” and “Intelligence”. “Threat” is an agent (that is, a menace or hazard) that takes advantage of the vulnerability. Whereas, “Intelligence” can be defined as:

- Act of generating new information by correlating information from different sources.
- Capacity to know or understand.
- Knowledge imparted through study, research or experience.

2.1.1 Consumption of Threat Intelligence

Indicators of Compromise (IOCs) are used by organizations to exchange threat intelligence. There are different standards for IOCs which includes OpenIOC, STIX, and snort rules. OpenIOC was introduced by Mandiant in 2011. It is primarily used in Mandiant products, but has also been released as an open standard. OpenIOC provides definitions for specific technical details including over 500 indicator terms. Adding new terms is easy because the terms are separated from the main schema. Most of the terms are host-centric.

Similarly, STIX is another format for storing IOCs. STIX allows defining threat information, including threat details, as well as the *context* of the threat. STIX is developed by Mitre [10], and is designed to support four cyber threat use cases: investigating cyber threats, stating indicator patterns, response activities management, and sharing threat information. STIX uses XML to define threat-related constructs such as exploit target, campaign, indicator, threat actor, and Tactics, Techniques and Procedures (TTP).

Snort is an open-source network intrusion prevention system (IPS) [63], which executes real-time traffic analysis and packet-logging on IP networks. It is also capable of performing protocol analysis, content searching, and matching, and can be used to detect a variety of attacks and probes, such as buffer overflows, stealth port scans, CGI attacks, SMB probes, OS fingerprinting attempts, and more. Snort uses a rule language to describe traffic that it should collect or pass, as well as a detection engine that uses a modular plug-in architecture. Rules created using this language are called snort rules.

Suricata rules are the defacto method for sharing and matching threat intelligence against network traffic. A suricata rule has three components: The action, header and rule-options.

2.2 Convolutional Neural Network Architecture

Recently, different models based on deep learning have achieved significant results in computer vision and speech recognition. Convolutional Neural Networks (CNNs) were responsible for major breakthroughs in image classification [45] and are the core of most computer vision systems today, from Facebook’s automated photo tagging [65] to self-driving cars [26].

In the case of Natural Language Processing (NLP), much of the work that uses deep learning methods involves learning word vector representations through neural language models and performing composition over the learned word vectors for classification tasks. Words are mapped into a lower dimensional vector space via a hidden layer. Hidden layer extracts the semantic features of the word and map them into the word vectors. Euclidean or cosine similarity is high between the word vectors of semantically close words.

Recently, researchers have also started applying CNNs with word embeddings to problems in Natural Language Processing (NLP) and got some interesting results. CNNs utilize layers with convolving filters that are applied to local features. CNNs are effective for NLP and have achieved excellent results in semantic parsing, search query retrieval, sentence modeling, and other traditional NLP tasks.

2.2.1 Convolution

Convolution [41] is a sliding window function applied to a matrix. The sliding window is called a kernel, filter, or feature detector. Here we use a $N \times N$ filter,

multiply its values element-wise with the original matrix, then sum them. To get the full convolution we do this for each element by sliding the filter over the whole matrix. Final matrix created after applying the filter is called convolved feature.

2.2.2 Convolution Neural Network

CNNs are fundamentally several layers of convolutions with nonlinear activation functions like rectified linear unit (ReLU) [54] or hyperbolic tangent (tanh) [37] applied to the results. In a traditional feedforward neural network, each input neuron is connected to each output neuron in the next layer. It is called as fully connected layer, or affine layer. But it's different in case of CNNs. CNNs use convolutions over the input layer to compute the output. This results in local connections, where each region of the input is connected to a neuron in the output. Each convolution layer applies different filters, characteristically hundreds or thousands, and combines their results. Then pooling is performed using the pooling (subsampling) layer. We will discuss more in Chapter 4. CNN automatically learns the values of its filters based on the task during the training phase. For example, in image classification a CNN learns to detect edges from raw pixels in the first layer, then uses the edges to detect simple shapes in the second layer, and then uses these shapes to detect higher-level features, such as facial shapes in higher layers. The last layer is then a classifier that uses these high-level features.

In the case of iGen, input to a CNN is sentences extracted from blog articles and security reports instead of image pixels. Each sentence is represented as a matrix. Each row vector of the matrix corresponds to one token, typically a word. That is, each row is vector that represents a word. These vectors might be, e.g., outputs from trained word2vec or GloVe [59] models. We will discuss more about the iGen CNN in Section 3.2.

Chapter 3

DATA COLLECTION

There are two sources of TI data: internal, e.g., malware analysis reports or network traces, or external, such as technical blogs or security reports. Examples of external sources include the Kaspersky whitepapers [7], Symantec blog [16] etc. Recently, thousands of threats are reported every day, which has been broadly reported by the security companies in the form of APT/security reports [31] or blog articles and aggressively collected by different organizations. For the scope of our research, we collected 1000 security reports and 1500 blog articles from different security organizations. All these reports and articles are published from year 2011 to 2017. Our scrapper collected information from more than 20 prominent security organization like Verizon, Symantec etc.

To bootstrap our research, we also collected around 35,000 malware reports from our in-house Cuckoo Malware Analysis System (CMAS) [58]. MAS provided some detailed results outlining what malware did when executed inside an isolated environment. While substantial volume of security reports, blog articles and malware reports

Table 1: Summary of Dataset

Dataset Type	Dataset Source	Number of articles/ reports
Security Reports	External	1000
Technical Blog Articles	External	1500
Malware Analysis Reports	Internal	35000

are collected and analyzed, however it only constitutes a small part of the bigger IOC landscape. Summary of our dataset is in Table 1.

SYSTEM DESIGN

Human are very intelligent in understanding natural language text and finding the information using reasoning. In our case, this information is IOCs from security reports, articles, and malware reports. However, it is very hard for the machines to do the information extraction with high accuracy and coverage.

iGen identifies sentences likely containing IOCs using a set of regular expressions. Then, IOCs are efficiently captured with high accuracy using CNN-based model. iGen tries to automatically identify the patterns or features which are frequent in IOC and non-IOC sentences. Architectural design of the iGen is shown in Figure 2.

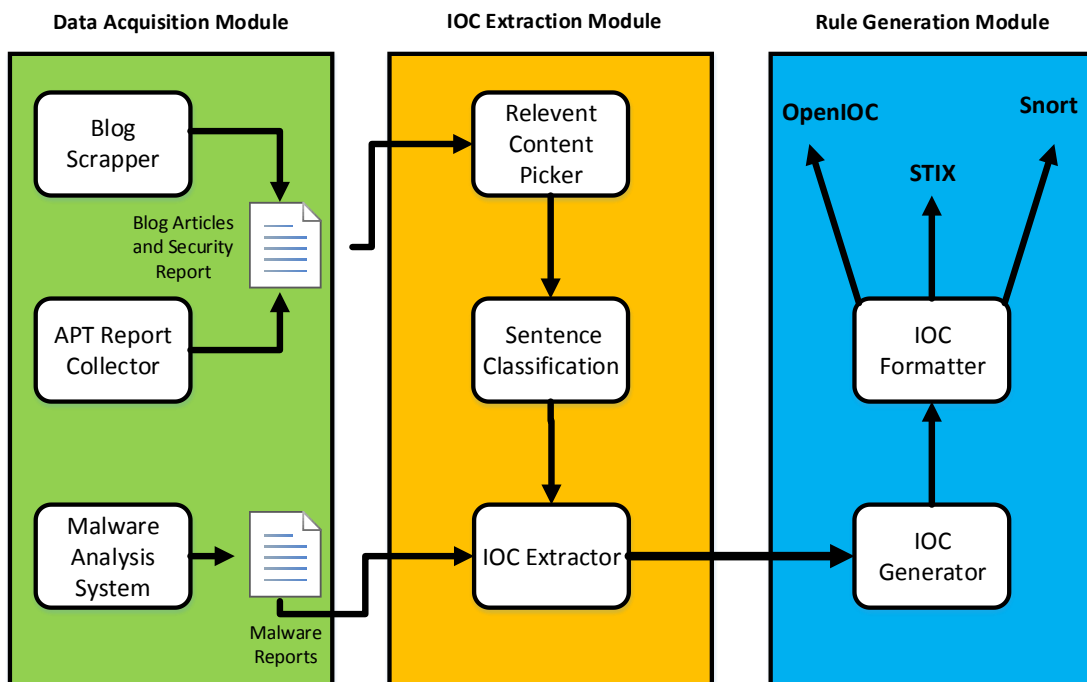


Figure 2: iGen Architecture Overview.

4.0.1 Architecture Overview

iGen consists of 3 modules from left to right: 1) data acquisition module, 2) IOC extraction module, and 3) rule generation module. Data acquisition module consists of three sub-module: Blog Scrapper (BS), APT Report Collector (ARC), and Cuckoo Malware Analysis System (CMAS). BS is a collection of crawlers based on the HTML structure of the different blogs. For instance, we designed a dedicated crawler which is continuously looking for new articles on Symantec public blog [16]. Similarly, ARC continuously looks for new whitepapers and reports about APT campaigns at “APTNotes” [2] database. CMAS collects malware reports from Cuckoo [24], which is an open source malware analysis system.

Then, these reports are passed to an IOC extraction module based on the type of the report. Because blog articles and security reports are written in English and follow natural language semantics, these are given as input to the Relevant Content Picker (RCP) component for further processing. RCP parses those reports, cleans the text, and selects the sentences having likely IOCs using regular expressions. Then, these sentences are fed to the Sentence Classification (SC) module for classification into two classes: IOC and non-IOC sentences. Whereas, malware reports follow JSON format with no natural language semantics involved. This makes it easier for IOC Extractor (IE) to extract IOCs based on the structure of the JSON report. Simultaneously, IE also performs regular expressions based extraction of IOCs from the sentences that are classified by Sentence Classifier (SC) as IOC sentences.

Next step is to organize IOCs extracted from each report into one event (malware or APT campaign) is done by IOC Generator (IG). For example, “Monkeys.exe” and “200.125.133.28” are the IOCs related to the event CozyDuke APT campaign [5]. Then, these events are stored and managed by IOC Formatter (IF) [9]. IF also

provide an additional functionality to iGen by converting these event into different output formats like STIX, OpenIOC, Snort etc.

4.1 Data Acquisition Module

This module collects information from different sources. Because of the scalable nature of the module, new data sources can be added as well. This module scrap related web pages, APT report from different sites. Simultaneously, it continuously collects malware reports from our internal CMAS as well. All the sub-modules of data acquisition module are described in the following.

4.1.1 *Blog Scrapper (BS)*

iGen blog scraper is designed to continuously monitor a list of security blogs and collect their articles. BS first scraps all its current articles before it is set to monitoring mode to identify new ones for each blog site. Most of the blog articles on most of the blog sites have a particular page id associated with it. BS keep track of the range of page ids that are already scrapped by it. In monitor mode, it looks for new articles. If there are no new articles posted on the blog sites, then BS changes its mode to idle and want back to monitor mode after n days. For our research, we have set the n as 7 days. BS also collects the reports if there is any PDF report link within the blog article.

4.1.2 *APT Report Collector (ARC)*

APT report collector consists of autonomous data crawler, which continuously check for new reports and whitepapers at APTNotes database. It also keeps track of the already scrapped reports to avoid duplication. Around 1000 APT campaign report are collected till date.

4.1.3 Cuckoo Malware Analysis System (CMAS)

For our research, we have a Cuckoo malware analysis system (CMAS) [24] for automating analysis of suspicious files. CMAS monitors the behaviour of the malicious processes while running in an isolated environment. In other words, if we submit any suspicious file to it and in a matter of seconds CMAS will provide us back some detailed results outlining what such file did when executed inside an isolated environment. Features of CMAS are :

- Analyze different type of malicious files (executables, document exploits, malicious websites etc.).
- Dump and analyze network traffic, even when encrypted.
- Trace API calls and general behavior of the file.

For scope of this research, around 35000 malware were submitted and analyzed by CMAS. Later, these reports were further used for IOCs extraction.

4.2 IOC Extraction Module

After the collection of dataset, the next step is to clean the data, select relevant content, and find IOCs with high degree of confidence. In this module, Relevant Content Picker (RCP) cleans the data by removing the terms with no semantic meaning such as stop words, etc. RCP also finds sentences which are likely to contain IOCs. Afterwards, Sentence Classification (SC) and IOC Extractor (IE) identify the IOC sentences and extract the IOCs from them respectively.

4.2.1 *Relevant Content Picker (RCP)*

Relevant Content Picker extracts text from PDF and HTML documents. In our implementation, the RCP uses the python library “pdfminer” [12] to extract the text from the PDF reports and “beautifulsoup” [3] to obtain content from HTML pages as well to find links to security reports in blog articles. RCP scans for URLs with an extension “.pdf” in the blog articles to identify security reports linked with the blog articles. If RCP finds any PDF report linked with blog article, it sends it to PDFminer for further processing. After the collection of content, the next step is to remove the terms or token with no semantic meaning. These words such as auxiliary verbs, conjunctions and articles can be ignored. These words are called stopwords. Stop words extremely common words which would appear to be of little value in terms of semantic significance like “a”, “an”, “the” etc. Stopwords are removed as a part of the data cleaning process.

Stemming is another very important pre-processing step. Size of feature set can be reduced by removing misspelled or words with the same stem. Algorithm called stemmer, removes words with the same stem and keeps the stem or the most common of them as feature. For example, the words "download", "downloads", "downloading" and "downloaded" can be replaced with "download".

Final step is to look for sentences which are likely to have IOCs. This part is done by finding the sentences which have patterns like IP, hostname, URL, filename, registry, email, filepath etc. These are patterns are identified using the set of total eleven regular expressions. Sample regular expressions are shown in Table 2.

Type of Regular Expression	Regular Expression
CVE	<code>\b(CVE\[0-9]{4}\-[0-9]{4,6})\b</code>
Email	<code>\b([a-z][_a-z0-9-.\+@][a-z0-9-.\+][a-z]+)\b</code>
Filepath	<code>\b[A-Z]:\\[A-Za-z0-9-_\.\+]\b</code>
IP	<code>\b(\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3})\b</code>
MD5	<code>\b([a-f0-9]{32} [A-F0-9]{32})\b</code>
Registry	<code>\b((HKLM HKCU)\\[\\A-Za-z0-9-_\+])\b</code>
SHA1	<code>\b([a-f0-9]{40} [A-F0-9]{40})\b</code>
SHA256	<code>\b([a-f0-9]{64} [A-F0-9]{64})\b</code>

Table 2: Sample Regular Expressions used by Relevant Content Picker (RCP) and IOC Extractor (IE).

4.2.2 Sentence Classification (SC)

Sentence classification sub-module takes sentences and output labels (IOC sentence or non-IOC sentence) for training. Training data is generated by manual labeling of sentences into IOC and non-IOC sentences. This task of sentence tagging was done by the security experts.

A sentence is a sequence of words [42]. So each sentence is usually represented by an array of words. The set of all the words of a training set is called vocabulary, or feature set. So a sentence can be presented by a binary vector, assigning the value 1 if the sentence contains the feature-word or 0 if the word does not appear in the sentence.

Feature selection is the next step in the sentence classification task. Feature selection can be done automatically or manually. Automatic feature selection is done by classifier like Convolutional Neural networks. Whereas, manual feature selection

is done using exploratory analysis, then these generated features and output labels are presented to the classifier.

The aim of manual feature-selection methods is to resolve the curse of dimensionality by removing features that are considered irrelevant for the classification. These procedures have several advantages like smaller dataset size, smaller computational requirements for the text categorization algorithms and considerable shrinking of the search space. The goal is the to reduce dimension to yield improved classification accuracy. Another advantage of feature selection is to reduce over-fitting, i.e. the phenomenon by which a classifier is tuned also to the contingent characteristics of the training data rather than the constitutive characteristics of the categories.

After feature selection, final step is to train the classification model on the training data consists of features and output labels. In our approach, we evaluated different classifiers like Naive Bayes [46], Support Vector Machines (SVM) [42], K-Nearest Neighbour (KNN) [38], and Convolution Neural Network [45].

Naive Bayes (NB) classifier is one of the popular approach in sentence classification. NB Classifier assign each sentence s the class $c^* = \arg \max_c P(c | s)$. *Naive Bayes* classifier is based on the Bayes' rule,

$$P(c | s) = \frac{P(c)P(s | c)}{P(s)},$$

where $P(s)$ has zero influence in selecting c^* . Estimation of the term $P(s | c)$ is done using Naive Bayes decomposition by assuming the f_i 's are conditionally independent given s 's class:

$$P_{\text{NB}}(c | s) = \frac{P(c) (\prod_{i=1}^m P(f_i | c))}{P(s)}.$$

Training phase consists use of add-one smoothing for relative-frequency estimation of $P(c)$ and $P(f_i | c)$.

Despite its simplicity and the fact that its conditional independence assumption clearly does not hold in real-world situations, Naive Bayes-based sentence categorization still tends to perform well.

Support vector machines (SVMs) have been proved to be highly accurate at sentence classification task, generally outperforming Naive Bayes [42]. They are *margin based classifier*, rather than probabilistic classifiers, in contrast to Naive Bayes. In our two-category (IOC or non-IOC sentence) case, the basic idea behind the training phase is to estimate a hyperplane, represented by vector \vec{w} , that not only separates the sentence vectors in one class from those in the other, but for which the separation, or *margin*, is maximum.

This search corresponds to a constrained optimization problem; letting $c_j \in \{1, 0\}$ (corresponding to IOC and non-IOC sentence) be the correct class of sentence s_j , the solution can be written as

$$\vec{w} := \sum_j \alpha_j c_j \vec{s}_j, \quad \alpha_j \geq 0,$$

where the α_j 's are obtained by solving a dual optimization problem. Those \vec{s}_j such that α_j is greater than zero are called *support vectors*, since they are the only document vectors contributing to \vec{w} . Classification of test instances consists simply of determining which side of \vec{w} 's hyperplane they fall on.

Similarly, K-Nearest Neighbour (KNN) method is considered one of the simplest and most effective sentence classification algorithms. It is based on the principle that given set of instances in a training set, the class of a new yet unseen occurrence is likely to be the same as the majority of its closest neighbour instances from the training set. Thus the k-Nearest Neighbour algorithm works by inspecting the k closest instances in the data set to a new occurrence that needs to be classified, and making a prediction

based on the k nearest neighbours. The notion of closeness is formally given by a distance function between two points in the attribute space. An example of distance function typically used is the standard Euclidean distance between two points in an n -dimensional space, where n is the number of attributes in the data set.

Recently, several researchers have used word embeddings and Convolution Neural networks (CNN) for sentence classification. These CNN based classifiers have outperformed other classifiers in most of the classification tasks. We will discuss more about CNN in the Chapter 5.

4.2.3 IOC Extractor (IE)

As shown in Figure 2, IOC extractor takes malware report as well the sentences which are classified as IOC sentences by SC sub-system as input. Since, malware reports doesn't have any natural language semantics, we skipped running RCP and SC sub-systems on them. Additionally, MAS is run only on the malware files and all the malicious indicator/tokens (e.g., file hashes, emails, etc.) are defined in a machine readable format.

In case of malware reports, structure of JSON file is used to extract the IOCs. Whereas, regex as shown in Table 2 were used to extract the IOC from the sentences. These are different type of IOCs collected by the IE: URL, host, IP, email, MD5, SHA1, SHA256, CVE, registry, filename, and filepath. Next step is to bundle the IOCs related to one event/malware together.

4.3 Rule Generation System

This system works on the managing the related IOCs and convert them into different output formats so that these IOCs based formats can be directly deployed on

the security infrastructure. Currently, iGen supports STIX, OpenIOC, snort, suricata, and BRO formats. Rule generation system has two sub-systems: IOC Generator (IG) and IOC Formatter (IF).

4.3.1 IOC Generator (IG)

IOC generator takes IOCs from IE and bind them together into one IOC event if they are coming from the same report (e.g., malware report, security whitepapers, blog article etc.). This sub-system creates event into IF acceptable JSON format. Then, this JSON is submitted to IF.

4.3.2 IOC Formatter (IF)

Malware Information Sharing Platform (MISP) [9] is used as IOC Formatter in iGen. MISP is an open source software for sharing, storing and correlating IOCs of targeted attacks. MISP benefits security community by showing collaborative knowledge about existing malware or threats. The aim of this trusted platform is to help improving the counter-measures used against targeted attacks and set-up preventive actions and detection. It stores data in a structured format and allows automated use of the database to feed detection systems or forensic tools. For example, it generates rules from IOCs (e.g. IP addresses, domain names, hashes of malicious files, filenames, filepath etc.) for different Network Intrusion Detection Systems (NIDS) like snort, suricata etc. MISP adds an extra feature of generating diverse rules from the IOCs collected from different data sources.

IMPLEMENTATION

We implemented iGen that collects IOCs from various external and internal IOC data sources. The automatic crawling and analysis systems were implemented in Python. For example, Blog Scrapper (BS) and APT report Collector (ARC) were built using Python based libraries like pdfminer, BeautifulSoup and textblob. Additionally, we are generating the malware reports using our in-house Cuckoo malware analysis system. Relevant Content Picker (RCP), IOC Extractor (IE) and IOC Generator (IG) uses re and nltk for information extraction. Sentence Classification (SC) uses numpy and Google’s tensorflow library. Different output formats are generated using the open-source tool Malware Information Sharing Platform (MISP).

We deployed our system iGen modules on our Openstack [11]. Data acquisition module and IOC extraction module on a dedicated Openstack instance with 8GB RAM. Whereas, rule generation module runs on other instance with 4GB RAM.

5.1 CNN based Sentence Classification

The model architecture of Sentence Classification (SC), shown in figure 3, is a minor variant of the CNN architecture suggested by Kim [44]. First step of SC involves converting all the words in the sentences into their respective word vector or embedding. Let $\mathbf{S}_i \in \mathbb{R}^d$ be the word vector of dimensionality d corresponding to the i -th word or token in the sentence. A sentence of length l is characterized as

$$\mathbf{S}_{1:l} = \mathbf{S}_1 \oplus \mathbf{S}_2 \oplus \dots \oplus \mathbf{S}_l, \quad (5.1)$$

where concatenation operator is defined as \oplus .

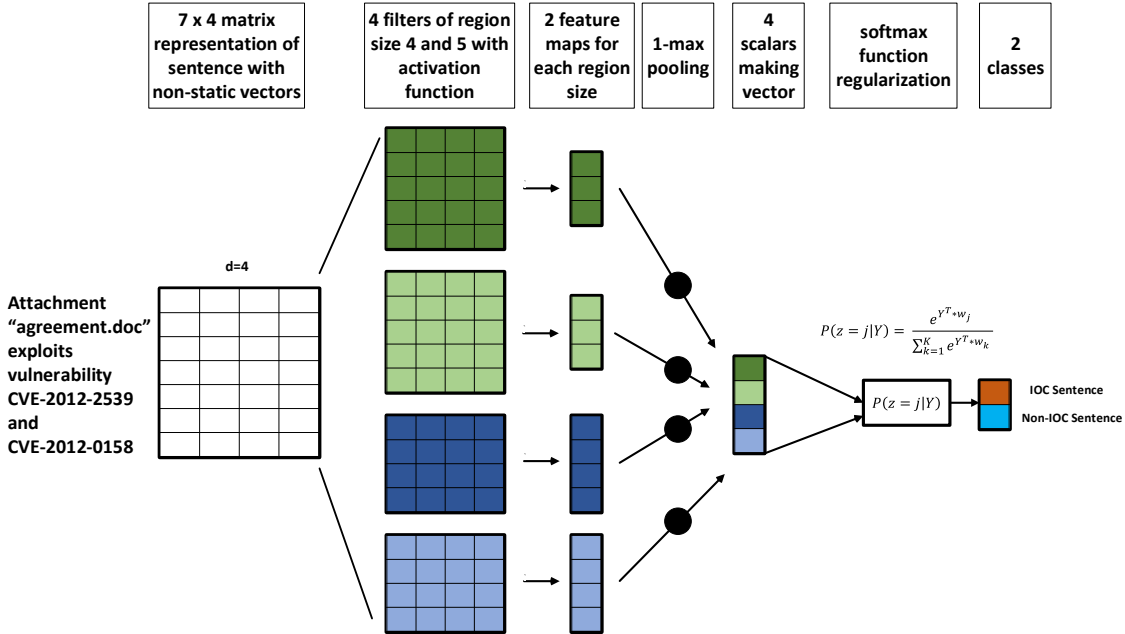


Figure 3: Architecture of iGen CNN based sentence classification model. We show two filter region sizes: 4 and 5, each of which has 2 filters. Filters accomplish convolutions and apply activation functions on the sentence matrix (formed using word embeddings of all the words in sentence) and generate (variable-length) feature maps; the largest number from each feature map is recorded using 1-max pooling function. Then scalar values generated from all 4 maps are concatenated to form a feature vector for the second to last layer. This feature vector is used by final softmax layer to classify the sentences; since iGen CNN is a binary classifier, only two output values (IOC sentence or non-IOC sentence) are possible [68].

The first step is to calculate word embedding of S_i based on our dataset. We call this layer as embedding layer, which maps vocabulary word indices into low-dimensional vector representations. X is our embedding matrix that we learn during training. We initialize it using a random uniform distribution [28]. `tf.nn.embedding_lookup` [17] function in tensorflow creates the actual embedding operation. The result of the

embedding operation is a 3-dimensional tensor of shape [None, sequencelength, embeddingsize]. TensorFlow’s convolutional conv2d [17] operation expects a 4-dimensional tensor with dimensions corresponding to batch, width, height, and channel. The results of our embedding does not contain the channel dimension, so we add it manually, leaving us with a layer of shape [None, sequencelength, embeddingsize, 1].

A word embedding $\mathbf{X}: words \rightarrow \mathbf{R}^n$ is a parameterized function mapping words in some language to high-dimensional vectors (perhaps 200 to 500 dimensions). For example

$$\mathbf{X}(\text{“malware”}) = (0.2, -0.4, 0.7, \dots) \quad (5.2)$$

Characteristically, the function is a lookup table, parameterized by a matrix, θ , with a row for each word: $\mathbf{X}_\theta(w_n) = \mathbf{X}_n = \theta_n$.

\mathbf{X} is initialized to have random vectors for each word. It learns to have meaningful vectors in order to classify the sentences into IOC and non-IOC sentences.

Let $\mathbf{S}_{i:j}$ refer to a sentence of length $j - i + 1$ or the concatenation of words vectors $\mathbf{S}_i = \mathbf{X}_\theta(w_i), \mathbf{S}_{i+1} = \mathbf{X}_\theta(w_{i+1}), \dots, \mathbf{S}_j = \mathbf{X}_\theta(w_j)$. Single convolution consist of a *filter* $\mathbf{w} \in \mathbb{R}^{h \times d}$, which is applied to a window of h words to automatically generate a new feature. For example, a feature f_i is generated from a window of words $\mathbf{S}_{i:i+h-1}$ by

$$f_i = f(\mathbf{w} \cdot \mathbf{S}_{i:i+h-1} + c). \quad (5.3)$$

Here bias term $c \in \mathbb{R}$ is added and f is a non-linear function such as the rectified linear unit (ReLU) or hyperbolic tangent (tanh). Filter is applied to each possible window of words in the sentence $\{\mathbf{S}_{1:h}, \mathbf{S}_{2:h+1}, \dots, \mathbf{S}_{l-h+1:l}\}$ to produce a *feature map*

$$\mathbf{f} = [f_1, f_2, \dots, f_{l-h+1}], \quad (5.4)$$

where $\mathbf{f} \in \mathbb{R}^{l-h+1}$. We then apply a 1-max pooling operation [60] over the feature map and take the maximum value $\hat{f} = \max\{\mathbf{f}\}$ as the feature corresponding to this

particular filter. Intuition behind this is to capture the most important feature—one with the highest value—for each feature map. This kind of approach works naturally with variable sentence lengths.

Above described process extracts *one* feature from *one* filter. The model uses multiple filters with different region size to obtain multiple features. These features creates the penultimate layer and the final layer is a fully connected softmax layer whose output is the probability distribution over labels (in our case, “IOC sentence” or “non-IOC sentence”).

In our model, we have fine-tuned word vectors via backpropagation [62]. For regularization [70] we used dropout [64] on the penultimate layer with a constraint on l_2 -norms [55] of the weight vectors [40]. Dropout helps in co-adaptation of hidden units by randomly dropping out—i.e., setting to zero—a proportion p of the hidden units during forward-backpropagation. That is, given the penultimate layer $\mathbf{y} = [f_1, \dots, f_n]$ where n is total number of filters). Instead of using

$$z = \mathbf{w} \cdot \mathbf{y} + b \tag{5.5}$$

during forward propagation for output unit y , dropout uses

$$z = \mathbf{w} \cdot (\mathbf{y} \circ \mathbf{r}) + b, \tag{5.6}$$

where \circ is the element-wise multiplication operator and $\mathbf{r} \in \mathbb{R}^n$ is a ‘masking’ vector of Bernoulli random variables with probability p of being 1. During the training, gradients[35] are backpropagated only through the unmasked units. During the testing, all learned weight vectors are scaled by p such that $\dot{\mathbf{w}} = p\mathbf{w}$, and $\dot{\mathbf{w}}$ is used to score new sentences without dropout.

Chapter 6

EVALUATION

6.1 Experiments

Datasets used in the experiments are shown in Table 1. For sentence classification, only blog articles and security reports were used. From these articles, we further extracted two kinds of sentences, those with IOCs (IOC sentences) and those without but involving non-malicious IOC-like strings (non-IOC sentences). More specifically, the 5,000 IOC sentences and 2,500 non-IOC sentences were used in the experiments. These sentences are manually tagged by security experts to maintain the quality of data.

In the experiments, the parameters of our iGen (most of them related to CNN sub-module) system were set as follow:

- *Sequence length.* The length of our sentences before inputting them for classification. Remember that we padded all our sentences to have the same length of 70. We selected 70 as sequence length, since longest sentence in our dataset has 60 word length.

Embedding Size(d)	Accuracy
64	89.3%
128	95.1%
256	95.1%

Table 3: Change in accuracy with embedding size.

Number of Filters	Accuracy
32	88.5%
64	95.1%
128	94.6%

Table 4: Change in accuracy with number of filters.

- *Embedding Size.* The dimensionality of our word embeddings or word vectors is kept as 128. Word embedding dimensionality was chosen as 128 based on the experiments where we compared accuracy with dimensionality of word embeddings as shown in Table 3. We did not increase the dimensionality beyond 128 as it has a bad effect on the performance.
- *Filter Sizes.* Filter size is the number of words we want our convolutional filters to cover. We have used three types of filter 3, 4 and 5 that slide over 3, 4 and 5 words respectively.
- *Number of Filter.* It is the number of filters per filter size. Based on the experiments shown in Table 4, it is chosen as 64. Total $3 * 64$ filters are used with 64 filters each of size 3, 4, and 5.
- *Dropout Keep Probability.* Dropout is the most popular method to regularize convolutional neural networks. A dropout layer stochastically disables a fraction of its neurons. This prevents neurons from co-adapting and forces them to learn individually useful features. The fraction of neurons we keep enabled is defined by the *dropout_keep_prob* input to our network. Based on the experiments, we set this to 0.5 during training, and to 1 (disable dropout) during evaluation.

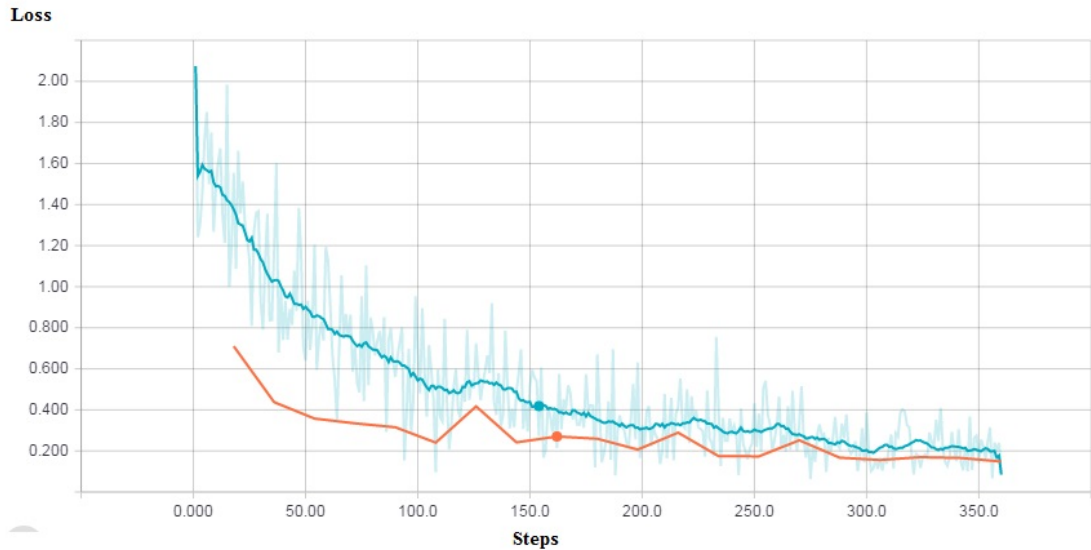


Figure 4: Change in loss function over steps (blue is training data, red is 10% dev data).

6.2 Results

6.2.1 Loss Function and Accuracy

The loss is a measurement of the error our network makes, and our goal is to minimize it. The loss function for categorization problems is the cross-entropy loss. `tf.nn.softmax.cross_entropy_with_logits` is a function that calculates the cross-entropy loss for each class, given our scores and the correct input labels. We have then taken the mean of the losses. We could also use the sum, but that makes it harder to compare the loss across different batch sizes and train/dev data. Figure 4 shows the change cross-entropy loss over steps.

We also used an expression for the accuracy, which is a useful quantity to calculate the performance during training and testing. Change in accuracy over steps is shown in Figure 5.

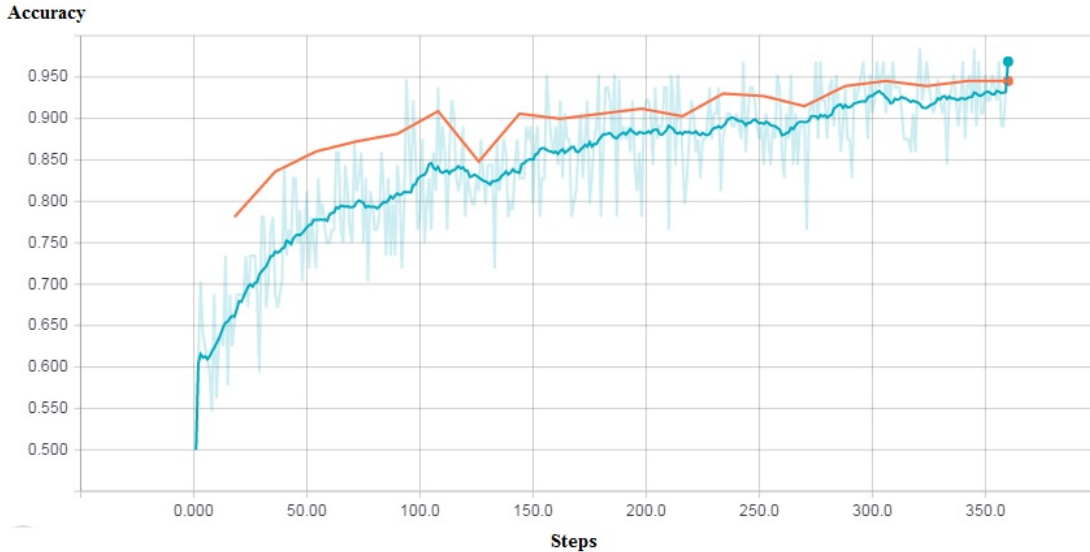


Figure 5: Change in accuracy over steps (blue is training data, red is 10% dev data).

6.2.2 Method Comparison

We compared iGen with other state-of-the-art methods like Support Vector Machine (SVM) [30], Naïve Bayes Classifier [36], and KNN (K-Nearest Neighbour) [48] classifier used for IOC sentence categorization. Three type of feature vector are used for each of the classifier described above. 10-fold cross validation technique is used to compare precision, recall and F-measure. Our comparative analysis is shown in Table 5.

iGen outperformed other methods by generating IOCs with a precision of 95% and a recall of 99%, while among other SVM with tfidf as feature vector performed well with a precision and recall of 90%.

Similarly, another approach iACE [47] proposed by Liao gave promising result in extracting IOC with high accuracy. But iACE used around 5,283 terms as features which were manually gathered from the dataset. This implies that it will work very

Methods	Precision	Recall	F-measure
iGen	95%	97%	96%
SVM ¹ (Bag of words)	89%	88%	88%
SVM (tf ²)	90%	89%	89%
SVM (tfidf ³)	90%	90%	90%
Naive Bayes (Bag of words)	81%	78%	79%
Naive Bayes (tf)	83%	79%	81%
Naive Bayes (tfidf)	85%	83%	83%
5-NN ⁴ (Bag of words)	80%	80%	80%
5-NN (tf)	81%	81%	80%
5-NN (tfidf)	83%	82%	82%

Table 5: Results of our iGen models against other methods. **SVM** (Support Vector Machine), **NB** (Naive Bayes), and **5 – NN** (5-Nearest Neighbours) classifiers were compared with feature vectors as bag of words (BOW), term frequency (tf), and term frequency - inverse document frequency (tfidf).

accurately on one dataset but may fail on sentences where different terminologies were used while writing the article. In case of iGen, features were not selected manually but features were selected by CNN automatically. This property makes iGen more adaptive towards any kind of data.

¹SVM — Support Vector Machine.

²tf — Term Frequency.

³tfidf — Term Frequency Inverse Document Frequency.

⁴5-NN — 5-Nearest Neighbour

DISCUSSION AND FUTURE WORK

Our experiments show that iGen not only fully automated the process of cyber threat intelligence collection but also generates deployable security rules from the unstructured data. With a large amount of IOCs automatically recovered from the wild and converted into a machine-readable form, these can be quickly and effectively utilized to counter emerging threats. For example, knowing the IP of C&C server from APT report and then finding the email associated with IP address from malware report can help us find the actor behind the attack. This will enable the defender to disable or block the servers associated with that email to stop future attacks. On the other hand, our current design is still preliminary. Here, we discuss the limitations of iGen and potential follow-up research.

7.0.1 Limitations

Although iGen extracts IOC with high accuracy, precision and recall, but it's necessary to recognize our system limitations. Some of the future work to enhance iGen framework is given below:

- First, iGen currently considering data sources like security reports, blog articles that are written in English. But lots of security documents are written in other language as well such as Arabic, Persian etc. Next step is to consider data sources that are written in other languages.
- Second, integration of more sources of data like social networks, open sources threat feeds, etc. More sources of data will help us to collect more threat

information about the attacks and malwares. More data will help iGen to build more accurate IOC rules .

- Third, test sentence classification sub-module with other state-of-the-art classifiers like random forest, decision tress etc. It will help to select the best classifier based on the evaluation metrics like precision, recall, f-measure etc.
- Fourth, expansion of iGen to accommodate other standards like TAXII, YARA rules, etc.

Further efforts are required to enhance the functionality of iGen framework.

Chapter 8

RELATED WORK

8.1 Knowledge discovery

The recent decades witnesses a rapid flow of information available in digital form on the Internet and intranets. Most of this information is present in the form of reports, blog posts, news articles, etc. This information is transmitted through unstructured documents and is thus difficult for machine to search in. This created a need for automatic, effective and efficient methods for analyzing the unstructured data and discovering relevant knowledge from it in the form of structured information, and led to more research in Information Extraction (IE) techniques. In particular, recent advances in the field of Natural Language Processing (NLP), specifically in the field of text processing techniques, have resulted to the spread of deployment of IE techniques in real-world applications for processing of vast amount of textual data.

One of the strong application of IE was discovered in 1986. Lytinen et.al. presented ATRANS system [50] which applies IR in the financial domain to extract information from messages regarding money transfers between banks. In particular, ATRANS tries to solve the problem of context localization in the absence of reliable syntactic clues, such as sentence boundaries.

Recently, NLP techniques are widely used in different fields of security. Darling et al. [33] developed a classification system based on lexical features of URLs for detecting phishing and malicious domains. Similarly, N-Gram based techniques is used to detect unknown network attacks [61]. System proposed by Zhu et al. [69] automatically engineers features for Android malware detection by mining scientific

papers. Some of the recent techniques also tried to extract IOCs using manual feature selection and graph similarity comparison [47]. Instead, we focus on automated feature selection using CNN for IOC identification and we created security rules using those IOCs.

8.2 Convolutional Neural Network

Deep learning has been used in different area for multiple applications recently, including object recognition [45], speech recognition [39] and natural language processing [44]. Among the different deep learning strategies, Convolutional Neural Networks (CNNs) have been successfully applied to different NLP tasks such as sentiment analysis [34], question classification [43], and hashtag prediction [66]. Kim [44] proposed an approach for sentence classification model which uses CNN trained on top of pre-trained word vectors. Their simple CNN with one layer of convolution performs remarkably well. Similar experiments for sentence classification task were performed by Zhang et al [68]. Similarly, iGen uses CNN for classification of sentences into IOC and non-IOC sentences. iGen is one of the few systems that is using CNN for information extraction in the field of security.

8.3 Sharing Threat Intelligence

Sharing TI is very important to help the organizations defend against the fast-evolving threat landscape. Big organization like IBM, Crowdstrike, McAfee, etc. are integrating their TI with existing products. Dandurand et al. [32] discussed that cyber security community requires efficient systems to facilitate threat information sharing and automation. Murdoch et al. [53] discussed the challenges and motivation for threat information sharing; like trust issues and other challenges like keeping the online community active to contribute. Evaluation and representation of large

quantities of information is also a major challenge in the management of threat sharing platforms.

Haass et al. performed a case study for threat information sharing to identify challenges and problems in organizational, technical and legal domains. Experiments indicated that there is a real need to reduce the number of false positives. Importance of multi-sector sharing to improve threat visibility is also discussed.

For sharing threat information, a lot of research has already been done in structuring information by proposing different data formats and transport mechanisms like STIX, IODEF, TAXII, etc.

Chapter 9

CONCLUSION

In this thesis, we introduced — iGen, a novel system for automatic extraction of IOCs from different sources of unstructured data. Our approach starts from collection of security data from diverse sources. Then, it cleans the data using advanced NLP techniques. CNN based IOC identification technique is found to be highly effective, immensely outperforming the top industry IOC analyzer tool in terms of accuracy, precision and recall. iGen generates ready-to-deploy rules from these IOCs that can be directly deployed on security infrastructure like NIDS etc.

iGen uses different sources for data to make fully automated cyber threat intelligence gathering more collaborative. iGen has collected around 400K IOCs till now with a precision of 95%, better than any state-of-art method.

REFERENCES

- [1] “Alien Vault Open Threat Exchange”, URL <https://otx.alienvault.com/> (2017).
- [2] “APTNotes database”, URL <https://github.com/aptnotes/data> (2017).
- [3] “Beautiful Soup”, URL <https://www.crummy.com/software/BeautifulSoup/> (2017).
- [4] “Common Attack Pattern Enumeration and Classification”, URL <https://capec.mitre.org/> (2017).
- [5] “CozyDuke Apt Report”, URL <https://www.f-secure.com/documents/996508/1030745/CozyDuke> (2017).
- [6] “Cyber Observable eXpression”, URL <https://cyboxproject.github.io/> (2017).
- [7] “Kaspersky White Papers”, URL <http://usa.kaspersky.com/enterprise-security/resources/white-papers> (2017).
- [8] “Malware Attribute Enumeration Characterization”, URL <https://maec.mitre.org/> (2017).
- [9] “Malware Information Sharing Platform (MISP)”, URL <http://www.misp-project.org/> (2017).
- [10] “The mitre corporation”, URL <http://www.mitre.org/> (2017).
- [11] “Openstack”, URL <https://www.openstack.org/> (2017).
- [12] “PDFMiner”, URL <http://www.unixuser.org/~euske/python/pdfminer/> (2017).
- [13] “Snort rules”, URL <https://www.snort.org/> (2017).
- [14] “Structured Threat Information eXpression”, URL <https://stixproject.github.io/> (2017).
- [15] “Suricata rules”, URL <https://suricata-ids.org/> (2017).
- [16] “Symantec blog”, URL <https://www.symantec.com/connect/blogs> (2017).
- [17] “Tensorflow Api docs”, URL https://www.tensorflow.org/api_guides/python/nn (2017).
- [18] “The OpenIOC Framework.”, URL <http://www.openioc.org/> (2017).
- [19] “Trusted Automated eXchange of Indicator Information”, URL <https://taxiiproject.github.io/> (2017).

- [20] “Verizon’s 2016 Data Breach Investigations Report”, Tech. rep., URL <http://www.verizonenterprise.com/verizon-insights-lab/dbir/2016/> (2017).
- [21] “Yara rules”, URL <https://virustotal.github.io/yara/> (2017).
- [22] Abadi, M., P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu and X. Zheng, “TensorFlow: A system for large-scale machine learning”, Google Brain p. 18, URL <http://arxiv.org/abs/1605.08695> (2016).
- [23] Andress, J., “Working with Indicators of Compromise”, in “Journal Information Systems Security Association (ISSA) 5”, (2015).
- [24] Bayer, U., I. Habibi, D. Balzarotti, E. Kirda and C. Kruegel, “A view on current malware behaviors”, Proceedings of the 2nd USENIX conference on Large-scale exploits and emergent threats: botnets, spyware, worms, and more p. 8, URL <http://portal.acm.org/citation.cfm?id=1855676.1855684> (2009).
- [25] Bengio, Y., R. Ducharme, P. Vincent and C. Janvin, “A Neural Probabilistic Language Model”, The Journal of Machine Learning Research **3**, 1137–1155 (2003).
- [26] Bojarski, M., D. Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. Jackel, M. Monfort, M. Zhang, X. Zhang, J. Zhao and Z. K., “End to end learning for self-driving cars”, in “arXiv preprint arXiv:1604.07316”, (2016).
- [27] Catakoglu O., B. M. and D. Balzarotti, “Automatic extraction of indicators of compromise for web applications”, (In WWW, 2016).
- [28] Claessen, K., J. Duregard and M. Palka, “Generating Constrained Random Data with Uniform Distribution”, Functional and Logic Programming, Lecture Notes in Computer Science **8475**, pp. 18–34 (2014).
- [29] Collobert, R., J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu and P. Kuksa, “Natural Language Processing (Almost) from Scratch”, Journal of Machine Learning Research **12**, 2493–2537 (2011).
- [30] Cortes, C., Vapnik and V., “Support-vector networks.”, Machine Learning pp. pp. 273–297 (1995).
- [31] Daly, M. K., “Advanced persistent threat”, Usenix, Nov **4**, 4, 2013–2016 (2009).
- [32] Dandurand, L. and O. S. Serrano, “Towards improved cyber security information sharing”, in “Cyber Conflict (CyCon), 2013 5th International Conference on”, pp. 1–16 (IEEE, 2013).
- [33] Darling, M., G. Heileman, G. Gressel, A. Ashok and P. Poornachandran, “A lexical approach for classifying malicious urls”, in “High Performance Computing & Simulation (HPCS), 2015 International Conference on”, pp. 195–202 (IEEE, 2015).

- [34] Dos Santos, C. N. and M. Gatti, “Deep convolutional neural networks for sentiment analysis of short texts.”, in “COLING”, pp. 69–78 (2014).
- [35] Friedman, J. H., “Stochastic gradient boosting”, *Computational Statistics and Data Analysis* **38**, 4, 367–378 (2002).
- [36] Friedman, N., D. Geiger and M. Goldszmidt, “Bayesian network classifiers.”, *Machine Learning* pp. pp. 131–163 (1997).
- [37] Goodfellow, I. J., D. Warde-Farley, M. Mirza, A. Courville and Y. Bengio, “Max-out Networks”, in “Proceedings of the 30th International Conference on Machine Learning (ICML)”, vol. 28, pp. 1319–1327 (2013).
- [38] Han, E.-H. S., G. Karypis and V. Kumar, “Text categorization using weight adjusted k-nearest neighbor classification”, in “Pacific-asia conference on knowledge discovery and data mining”, pp. 53–65 (Springer, 2001).
- [39] Hinton, G., L. Deng, D. Yu, G. E. Dahl, A.-r. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath *et al.*, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups”, *IEEE Signal Processing Magazine* **29**, 6, 82–97 (2012).
- [40] Hinton, G. E., N. Srivastava, A. Krizhevsky, I. Sutskever and R. R. Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors”, *ArXiv e-prints* pp. 1–18, URL <http://arxiv.org/abs/1207.0580> (2012).
- [41] Hu, B., Z. Lu, H. Li and Q. Chen, “Convolutional neural network architectures for matching natural language sentences”, *NIPS* pp. pp. 2042–2050 (2014).
- [42] Joachims, T., “Text categorization with support vector machines: Learning with many relevant features”, in “Proceedings of the 10th European Conference on Machine Learning”, *ECML ’98*, pp. 137–142 (Springer-Verlag, London, UK, UK, 1998), URL <http://dl.acm.org/citation.cfm?id=645326.649721>.
- [43] Kalchbrenner, N., E. Grefenstette and P. Blunsom, “A convolutional neural network for modelling sentences”, *arXiv preprint arXiv:1404.2188* (2014).
- [44] Kim, Y., “Convolutional neural networks for sentence classification”, in “Proceedings of the Empirical Methods in Natural Language Processing”, (2014).
- [45] Krizhevsky, A., I. Sutskever and G. Hinton, “ImageNet classification with deep convolutional neural networks”, *NIPS* (2012).
- [46] Lewis, D. D., “Naive (bayes) at forty: The independence assumption in information retrieval”, in “European conference on machine learning”, pp. 4–15 (Springer, 1998).
- [47] Liao, X., K. Yuan, X. Wang, Z. Li, L. Xing and R. Beyah, “Acing the IOC Game : Toward Automatic Discovery and Analysis of Open-Source Cyber Threat Intelligence”, in “CCS ’16 Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security”, pp. 755–766 (2016).

- [48] Liao, Y. and V. Vemuri, “Use of K-Nearest Neighbor classifier for intrusion detection.”, *Computers and Security*. **21**, 5, pp. 439–448 (2002).
- [49] Lunt, T. F., “A survey of intrusion detection techniques”, In: *Computers & Security* **12**, 4, pp. 405–418 (1993).
- [50] Lytinen, S. L. and A. Gershman, “Atrans: Automatic processing of money transfer messages”, in “Proceedings of the Fifth AAAI National Conference on Artificial Intelligence”, AAAI’86, pp. 1089–1093 (AAAI Press, 1986), URL <http://dl.acm.org/citation.cfm?id=2887770.2887949>.
- [51] McMillan, R., *Open Threat Intelligence* (2013), URL <https://www.gartner.com/doc/2487216/definition-threat-intelligence>.
- [52] Mikolov, T., I. Sutskever, K. Chen, G. Corrado and J. Dean, “Distributed Representations of Words and Phrases and their Compositionality”, in “Proceedings of Neural Information Processing Systems”, (2013).
- [53] Murdoch, S. and N. Leaver, “Anonymity vs. trust in cyber-security collaboration”, in “Proceedings of the 2Nd ACM Workshop on Information Sharing and Collaborative Security”, WISCS ’15, pp. 27–29 (ACM, New York, NY, USA, 2015), URL <http://doi.acm.org/10.1145/2808128.2808134>.
- [54] Nair, V. and G. E. Hinton, “Rectified Linear Units Improve Restricted Boltzmann Machines”, *Proceedings of the 27th International Conference on Machine Learning* , 3, 807–814 (2010).
- [55] Ng, A. Y., “Feature selection, l1 vs. l2 regularization, and rotational invariance”, in “Proceedings of the Twenty-first International Conference on Machine Learning”, ICML ’04, pp. 78– (ACM, New York, NY, USA, 2004), URL <http://doi.acm.org/10.1145/1015330.1015435>.
- [56] Obrst, L., P. Chase and R. Markeloff, “Developing an ontology of the cyber security domain”, in “CEUR Workshop Proceedings”, vol. 966, pp. 49–56 (2014).
- [57] O’Gorman, G. and G. McDonald, *Ransomware: A growing menace (Technical report)* (Symantec Corporation, 2012).
- [58] Oktavianto, D. and I. Muhandianto, *Cuckoo Malware Analysis* (Packt Pbl. Ltd, 2013).
- [59] Pennington, J., R. Socher and C. D. Manning, “Glove: Global vectors for word representation”, in “Proceedings of the Empirical Methods in Natural Language Processing”, (2014).
- [60] Ranzato, M., F. J. Huang, Y. L. Boureau and Y. LeCun, “Unsupervised learning of invariant feature hierarchies with applications to object recognition”, in “Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition”, (2007).

- [61] Rieck, K. and P. Laskov, “Detecting unknown network attacks using language models”, in “International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment”, pp. 74–90 (Springer, 2006).
- [62] Rumelhart, D., G. Hinton and R. Williams, “Learning Representations by Back-Propagating Errors”, *Nature* **323**, pp. 533–536 (1986).
- [63] Sekar, R. and P. Uppuluri, “Synthesizing Fast Intrusion Prevention/Detection Systems from High-Level Specifications”, in “Proceedings of the USENIX Security Symposium”, (1999).
- [64] Srivastava, N., G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting”, *JMLR* **15**, pp. 1929–1958 (2014).
- [65] Stone, Z., T. Zickler and T. Darrell, “Autotagging Facebook: Social network context improves photo annotation”, in “2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops”, (2008).
- [66] Weston, J., S. Chopra and K. Adams, “# tagspace: Semantic embeddings from hashtags”, (2014).
- [67] Yih, W., K. Toutanova, J. Platt and C. Meek, “Learning discriminative projections for text similarity measures”, *Proceedings of the Fifteenth Conference on Computational Natural Language Learning* pp. 247–256, URL <http://dl.acm.org/citation.cfm?id=2018965> (2011).
- [68] Zhang, Y. and B. C. Wallace, “A sensitivity analysis of (and practitioners’ guide to) convolutional neural networks for sentence classification”, *CoRR* **abs/1510.03820**, URL <http://arxiv.org/abs/1510.03820> (2015).
- [69] Zhu, Z. and T. Dumitras, “Featuresmith: Automatically engineering features for malware detection by mining the security literature”, in “Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security”, pp. 767–778 (ACM, 2016).
- [70] Zou, H. and T. Hastie, “Regularization and variable selection via the elastic net”, *Journal of the Royal Statistical Society. Series B: Statistical Methodology* **67**, 2, 301–320 (2005).

APPENDIX A
LIST OF LINKS TO WEBSITES

Some quick links to the project website, and to this document:

- Project Repository
 - URL: <http://10.90.90.103/threatintelligenceanalytics/iGen-CNN>
- Links to this Document on the Web
 - URL: <http://10.90.90.103/threatintelligenceanalytics/iGen-Paper>