



# UNIVERSITY OF TWENTE.

Faculty of Electrical Engineering,  
Mathematics & Computer Science

Faculty of Behavioural,  
Management & Social Sciences

## Detecting Combosquat Domains using Active DNS Measurements

## Communication of Incident Severity between Customers and Analysts in a SOC

Joost Jansen  
Combined MSc Thesis  
August 2019



**FOX IT**  
part of nccgroup

---

### Supervisors:

prof. dr. M. Junger (chair) (IEBIS/BMS)  
prof. dr. ir. A. Pras (DACS/EEMCS)  
O. van der Toorn MSc (DACS/EEMCS)  
S. Rog MSc (Fox-IT/MSS)  
M. van Hensbergen MSc (Fox-IT/TI)

---



# Preface

In front of you lies the combined master thesis that was written to finalize two master's programmes attended at the University of Twente; Computer Science (CS) with specialization 4TU Cyber Security and Business Information Technology (BIT) with specialization IT Management & Innovation. Although this combined thesis is written as a whole, **it contains two distinct parts that reflect the two distinct programmes**. For the CS part, it was investigated whether a detection model for com-bosquat domains was possible based on active DNS measurements. For the BIT part, the communication between analysts and customers in a Security Operations Center was analyzed, specifically focusing on possible differences in the perception of 'severeness' of incidents between these two groups. I have started working on the CS part of the thesis from September 2018, and from February 2019 on I have worked on the BIT part. When I first approached Fox-IT in the spring of 2018, and got in contact with Christian & Krijn, I could only hope that things would turn out the way they did. Together with the graduation committee, which consists of prof. dr. M. Junger, prof. dr. ir. A. Pras and O. van der Toorn MSc, multiple research topics were discussed that eventually resulted in the two subjects included in this thesis.

I would like to thank Christian & Krijn for helping me find a graduation assignment and their monthly feedback. Furthermore, I would like to thank Martin, Sanne & Ruud for their helpful feedback and guidance on a daily basis, countless colleagues at Fox-IT who made it fun to go to Fox-IT every day, and from whom I learned a lot over the past year. Special thanks for my graduation committee, who regularly provided really helpful feedback and provided me with new insights when I got stuck. Besides the feedback on this thesis, I really enjoyed the conversations we had about all kinds of topics related to cyber security. Lastly, I would like to thank my family, friends and my girlfriend Dorian who have always motivated me, especially at times when I needed it the most.

I hope you will enjoy reading this combined thesis and that you will gain new insights that can be used to make the (digital) world more secure!

Joost Jansen

Delft - August 12, 2019

Since Part II has been marked Confidential for the next five years after the hand-in at the University of Twente library, Part II has been removed from this Computer Science thesis submission. However, for e.g. accreditation reasons the full version of the combined thesis (including Part II) can be requested at the University of Twente library. The full combined version will be publicly released around 12-08-2024.

# Summary

Domain squatting is a phenomena where attackers register domains that mimic popular domains and/or trademarks, in order to trick people into believing they are visiting a legitimate website. A distinct form of domain squatting is combosquatting; adding one or more words to an existing domain/trademark to craft a new domain. Think of *http://utwente-login.nl* as a combosquat domain for the original domain *utwente.nl*. A literature study revealed that a lot of research was performed in the field of malicious domain detection, however not specifically tackling the problem of combosquatting domains. Given this information, combined with the active DNS measurements available from the OpenINTEL project, a research was initiated that aimed at creating model to detect these combosquat domains.

At first, it was investigated whether a generic detection model for combosquat domains existed. After a validation, implementation and evaluation phase involving a ground truth dataset of 10.548 labeled domains, it became clear that no generic fingerprint of combosquat domains could be created given the data that was available. This led to the conclusion that it is extremely difficult to construct a generic model for detecting combosquat domains without a predefined list of trademarks.

The next part of the research focused on the lifecycle of combosquat domains, more specifically in which stages of the killchain they reside and which features could be used to determine when a combosquat domain turns into a malicious state.

Finally, a model that was trained on the information from the sub-questions was designed and validated in a real-world context. The results showed that the detection of combosquat domains turning malicious based on active DNS measurements is not sufficient. Future work includes the use of additional data sources and a bigger responsibility for registrars.



# Contents

<b>Preface</b>	<b>iii</b>
<b>Summary</b>	<b>v</b>
<b>List of acronyms</b>	<b>xi</b>
<b>1 General Introduction</b>	<b>1</b>
1.1 Thesis structure . . . . .	3
<b>I Detecting Combosquat Domains using Active DNS Measure- ments</b>	<b>5</b>
<b>2 Introduction</b>	<b>7</b>
2.1 Introduction into Cyber Threat Detection . . . . .	7
2.2 Motivation . . . . .	8
2.3 Requirements . . . . .	9
2.4 Research Questions . . . . .	10
2.5 High level approach . . . . .	10
<b>3 Background Information</b>	<b>11</b>
3.1 Domain Name System . . . . .	11
3.2 Detecting Malicious Domains . . . . .	15
3.3 Domain Squatting . . . . .	17
<b>4 Approach</b>	<b>21</b>
4.1 Main research question . . . . .	21
4.2 Generic model design . . . . .	27
4.2.1 Problem investigation . . . . .	27
4.2.2 Treatment design . . . . .	27
4.2.3 Treatment validation . . . . .	30
4.2.4 Treatment implementation . . . . .	31

4.2.5	Implementation evaluation . . . . .	31
4.3	Domain lifecycle analysis . . . . .	32
4.4	Feature selection . . . . .	33
<b>5</b>	<b>Generic model design</b>	<b>37</b>
5.1	Designing a ground truth . . . . .	37
5.2	Defining features . . . . .	37
5.2.1	Defining internal features . . . . .	38
5.2.2	Defining contextual features . . . . .	39
5.3	Prototype construction . . . . .	41
5.4	Prototype validation . . . . .	43
5.5	Real world validation . . . . .	43
5.6	Discussion . . . . .	45
<b>6</b>	<b>Domain lifecycle analysis &amp; feature feature selection</b>	<b>49</b>
6.1	Defining the killchain phases . . . . .	49
6.1.1	Discussion . . . . .	52
6.2	Feature selection . . . . .	56
6.2.1	Discussion . . . . .	57
<b>7</b>	<b>Detection model design &amp; validation</b>	<b>59</b>
7.1	Treatment design and validation . . . . .	59
7.2	Treatment implementation and evaluation . . . . .	59
<b>8</b>	<b>Conclusion</b>	<b>63</b>
8.1	Conclusion . . . . .	63
8.2	Future work & recommendations . . . . .	65
	<b>References</b>	<b>67</b>
	<b>Appendices</b>	
<b>A</b>	<b>Generic model data</b>	<b>71</b>
A.1	Manual trademark selection . . . . .	71
A.2	Frequent combosquatting words . . . . .	71
A.3	Public blacklists . . . . .	72
<b>B</b>	<b>Python code snippets</b>	<b>73</b>
B.1	OpenINTEL queries . . . . .	73
B.2	Validation of real-world domains . . . . .	74
B.3	Running the prototype, training & test phase . . . . .	75



---

<b>C</b>	<b>Scraped blacklists</b>	<b>83</b>
<b>D</b>	<b>Trademark distribution</b>	<b>85</b>
D.1	Total days before blacklisted . . . . .	85
D.2	Trademark frequency . . . . .	85



# List of acronyms

<b>2LD</b>	Second-Level Domain
<b>APT</b>	Advanced Persistent Threat
<b>CDM</b>	Combosquatting Detection Model
<b>CKC</b>	Cyber Kill Chain
<b>CTD</b>	Cyber Threat Detection
<b>DNS</b>	Domain Name System
<b>ML</b>	Machine Learning
<b>MSS</b>	Managed Security Services
<b>OSINT</b>	Open Source Intelligence
<b>RR</b>	Resource Record
<b>TI</b>	Threat Intelligence
<b>TLD</b>	Top-Level Domain



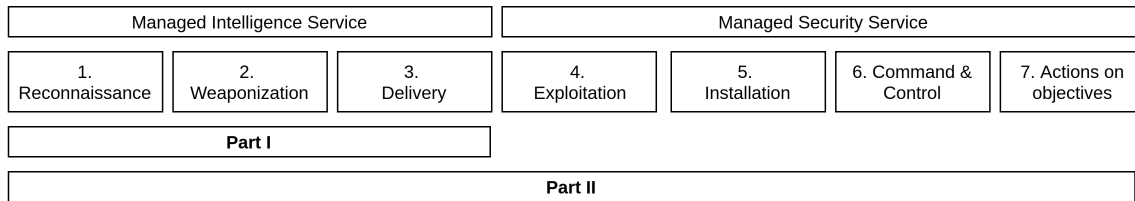
# **General Introduction**

In the current information age, being connected to the internet is a major part of life. Social media, online shopping, movie streaming; all examples of services that are frequently used on the modern-day internet. Despite providing a lot of convenience to people, these services also bring unwanted side-effects [1]. Privacy and identity are at stake when personal data gets compromised, stolen credit card information may result in illegitimate transactions and ransomware may infect one's device, encrypt the files and ask for ransom.

Not only individuals are targeted in this harsh world of cyber crime. Businesses and governments are attacked on a regular basis by a diverse set of attackers. These attackers can be individuals looking for personal gain, hackers attacking for propaganda purposes, organized crime groups looking for financial benefits or even military cyber units, disrupting and degrading an adversary's capabilities [2]. Although it is hard to accurately calculate the global cost of cyber security, several models estimate the costs into hundreds of billions of US dollars. In the future, the global costs of cybersecurity will grow even more. Therefore, the academic community as well as business around the world are providing knowledge, methods or services to minimize the impact of cybercrime.

On the one hand, businesses need to protect themselves in advance to minimize the risk of a being a victim of cybercrime. On the other hand, a business should also acknowledge that 100% security can never be achieved and thus, sooner or later the business might become victim.

Fox-IT [3] provides managed & professional services to protect businesses and governments against cyber attacks. The Threat Intelligence (TI) department is responsible for detecting and reporting threats on external networks, that is, out on the internet (e.g. DDoS attacks, phishing campaigns, Advanced Persistent Threat (APT)). More specifically, the Managed Intelligence Service of Fox-IT automatically collects Open Source Intelligence (OSINT) and performs the corresponding analysis/triage for customers. The Managed Security Services (MSS) department is



**Figure 1.1:** Relation of both parts of the thesis to the CKC and Fox-IT departments

responsible for detecting and reporting attacks and other anomalies on an internal network and individual endpoints (e.g. propagating viruses & trojans, insider crime). Part of this is the Security Operation Center, which actively monitors customers' networks and reports if anything suspicious is discovered. New technical developments and the emergence of new malicious actors in the cyber crime domain force Fox-IT to evaluate and improve their products on a continuous basis.

Since this thesis covers two separate but related parts, it's useful to plot the scope of both parts on some scale. A common way to do this in the academic & professional community is to use an attack modelling technique. An attack modelling technique is useful to understand the characteristics of an attack and the objectives of the attackers, in order to gain information about how and when the attack can be stopped. Over the years, several attack modelling techniques have been developed [1], [4], one of which is the Cyber Kill Chain (CKC). The CKC that has been introduced by Lockheed Martin [5] provides 7 common stages of an attack. This CKC will be a connecting thread throughout the thesis so that at any stage of detecting & reporting an incident, a reference to the CKC can be made. Within Fox-IT, the CKC is also widely adapted. In Figure 1.1, the 7 stages of the CKC are shown in relation to two parts of the thesis and the aforementioned Fox-IT departments in order to provide a high-level overview of the thesis.

The CKC consists of 7 stages, which are briefly described below:

1. *Reconnaissance*: An attacker searches for any publicly available information on the victim, in order to prepare the attack.
2. *Weaponization*: An attacker selects the malicious payload that will be sent to the victim. This payload usually contains code that is capable of performing some action on the victims machine, e.g an .exe file on Windows machines.
3. *Delivery*: An attacker delivers the malicious payload to the victim using any communication medium, e.g. providing a link to download the payload or attaching the payload to an email message.
4. *Exploitation*: The victim accidentally or deliberately stores the malicious payload on the victims machine.

5. *Installation*: The malicious payload on the victims machine gets executed, either automatically or by the victim performing some action.
6. *Command and control*: An attacker creates a communication channel to the victims machine to control the status and remotely execute commands. At this stage, the attacker is in control of the victims machine.
7. *Action on objectives*: An attacker performs the actions required to achieve his/her goals on the victims machine or the connected network. From the victims machine, the attacker can also launch a new attack to achieve a goal that requires more access.

At any stage in the killchain the attack can be stopped. In general, the earlier the attack is stopped in the killchain, the less damage is inflicted on the victim. The CKC will be addressed as the 'killchain' throughout the thesis for legibility.

## 1.1 Thesis structure

The thesis is split into two distinct parts: Part I covers the Computer Science-focused research on combosquat domains and subsequently, Part II covers the BIT-focused research on the communication of incident severity in a Security Operations Center. The reasoning behind this is the fact that besides the overlap of literature and killchain stages of the parts, the research objectives and methodology is unique for each part. Furthermore, each part consists of a distinct research question and corresponding subquestions. Still, both parts have a similar structure, as outlined in Table 1.1.

Part	I	II
Introduction	Chapter 2	N/A
Background information	Chapter 3	N/A
Methodology	Chapter 4	N/A
Results	Chapter 5 Chapter 6 Chapter 7	N/A
Conclusion	Chapter 8	N/A

**Table 1.1:** Chapter structure of the thesis





## **Part I**

# **Detecting Combosquat Domains using Active DNS Measurements**



# Introduction

This introductory chapter will provide information about the current state of research into Cyber Threat Detection (CTD). Section 2.1 will provide an introduction in the research field of CTD and explain some of the basic (technical) concepts related to CTD. Section 2.2 describes the shortcomings in current literature and outlines the ideas that led to the motivation of this research. Next, in Section 2.4 the research questions that are used in this research are displayed. Finally, in Section 2.3 the requirements for the detection model are discussed.

## 2.1 Introduction into Cyber Threat Detection

The CTD research field covers a lot of subjects [6]. From a technical perspective, CTD looks at e.g. email indicators (email traffic, attachments, subject lines), host-based indicators (malware hashes, binaries, DLL's, registry keys) and network indicators (malicious URLs and domain names) to discover new or emerging threats. Furthermore, OSINT information can be used to identify potential threats in a stage long before it develops into an actual incident. CTD based on network indicators is a subject that has been researched for quite some time and is used to detect a wide range of threats in an early stage. The detection of malicious domain names implies the involvement of the Domain Name System (DNS), the backbone of the internet. DNS is a globally used protocol & system. In short, DNS translates human-readable domain names (e.g. *utwente.nl*) into IP-addresses used in the worldwide TCP/IP infrastructure (e.g. *130.89.3.249*). The aforementioned domain names are first registered and configured through DNS, before attackers can make use of it. This provides DNS with the unique opportunity to detect malicious domains in the earliest stages of the killchain (*weaponization* and *delivery*) and prevent the attackers from exploiting their malicious domain name.

A subset of malicious domains are 'squatting' domains. These domains are de-

signed in such a way that end users are tricked into believing they are connecting to a legit domain. Domain squatting exists in many forms (see Section 3.3), one of them being combosquatting. In short, combosquatting is the act of combining one or more arbitrary words with an existing trademark, to craft a seemingly legitimate domain. An example would be `utwente-login.nl`, which acts like a login-page for the University of Twente but in fact passes on these credentials to an attacker. Although the loss of university credentials might not be the end of the world, imagine losing credentials for a service authorized to perform financial transactions.

A study on combosquatting by Kintis et al. [7] suggests that combosquatting domains are currently observed 100 times more than typosquat domains. The lack of a generic model for combosquatting domains contributes to this problem, according to Kintis et al. In their large-scale empirical study, they furthermore found out that combosquat domains often remains undetected for a long period of time and that the abuse of combosquat domains is increasing by the year. An analysis of the attacker's usage of combosquat domains resulted in a list of many different forms of abuse, e.g. phishing, social engineering and affiliate abuse. Because of their findings, they called for more research into combosquatting domains & abuse.

It should be clear that malicious combosquat domains pose a serious threat to the mostly uninformed end users. Though some forms of domain squatting have been thoroughly studied, research into combosquatting is (besides the study by Kintis et al) still in its infancy. This leads us to the motivation of this research.

## 2.2 Motivation

The recent study by Kintis et al. called for the urge of further research after performing an empirical study and finding out that combosquatting domains are a growing threat. Before that, only one other empirical study was performed on combosquatting, which is an industry whitepaper published in 2008 [8]. Moreover, no actual detection models have been proposed in current literature. Kintis et al. state that in comparison to typo squatting, for which detection models have been created and validated, combosquatting lacks a generic model because of its nature; there are infinite amounts of possible combinations. While at the first glance this statement may look credible, no proof is provided for this claim. A frequently heard and simple solution to this would be a trademark search on newly registered domain names. This would however result in a lot of false positives, as Kintis et al. already stated. As an example, imagine the Dutch bank ING; if a substring search is performed on 'ing', the domains 'burgerking.com' and 'bing.com' would be flagged as malicious. Since investigating positives is a costly & time-consuming activity, a more refined model using multiple features is needed to minimize the amount of false positives.

Furthermore, this trademark search on new domain names limits the detection of combosquat domains to the trademarks included in a predefined list. This raises the desire for a generic model that is capable of detecting combosquatting domains regardless of the trademark involved.

For several types of domain abuse, not limited to domain squatting, DNS data can be used to create detection models. A recent study on the detection of snowshoe spam using active DNS data resulted in several characteristics that were useful to detect these domains [9]. Therefore, the question arose whether a model based on active DNS data could also be created for combosquatting domains.

Since the study by Kintis et al. (only) covers an empirical research, no real detection methods have been proposed for combosquat domains. Since attackers seem to be able to keep these domains off the blacklist for a long period of time, the early detection of combosquat domains is beneficial. While Kintis et al. present a temporal analysis of combosquat domains and their presence on blacklists, they do not specify which changes in DNS resource records correlate to the addition on a blacklist. The suspicion arises that a change of IP addresses related to the combosquat domain may be an indicator of a combosquat domain turning malicious. The majority of the detection models described in literature use *passive* DNS data, implying that the attackers have set-up and abused a malicious domain before it is detected. The latter also applies to the TI platform of Fox, which is fed by *passive* data on current attacks, e.g. malicious domains appearing in phishing campaigns. This means that businesses are warned only after the attack has been successfully set up. The desire is to actively detect combosquatting domains in order to warn the involved businesses in an earlier stage and thus, reduce the impact of attacks involving combosquatting domains.

Concluding, the motivation was to design a detection model that was better able to detect combosquatting domains using active DNS data, when compared to existing generic detection models.

## 2.3 Requirements

Section 3.2 gives an overview of detection methods that have been proposed to detect a wide array of domain abuse, using different data sources. These studies achieve False Positive rates of as low as 1%, and precision rates as high as 98%. However, no reference percentage is known for specifically detecting combosquat domains with active DNS measurements. Since this is an experimental study the lowest False Positive rate and highest precision rate as possible are desired, but as a bare minimum for the model to be of practical use, the following requirements are set:

*REQ1: The model should have a False Positive rate of at most 5%.*

*REQ2: The model should have a precision rate of at least 90%.*

## 2.4 Research Questions

Because there is a need for detecting combosquatting domains in an earlier stage, this will be the main focus of the thesis. An additional challenge is the fact that seemingly, it is difficult to design a generic model that can distinguish combosquat domains from legitimate domains. This leads us to the following technical research problem:

**CTD:** *How to develop an active combosquatting detection model that meets the requirements set in Section 2.3, so that businesses can be warned in an earlier stage within a threat intelligence platform?*

In order to provide an answer to this problem, first some sub questions have to be answered, which provide the basis for the Combosquatting Detection Model (CDM) design.

**CTD1:** *Is it possible to construct a generic model for detecting combosquat domains?*

**CTD2:** *In which stages of the killchain can a combosquat domain reside?*

**CTD3:** *Which features define the transitions between killchain stages?*

## 2.5 High level approach

This research approach is based on the engineering cycle, described in [10]. The main research question is a design problem; a treatment needs to be designed in order to solve a problem in a specific context. In this case, the treatment to be designed is a model that is able to detect combosquat domains as they turn malicious, and the context is a threat intelligence platform.

Furthermore, the first sub-question about a generic model for combosquat domains is also a design problem, which has to be solved separately in the early phase of the research.

After all sub-questions have been answered, the approach consists of constructing a prototype, placing it in a model of the intended context and apply some scenarios to observe the responses, after which the model can be validated. These requirements are described in Section 2.3. It should also be noted that these requirements are not definitive; they may change based on the outcomes of the design cycles.

# Background Information

This chapter will provide background information on detecting squatted domains: domains that impersonate existing trademarks in order to let members of the public think the domain name legitimately belongs to the existing trademark. The DNS, discussed in Section 3.1, is an OSI Layer 7 level protocol & global system that seems inseparable with the detection of malicious domains. Afterwards, Section 3.2 will provide insight in the latest research into malicious domain detection in general. In Section 3.3 will this literature study will continue, but will be focused on squatted domain specifically.

## 3.1 Domain Name System

The Domain Name System DNS is a global protocol & system that is a major part of the internet [11]. Its most basic function is to translate human readable domain names (e.g. `people.utwente.nl`) into IP-addresses (e.g. `130.89.252.58`) used by the global TCP/IP network. DNS is set up as a distributed, hierarchical client-server system to ensure scalability and high availability. A client can perform a DNS query, which will result in a response from a resolver containing the requested information. Throughout this background information, the `people.utwente.nl` domain and corresponding IPv4-address `130.89.252.58` will be used as a reference.

A domain name consists of multiple levels, separated by a single dot.

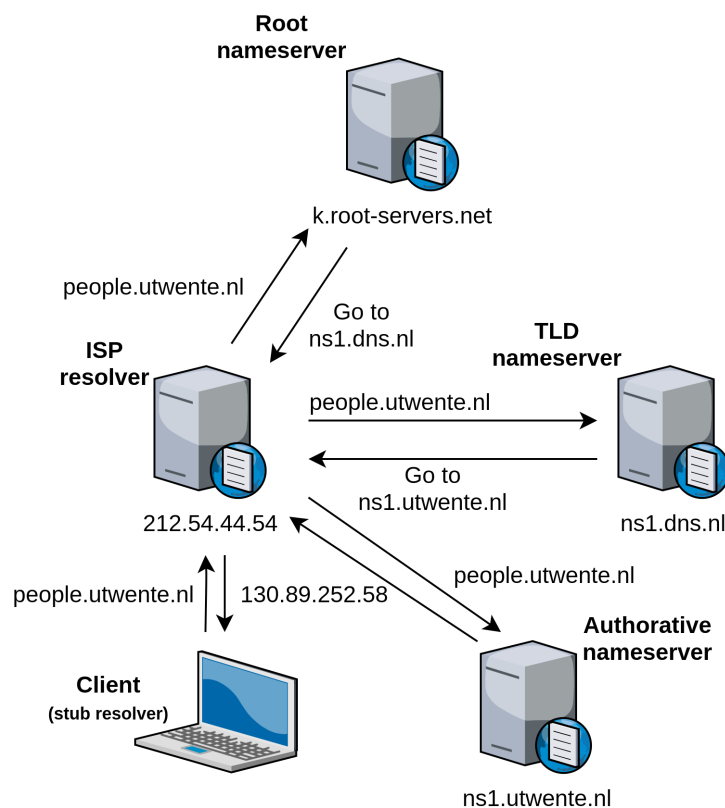
The **Top-Level Domain (Top-Level Domain (TLD))** is the rightmost level of the domain name; often called the domain suffix. TLD's can further be split up into country-code TLD's (ccTLD) and generic TLD's (gTLD). Where ccTLD's are allocated to specific countries (e.g. `nl` for The Netherlands), gTLD's are not bound to a country and are thus independent. In the example, `nl` is the TLD.

After the TLD, the domain is further identified by the **Second-Level Domain (2LD)**.

This level of the domain usually corresponds to the organization that has registered the domain name. Every individual person or business can register a 2LD under a TLD, only bound by regulation regarding trademark names. In the example, the 2LD is `utwente`.

After the 2LD, many more domains levels are possible. In the example, `people` refers to a 3LD. This makes it possible for organizations to have multiple 3LD or even 4LD domains for different services within the organization (e.g. `ftp.utwente.nl` for a FTP-server, `www.utwente.nl` for a webserver)

The functions and components of DNS will be described according to the DNS resolve sequence shown in Figure 3.1.



**Figure 3.1:** An example iterative DNS resolve for *people.utwente.nl*

### Resolver

A resolver is the client-side application in the DNS system; it is responsible for initiating and finishing a full DNS query. Resolvers come in different forms; a simple stub resolver is a piece of software that can check whether the answer to a DNS query is available locally, or can pass the query onto another resolver. This second resolver can be hosted at the user's ISP (e.g. `212.54.44.54`); it may be a more complex system and can perform more difficult tasks. The resolver can make iterative or recursive requests to nameservers, which are explained next.



Since iterative requests are mostly used in resolvers, this type of request will be discussed. In Figure 3.1, the client's stub resolver cannot find the IP-address of `people.utwente.nl` locally, so it forwards the request to the recursive resolver. The ISP resolver is now responsible for returning the query to the client. The ISP resolver consecutively queries the root nameserver, TLD nameserver and authoritative nameserver.

### Root nameserver

A root nameserver's function is to refer the resolver to the correct TLD nameserver. A total of 13 root nameservers exist geographically distributed around the globe, each one operated by a separate organization. In the example, the root server on `k.root-servers.net` is queried and responds with a TLD nameserver (`ns1.dns.nl`) for the `nl` TLD, to which the resolver heads next. One could say that the TLD part of the domain name has now been resolved.

### TLD nameserver

The TLD nameserver holds references to all authoritative nameservers for a certain TLD. Usually TLD nameservers for ccTLD's are hosted by state-owned organizations, whereas gTLD nameservers can also be hosted by other organizations. A TLD nameserver responds to a query by providing the authoritative nameserver for a certain TLD. In the example, the nameserver for the `nl` TLD responds to the resolver that one of the authoritative nameservers for `utwente` is located at `ns1.utwente.nl`.

### Authoritative nameserver

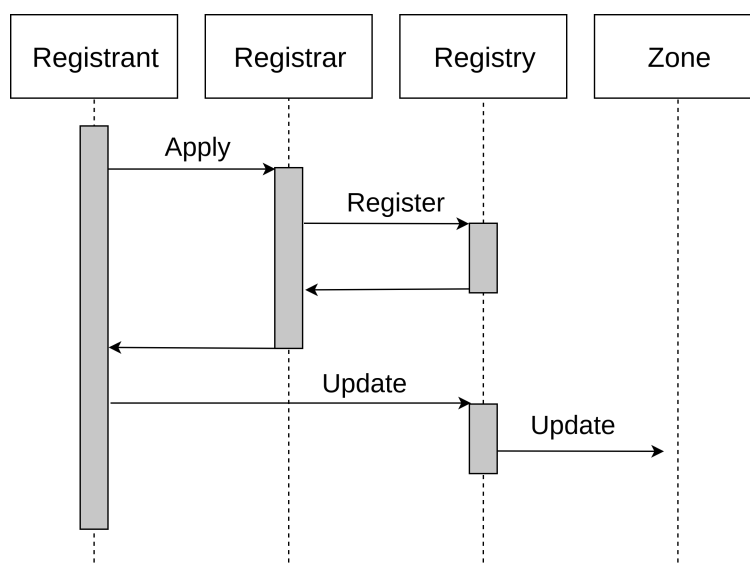
The authoritative nameserver holds all the information for certain TLD, in the form of a Resource Record (RR). These records have a fixed layout, but the information in each of the RR's may be different. The most common RR's are:

Type	Description	Function
A	IPv4 address	Returns the IPv4 address for the domain
AAAA	IPv6 address	Returns the IPv6 address for the domain
MX	Mail exchange	Returns the location of the mail server for the domain
NS	Name server	Returns the authoritative name server for the domain
CNAME	Canonical name	Returns an alias to refer from one domain to another

The authoritative nameserver responds the requested RR's to the resolver, which in turn responds the completed DNS query to the client. In the example, the authoritative nameserver for `utwente`, which is `ns1.utwente.nl`, has a RR of type A for the `people` TLD; `130.89.252.59`. Note that to the client, the process is recur-

sive; the client only has to perform one action, after which the resolver iteratively queries the different nameservers on behalf of the client.

From an organizational perspective, there are three major roles when it comes to managing & registering new domain names. Figure 3.2 displays the sequence diagram provided in the paper by Kidmose et al. [12], along which the different roles will be explained.



**Figure 3.2:** Sequence diagram for registering a new domain name

### Registrant

The registrant is an individual person or organization that wants to register a 2LD. In the example, this is the University of Twente. Usually, a registrant registers a domain name at a registrar.

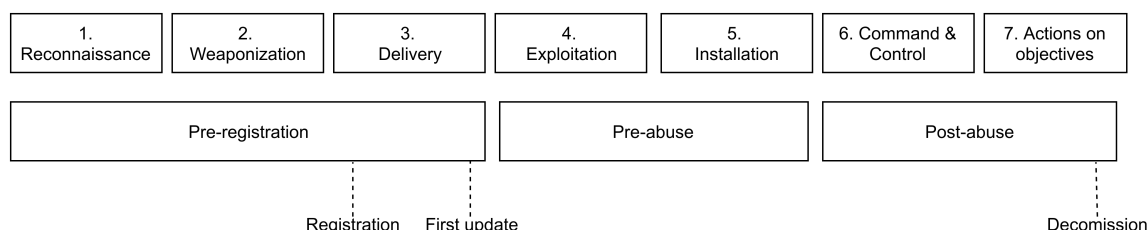
### Registrar

A registrar is an organization that sells domain names on behalf of one or more registries. Registrars are typically webhosting providers or businesses providing internet services. Registrars have direct contact with the registries and function as an intermediary agency for the registrant.

### Registry

A registry is the operator of a TLD and is responsible for taking care of the technical aspects of that operation. Furthermore, it takes care of meeting the requirements set by the ICANN and making the TLD available for commercial use. The latter is most often outsourced to the registrars. In the example, the `nl` TLD is operated by SIDN [13].

As can be seen in Figure 3.2, a registrant applies for a domain name at a registrar. The registrar passes the request on to the registry, which makes sure the domain name does not violate any of the abuse rules. The registry approves and charges a fee for the registration; the registrar on its turn charges a fee to the registrant. Now that the registration is complete, the new domain gets included in the *zone file* update of the registry. Kidmose et al. call this stage of the registration the *pre-registration* stage. After the update is published, the new domain name gets propagated over the other nameservers on the internet and can be resolved everywhere around the globe as displayed in Figure 3.1. This stage is called the *post-registration* phase. The *post-registration* can, in regard to abuse, be split into a *pre-abuse* and *post-abuse* stage. Combining these different stages in the domain name registration process with the CKC introduced in Chapter 1 results in Figure 3.3.



**Figure 3.3:** DNS registration process combined with CKC

## 3.2 Detecting Malicious Domains

A starting point for the literature review on the detection of malicious domains is the study performed by [12]. The authors perform an analysis on existing frameworks and theories. In addition, the study by Zhauniarovich et al. [14] provides a systematic review of malicious domain detection approaches based on this DNS data. This is a good starting point for enumerating state of the art research into malicious domain detection using DNS data. Several studies discussed in these overviews will now be discussed. EXPOSURE [15] is a passive DNS analysis service. It focuses on detecting DGA & C&C domains. The service was built and tested on billions of DNS requests, and during the 17 months of operation it detected over 100.000 malicious domains using a J48 decision tree algorithm. The classifier used multiple feature categories; time series based, DNS answer based, TTL value based and domainname based features were used. During the evaluation phase, the service managed to achieve a high detection rate on the training data (99.5%), as well as a low False Positive rate (0.3%).

Phoenix [16] is a system that focuses on detecting and distinguishing DGA and non-

DGA domains, and furthermore is able to find groups of DGA-generated domains that are used alongside in botnets. The system uses passive DNS data. It uses a combination of linguistic and IP-based features to do this 'fingerprinting' of botnet DGA domain groups. During the evaluation phase, the system was able correctly distinguish DGA-generated domains from non-DGA-generated domains in 94.8% of the cases. Furthermore, the system was also able to detect these DGA domain groups in a real-world setting.

DFBotKiller [17] is another system that is able to detect botnet traffic to malicious C&C servers. This is done by analyzing passive DNS data within the network. Its main task is to assign a negative reputation score to a domain, that takes into account three suspicious measurements. These three metrics are then, along with a number of failed DNS queries, processed into a verdict. The system was evaluated in a test settings and resulted in interesting scores.

A system that acts in the domain pre-registration stage is the PREDATOR framework [18], which has a detection rate for new malicious DNS entries of 70%, in combination with a low false positive rate of 0.35%. This means that as early as in the pre-registration phase, a majority of the malicious domains can already be identified before they can be abused. The system uses features based on registrar data, characteristics of the domain name, previous registration history and correlation with registration bursts. This data is used, since no active nor passive DNS data is present in the pre-registration phase.

The study performed by Vissers et al. [19] is interesting since it provides an automated clustering process that analyzes the registration of malicious registrations in any TLD during the pre-registration stage. Using DNS registration data and publicly available blacklists, they found that at least 80.04% of the data corresponded to 20 malicious DNS registration campaigns. A remaining 19.30% of the traffic could also be related to these campaigns, after more rigorous inspection of the individual campaigns features, resulting in a false positive rate of only 0.92%.

Finally, another study outlines the malicious domain registration ecosystem and puts the different types of abuse in perspective regarding absolute numbers and total costs [20]. This study does not propose a concrete detection system, but gives insights into a malicious actor's preferences and economic incentives.

Having covered most of the existing theory, Kidmose et al. state that future research should be into detecting malicious & abusive domains in the *pre-registration* stage and should not be limited to spam domains. They suggest to use new, currently unused, features in this stage of detection, namely a *lexicology analysis of a domain name*, the *registration history of domain name*, the *registrant information*, *contents of first zone update* and the *reputation of the registrar*.

### 3.3 Domain Squatting

A specific type of domain name abuse is called is domain squatting. In the case of domain squatting, an attacker registers a domain name that appears to be the legitimate domain name of a trademark, while in fact it hosts some malicious or abusive content. Domain squatting is particularly hard to detect since it involves no technical errors or flaws in the DNS protocol; in the end it is up to the user to spot a 'squatted' domain. Several types of domain squatting can be identified, each type with its own characteristics and features. To illustrate the different types of domain squatting, Figure 3.4 displays some examples. Additionally, the domain squatting types are briefly explained in Figure 3.4, with combosquatting outlined more in detail.

Type	Example	Literature
Typosquatting	utwent.nl	[21], [22]
Bitsquatting	utwenpe.nl	[23]
Homophone-Based squatting	youtwente.nl	[24]
Homograph-Based squatting	utvvente.nl	[25], [26]
Abbrevsquatting	ut.nl	[27]
Combosquatting	utwente-login.nl	[7], [8]

**Figure 3.4:** Different types of domain squatting targeting *utwente.nl*

#### Typosquatting

Typosquatting is a type of domain squatting where the domain names consist of typo variations of popular websites. This method of domain squatting requires an end user to make a mistake when entering a domain name in the browser. In the example, a user wants to visit *utwente.nl*, but accidentally forgets to type the *e* in the domain name. The user then ends up at a completely different website, which could be used for malicious purposes.

Typosquatting is phenomena that has been around for many years. The study by Wang et al. [21] from 2006 already presents a tool that is able to detect and monitor typosquat domains. In their study they list five different forms of typos, ranging from a missing dot typo (e.g. *wwwutwente.nl* to the character-omission typo mentioned in the example. With respect to combosquatting, it is worth mentioning that the amount of typosquat domains for a given popular domain is fixed; at a certain point, no further variations can be computed.

#### Bitsquatting

Bitsquatting is a type of domain squatting where the attacker anticipates on random bit-errors originating from the hardware in client devices (e.g. *comput-*

ers & smartphones). The study by Nikiforakis et al. [23] shows that these domains are actively being registered and used for abuse purposes. End users are often unaware of being redirected to a malicious website, since they have not performed a faulty action but rather are the victim of hardware errors and the attackers who anticipated for this. In the example, the users requests `utwente.nl`, but due to a random bit-error the client actually resolves `utwenpe.nl`.

### **Homophone-Based squatting**

In another study by Nikiforakis et al. [24], homophone-based squatting or squatting based on words that sound exactly like the original domain, but are written in a distinctive matter. In the example, imagine someone telling the end user to visit the 'utwente' website, which then misinterprets the URL as `youtwente.nl`. Another example would be `weather.com` and `whether.com`; two URL's who sound exactly the same but would resolve to different hosts.

### **Homograph-Based squatting**

Homograph-based squatting is a type of domain squatting where an attacker registers a domain that is visually (almost) indistinguishable from the original popular domain. Research on this topic has been done by Holgers et al. [25]. An example attack (using this thesis' font) would be replacing a Latin lower case letter `l` with a number `1`; this would make `paypal.com` hard to distinguish from `paypa1.com`. A study that used homograph-based domain squatting was conducted at the University of Twente, where the original domain was replaced by `utvvente.nl` [26]; in this study the double `v`'s in the domain name were impersonating a `w`. Since the adoption of International Domain Names, which allow non-ASCII characters to be used, homograph-based squatting has become harder to detect since many Latin letters have similar looking characters in different alphabets.

### **Abbrevsquatting**

Abbrevsquatting is a type of domain squatting where an attacker uses the abbreviation of popular domains to trick users into believing they are visiting the legitimate website. In the example, an attacker registers `ut.nl` because the University of Twente is often abbreviated as UT. The study by Lv. et al [27] shows that attackers are aware of the principles of abbrevsquatting and are already leveraging them in malicious ways.

### **Combosquatting**

Combosquatting is the act of combining one or more arbitrary words with an existing trademark, to craft a seemingly legitimate domain. The first research

into combosquatting dates back to 2008, when an industry whitepaper by Fair-Winds Partners, LLC was published [8]. An initial set of 30 trademarks was selected based on their "strength" and number of search terms that were regularly associated with the trademarks. A keyword suggestion tool was used to generate the top 50 most popular keywords associated with the trademarks, and these were then combined into a total of 1500 domain names. Using the major search engines at that time (Google, Yahoo! and MSN), the daily searches for these domain names were analyzed per month. Furthermore the traffic for each domain name was registered. Afterwards, the domains were ordered by their traffic/search ratio and the top and bottom 500 were excluded, leaving 500 domains for manual testing purposes. Results showed that 50.6% domains contained Pay Per Click (PPC) advertisement, 22% of the domains were legitimately used for trademarking purposes and 75% of the trademark+keyword combinations that were not owned by the trademark contained PPC advertisements.

Nine years later, in 2017, the study by Kintis et al. [7] was published. This was and until now is the only academic research into combosquatting. For the first time, they introduce a definition of a combosquat domain; the domain contains a trademark and the domain cannot result by applying the five typosquatting models of Wang et al. [21]. Furthermore, they perform an empirical study on the presence of combosquat domains on the internet, using several large-scale datasets.

They conclude that current domain squatting detection techniques are not detecting combosquat domains properly due to the different threat models involved. They also state that no generative model can be constructed, since in theory an infinite amount of trademark+keyword combinations exist. This makes detection harder to perform than for example typosquatting, for which an exhaustive list of mutations regarding a trademark can be made. A temporal analysis of the detected combosquat domains shows that most domains were active for several months, before the domains were blacklisted. Combosquat domains are used in phishing campaigns, affiliate abuse and other types of abuse. All of these findings result in their call for future research regarding combosquat domains. Kintis et al. argue that not only registrants and registrars can help resolving this threat, but there is also a task for third parties to search for and monitor new combosquat domains.



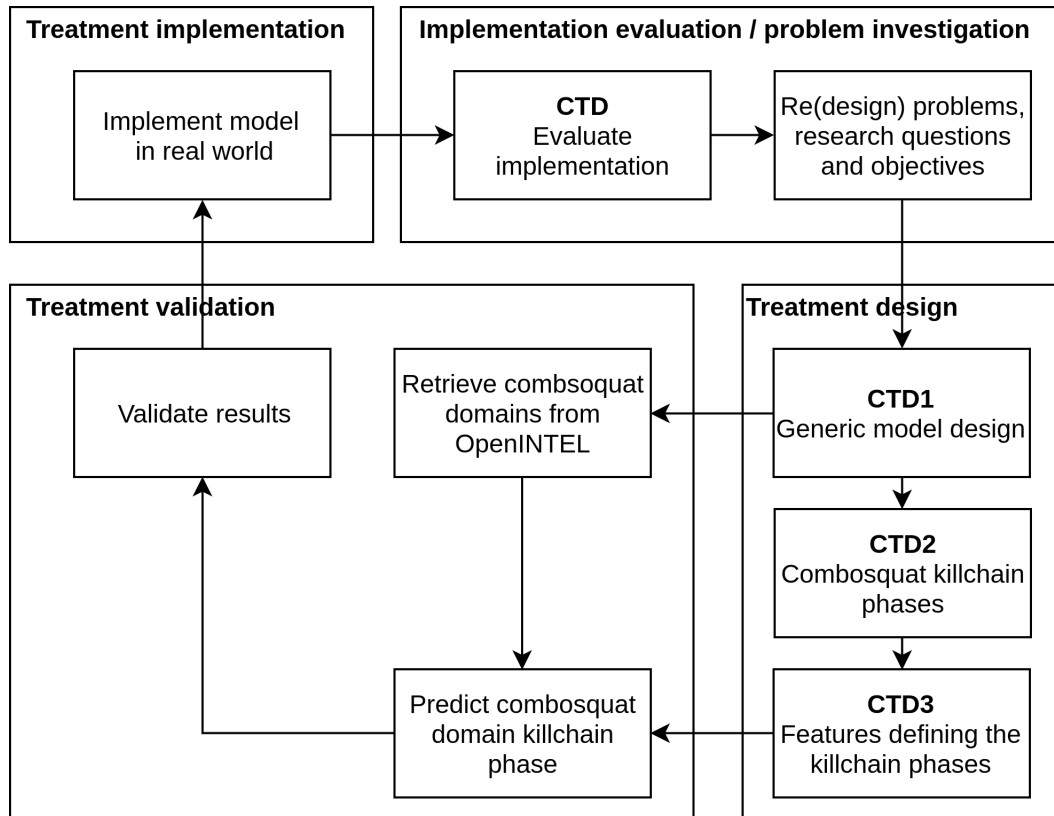


# **Approach**

This chapter describes the approach that was used to provide the answer to the research questions. In Section 4.1 the approach that was used to answer the main research question is explained. In Section 4.2 the approach that was used to find out if a generic model could be designed is explained. Section 4.3 described how the different killchain phases are measured and assigned to the stages a combosquat domain can reside in. Finally, Section 4.4 describes the approach to find the most relevant features for the detection model.

## **4.1 Main research question**

Following the problem statement & objectives in Chapter 1, this chapter will describe the research approach that was used during the research. This provides structure in the research and gives an overview of the steps that were taken. In Figure 4.1 the research approach is displayed. This research approach is based on the engineering cycle, described in [10]. Since this approach focuses on answering knowledge questions and solving design problems, it fitted the needs of this research. The five individual phases of the engineering cycle, interpreted in the context of this research, are outlined below. This approach also keeps in mind the framework provided by [14], which outlines a general framework to design a detection model primarily based on DNS data. Steps included in the framework are data collection, data enrichment, algorithm design & evaluation; these steps will be identifiable in the engineering cycle as well. Before continuing, let us first define a combosquat domain. The definition of a combosquat domain is based on the definition provided by Kintis et al. as shown in Section 3.3, but is extended to fit the needs of this research. Below, the formal definition is outlined alongside a few examples to make the definition more tangible. This definition is expressed in Python code in Section B.2, which is used for validation purposes throughout the thesis.



**Figure 4.1:** Schematic view of the research approach

Domain name C is considered a combosquat domain of trademark T, if:

- 1) T is the original trademark name, without a spelling deviation
- 2) T is left intact within a set of other characters, in this case C
- 3) T is a standalone word in C
- 4) The owners of domain names C and trademark T are different
- 5) C can not be classified as any other form of domain squatting as listed in Section 3.3

Next, the four phases of the design cycle will be explained in the context of this research.

### Problem investigation

During the first stage of the research, the initial problems are defined, research questions are created and objectives are set. This stage is described in Chapter 2.

### Treatment design

In the treatment design phase, the sub questions will be answered. It should be noted that since CTD1 in itself is an extensive design problem it is only answered

once when iterating over the design cycle multiple times. The different methodologies to answer the sub questions are described in the subsections of this chapter. The results from CTD2 (a list of killchain phases that a combosquat domain can reside in) will function as input for CTD3.

### Treatment validation

After the sub questions have been answered, it is time to construct the prototype and validate the results. This is done by checking if the prototype matches the requirements. More specific, this means that it is validated that the classifier can distinguish between features that define a 'benign' combosquat and a malicious combosquat. The validation of the prototype is done by constructing a *confusion matrix*, as shown in Table 4.1. This is a common method to validate results regarding predicted and actual values.

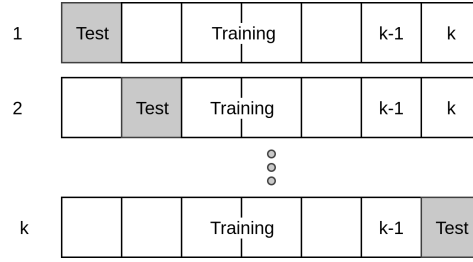
		Prediction outcome		Total
		p	n	
Actual value	p'	True Positive	False Negative	P'
	n'	False Positive	True Negative	N'
Total		P	N	

**Table 4.1:** Example confusion matrix

From the confusion matrix, three important metrics can be calculated:

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

The *FP rate* represents the amount of wrongly predicted positives compared to the total amount of actual negatives. In this research, a low False Positive rate is desirable since every domain that is predicted as combosquat needs to be investigated; when the domain turns out to be a False Positive, the time spent on the analysis is 'wasted'.



**Figure 4.2:**  $k$ -fold cross validation

$$Accuracy = \frac{TP + TN}{P + N}$$

The *accuracy* score represents the amount of correctly classified labels out of the total. A high accuracy score means that the classifier does not make many errors in relation to the total amount of predictions.

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

The *precision* score represents the amount of labels that are correctly labeled positive relative to the wrongly labeled positives. A high precision score means that the classifier makes little mistakes when labeling positives. On the other hand, a low precision score implies that a lot of positives labels are predicted while in facts, these are negative. In this research a high precision score is one of the main requirements, in order to keep the amount of False Positives as low as possible.

A common problem with ML classifiers is overfitting: a classifier performs perfectly on the training data, but it does not perform well on newly, previously unseen data. To get a better picture of the performance of the classifier, a process called *k-fold cross validation* is performed on the training data. This process is shown in Figure 4.2.

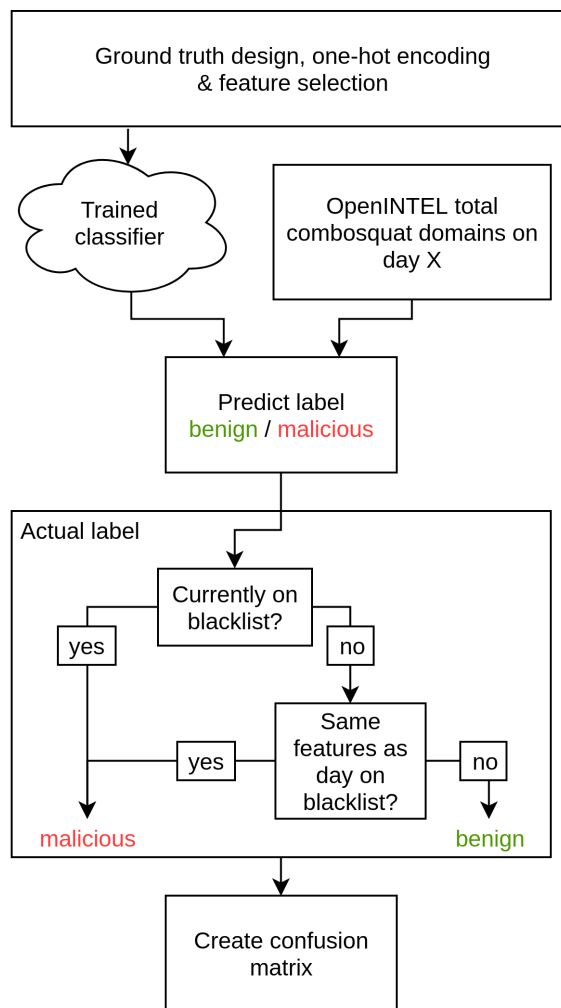
The training data upon which the model is trained is split into  $k$  sections. In each one of the  $k$  iterations, a different section is used for training & testing purposes. This approach has another advantage, namely that the scores for the classifiers are calculated over the total training data and do not rely on a randomly chosen test sample. Each one of the  $k$  iterations produces a *False Positive*, *accuracy* and a *precision* score, which are in the end summed up and divided by  $k$  to get the average scores of the *False Positive rate*, *accuracy* and *precision*.

### **Treatment implementation & Implementation evaluation**

In this phase, the validated prototype is placed in its intended context: a threat intelligence platform. An external Virtual Private Server, functioning as an abstraction of such a platform was chosen for this purpose.

Afterwards, the implementation is evaluated and will answer the main research question. A question to be asked is: Did the CDM setting function according to its requirements in the real world context? In most design researches, the design cycle is iterated over multiple times. When a new iteration is started, this phase redefines the problems, research questions and objectives in order to improve the quality of the CDM. The treatment implementation phase, along with the implementation evaluation is shown in Figure 4.3.

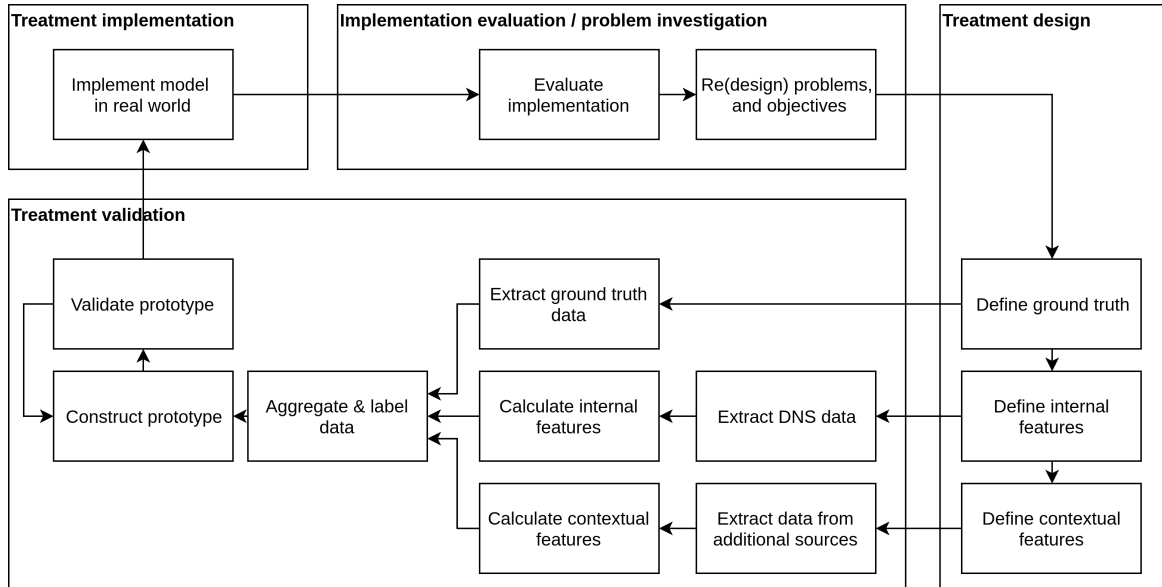
As can be seen, all combosquat domains from a specific day  $x$  are retrieved from OpenINTEL and fed to the classifier, which predicts the domain as either 'benign' or `malicious`. To evaluate this decision, it is checked whether the domain was actually on a blacklist on day  $x$  or not. Note that this validation completely relies on the presence of a domain on a blacklist; if the classifier manages to predict a malicious domain while it was undetected at that moment, it cannot be verified. Therefore, an extra check is performed; if the domain is not listed on a blacklist on day  $x$ , the features of the day it actually got detected are obtained and verified against the features of day  $x$ . If the features match the domain is labeled `malicious` and 'benign' if the features do not match. If no appearance on a blacklist can be observed after the prediction, the domain is also labeled 'benign'.



**Figure 4.3:** Treatment implementation & evaluation

## 4.2 Generic model design

This subsection describes the approach that was used to answer sub question CTD1: *Is it possible to construct a generic model for detecting combosquat domains?* Since this sub question in itself is a design problem, another design cycle has been constructed with the sole purpose of providing an answer to this sub question. This design cycle is shown in Figure 4.4.



**Figure 4.4:** Design cycle used to answer CTD1

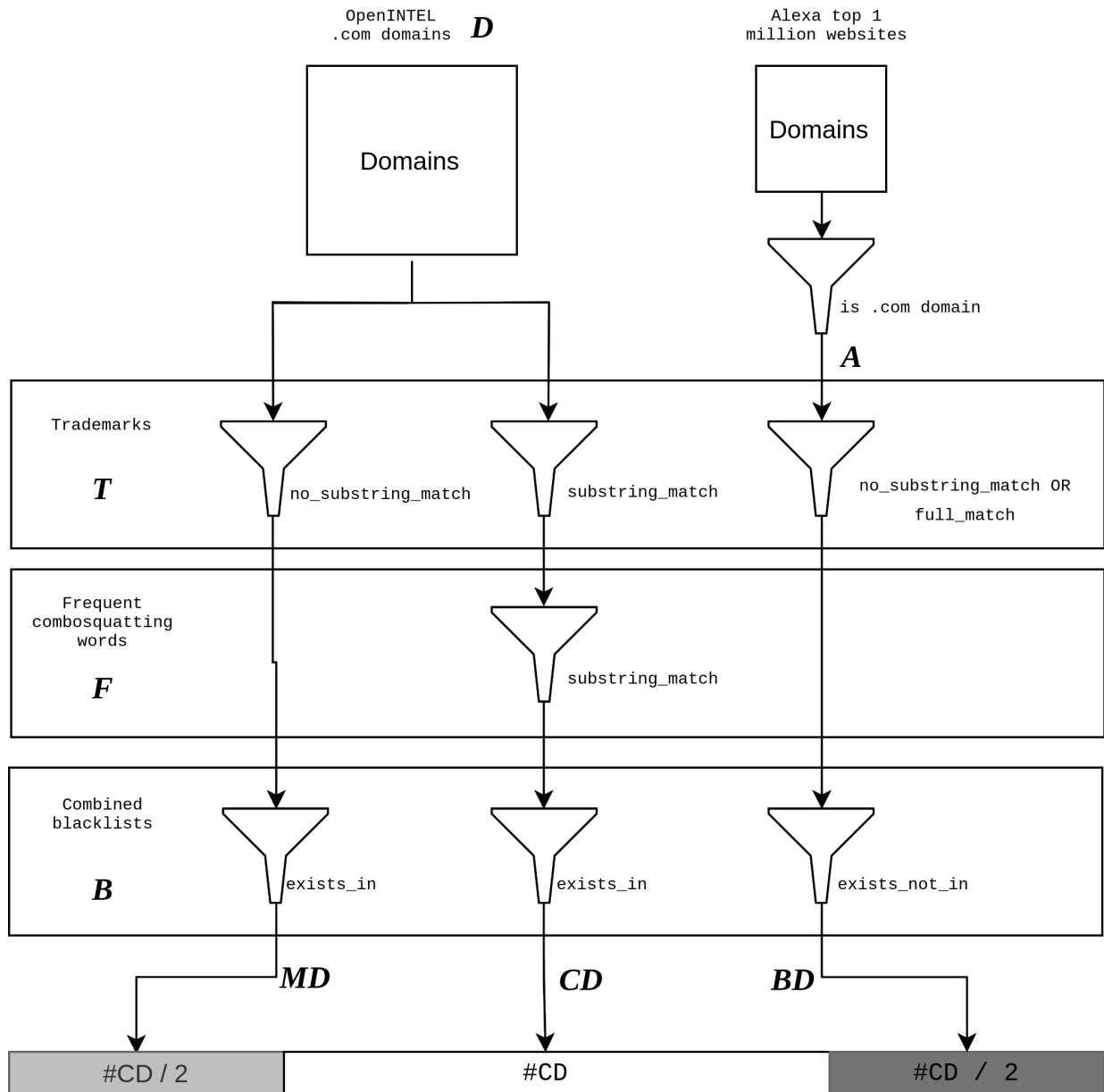
### 4.2.1 Problem investigation

This problem is also described in Chapter 2; it is a direct result of the claim that there is no generic model possible for combosquat domains [7]. Since no proof was provided to support this claim, an attempt was made to design a generic model to check if it was indeed the case.

### 4.2.2 Treatment design

The design of the generic model is based on the approach used by van der Toorn et al. [9]; first a ground truth has to be composed of combosquat domains and non-combosquat domains because labeled data is needed for the model to be based on. Since no labeled dataset was available regarding combosquatting domains, this was created manually using the approach described below.

Before diving into detail regarding the ground truth creation, the different datasets need to be defined. **D** is the set of domain names in the .com TLD. **T** is the set



**Figure 4.5:** Approach to construct the ground truth



of trademarks. This set is based on the global Alexa top 500 domains, retrieved from [28]. Ambiguous domain names are excluded, as well as short names ( $< 4$  chars). This manual selection of domain names resulted in 106 unique domain names, displayed in Appendix A.1. **F** is the set of frequently used combosquatting words, as displayed in the paper by Kintis et al. [7]. **B** is the set of the blacklisted domains. This set is constructed out of the blacklists as shown in Appendix A.3. **A** is the set of .com domains appearing on the Alexa top 1M list [28].

Next, two string operations should be clarified. Consider a string as a sequence of characters, represented as:  $S = c_1, c_2 \dots c_n$ . Then, a *substring*  $B$  is formally defined as  $B = c_{1+i} \dots c_{m+i}$  where  $0 \leq i$  and  $m + i \leq n$ . For example, `ent` is a valid substring of `utwente`, but `twete` is not. In the same way, a *full match* is when two strings are equal; `utwente` is a valid full match of `utwente`.

Now that the different datasets and functions are defined, the ground truth is created using a filtering process. The approach used in this filtering process is shown in Figure 4.5.

First, a list of combosquat domains (**CD**) is created; for each domain in **D**, it is checked if there exists at least one substring match with a trademark from **T** and a word from **F**. Afterwards, for each of these domains, a check is performed whether it is present on a blacklist. If this is the case, then it is added to **CD**. In order to create a proper ground truth that is usable for training & testing purposes, an equally long list of 'non-combosquatting' domains should be appended. This 'whitelist' is built out of two sources. 50% of the 'whitelist' consists of malicious domains which do not contain a trademark, called Malicious Domains (**MD**). The other 50% consists of Benign Domains (**BD**). These domains are first extracted from **A**. Then, the domains containing a trademark from **T** are excluded, except when there is a full match. The rationale behind this is that according to Kintis et al., popular combosquat domains frequently make their way into the Alexa top list. Since the **BD** set should definitely not include any combosquat domains, these popular combosquatting domains are filtered out. For example, `youtube.com` must be included in the set of benign domains, but `youtubedownloader.com` should not. In the end, a ground truth consisting of **#CD** combosquat domains, **#CD/2** 'whitelisted' domains and **#CD/2** malicious but non-combosquat domains is present.

After the groundtruth is constructed, the ground truth needs to be enriched with features. Features are essentially datapoints used to train Machine Learning (ML) classifiers. ML is a technique that enables automated binary classification and prediction of entities and is commonly used to pick new features out of a large data pool. Following the terminology introduced by Zhauniarovich et al. [14], *internal features* are features based on DNS data, while *contextual features* are features based on

additional, non-DNS data.

The main data source for the *internal* features will be historical measurements of domains extracted from the OpenINTEL project. The available variables (for example `domain_name`) must be transformed into features ready to be processed (for example `number_of_characters`, `number_of_digits`). Moreover, the data source may contain a large amount of features and since they are not equally interesting to the model, the most relevant features have to be selected. These features can be found in literature, but also arise as a result of statistical analysis on the ground truth data. A list of most valuable internal features will be the answer to this question. For the purpose of this research the internal features are further split up in *lexical features*, which are extracted from the domain name itself, and *DNS features*, all other features extracted from the OpenINTEL project. The study by Kintis et al. [7] already provides several lexical features that are commonly observed at combosquatting domains, which can be of use.

The *contextual* features are extracted from the Certificate Transparency Log, as well as the WHOIS service. The Certificate Transparency Log is an append-only Merkle-hash tree, which is used to verify the validity of a SSL/TLS certificate. Since new certificates are added on a continuous basis, Google stores information about these certificates and provides reports for all domains [29]. The WHOIS-servers provide information about a domain name through a special WHOIS query. Usually, these queries hold information about for example the registrant, registrar and name-servers. While there are standards in place for WHOIS queries & responses, it is up to the registries and registrars to determine what is inserted in the fields. While the WHOIS-information might not be fully reliable, it may still be a significant feature. Because WHOIS information is also used in relevant literature, it is initially included.

### 4.2.3 Treatment validation

In this phase, the a prototype of the CDM will be constructed and validated. All of this is performed on the ground truth sample dataset, functioning as a model of the intended context. The prototype is validated against the requirements that will be shown at the end of this section. The construction & validation of the prototype is itself an iterative process; by applying small changes to the prototype, the outcome of the validation model will slightly be changed.

The corresponding requirements which the generic model has to meet are specified as follows:

#### **Functional requirements**

*FRQ1:* The CDM should be able to classify a combosquat domain into the correct killchain phase.

*FRQ2:* The CDM should not include any other form of domain abuse other than combosquatting.

*FRQ3:* The CDM should be able to detect combosquat domains without the use of a predefined list of trademarks.

### **Non-functional requirements**

*NRQ1:* The CDM should have a False Positive rate of at most 5%.

*NRQ2:* The CDM should have a precision rate of at least 90%.

*NRQ3:* If either NRQ1 and NRQ2 can be met, NRQ1 should be given a higher priority when selecting a single classifier.

The 5% and 90% percentages are set as a bare minimum. Since no reference percentages are known due to the lack of research into generic combosquatting detection models. The values will however be adjusted according to performance of the first prototype, since a generic model is new ground and realistic threshold values are not known up front.

## **4.2.4 Treatment implementation**

In this phase, the validated prototype is placed in its intended context: a threat intelligence platform. The initial idea was to use Fox's threat intelligence platform for this purpose. In the end a simplified solution was chosen, in which the model was hosted on a Virtual Private Server, with a live connection to the OpenINTEL system. This VPS was used as an abstraction of a real threat intelligence platform.

## **4.2.5 Implementation evaluation**

In this stage, a check was performed whether the model in context met the requirements set earlier in this section. For this research, it meant that the model should be able to distinguish newly added combosquats from newly added non-combosquats according to the requirements. This stage of the design cycle was used to answer directly sub question CTD1; whether it is possible to create a generic model to detect combosquat domains.

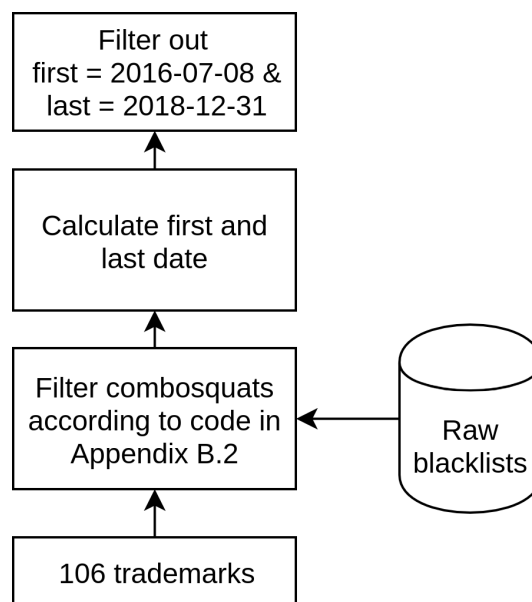
### 4.3 Domain lifecycle analysis

Since a generic model is not available for detecting combosquat domains, at this point the list of 106 trademarks as listed in Appendix Section A.1 is used; combosquat domains impersonating those trademarks are taken into account.

A starting point for answering this subquestion is a collection of blacklists, that have been scraped from 2016-07-08 until 2019-01-11. For the analysis, the blacklists up to 2018-12-31 were included to leave some data untouched for later validation usages. More information about the blacklists that have been scraped can be found in Appendix C.

The blacklisted domains functioned as the starting point for the analysis. First, out of the total domains on the blacklists, the combosquat domains were filtered according to the definition of a combosquat domain provided in Chapter 4. This means that each domain on the blacklist was checked against the five specified requirements, and those who did not match all of the requirements were dropped. The code that performed this combosquat filtering is shown in B.2.

Since in this analysis the complete combosquat domain is being researched, combosquat domains that were present on a blacklist on the first and last day of the selected blacklists were left out.



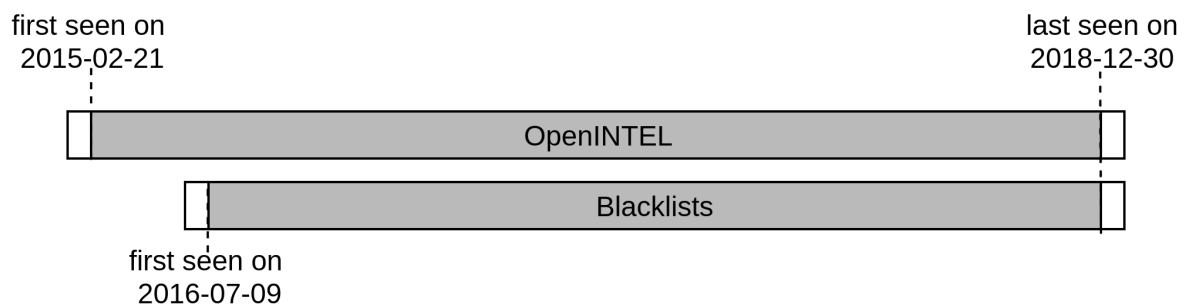
**Figure 4.6:** Blacklist filtering process

This resulted in a list of blacklisted combosquat domains, not present on the first and last day of the blacklist collection period. The domains on this list were then enriched with data from OpenINTEL. For every day, starting on 2015-02-20 (the starting day of OpenINTEL .com measurements) until 2018-12-31 (the fixed end date for

the analysis) the presence of the domain in OpenINTEL was checked. The check was performed by checking if any records for the domain were present in OpenINTEL. It should be noted that a domain can be registered without any record being related to it; if absolutely no records related to the domain are present in the zone file it is not measured by OpenINTEL. In this case, while the domain is registered, it is considered inactive and thus, not present in OpenINTEL. Afterwards, a list of domains and the corresponding first and last day in OpenINTEL is present.

Before continuing, only the domains of which the full lifecycle could be observed were taken into account; domains that were already active at the start of the OpenINTEL measurements, or still active at the last day of the analysis were filtered out. This means that the domains whose first date was 2015-02-20 and/or last date was 2018-12-31 were removed from the dataset.

For convenience, the OpenINTEL data combined with the blacklists is shown in Figure 4.7; the parts marked in grey are included in the dataset.



**Figure 4.7:** Overview of the selected dataset

At this point, a dataset containing only combosquat domains that could be observed throughout their full lifecycle are present. This dataset was used to calculate several graphs and metrics regarding the lifecycle of the domains. The results are described in Section 6.1.

## 4.4 Feature selection

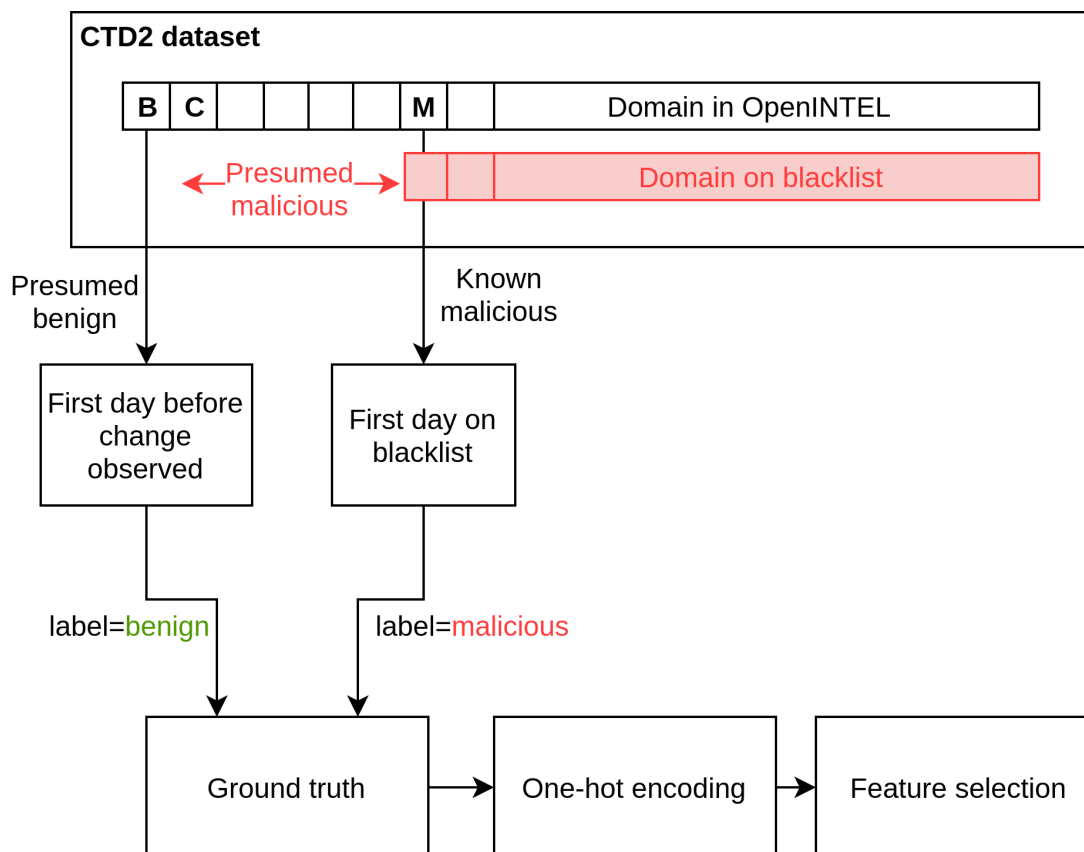
Following the analysis in sub question CTD2, a distinction can be made between combosquat domains in the different killchain phases. More specifically, combosquat domains can be detected in two phases; the *Weaponization* and *Delivery* phase. It was also stated that domains in the *Weaponization* phase are considered benign and domains in the *Delivery* phase are considered malicious.

The dataset containing only combosquat domains that could be observed throughout their full lifecycle, produced by CTD2, is again used and enriched with the features described in Section B.1. This dataset consisted of 13693 domains. All these

domains appeared in OpenINTEL, have gotten blacklisted and disappeared from either OpenINTEL or the blacklist during the measurement period.

Historical DNS & blacklist data was used in favor of creating a new dataset containing more data sources (for example HTTP), since the limited time available for this research would result in a dataset that only contains domains with a lifetime up to a few months, excluding the undetected long-living domains pointed out by Kintis et al. Note that this is only a preliminary study; it is only a starting point to see if additional measurements are feasible.

This means that the label `benign` or `malicious` has to be extracted from the blacklists alone. Since historical data is being processed, it cannot be verified whether the blacklisted domains were actually malicious. This also means that this detection can only be based on changes in the DNS records of that domain. Therefore, **a change in the DNS records of a domain is considered an indicator of an attacker's activity**. Keeping this in mind, a ground truth can be created as shown in Figure 4.8.



**Figure 4.8:** Constructing the ground truth and extracting features

Regarding the lifecycle of a combosquat domain, three important days can be marked. **M** is the day that a domain is known to be **Malicious**; this is the first day the domains appeared on a blacklist. **C** is the day when the **C**hange to the current DNS

settings was observed. This means that the features on day M and C are equal. **B** is the day **B**efore the change was observed; the features of this day differ from M and C. The assumption that was previously made can not be narrowed down; **a change in the DNS records of a domain, with the new DNS records matching the DNS records at the time of blacklisting, is an indicator of a benign domain turning malicious**. The time between B and M is the time to be 'won' with earlier detection; in this period the domain is presumed to be malicious.

The features obtained from OpenINTEL on day B are then labeled as `benign`, while features obtained on day M are labeled as `malicious`. This was then applied to all 13693 domains of whose B and M days, with a ground truth dataset as a result. One-hot encoding is performed to transform the data for Machine Learning, after which the most important features will be extracted by using a *DecisionTreeClassifier*.

Note that in this approach, `combosquat` domains that are being registered and do not change their DNS records before being blacklisted are excluded; since no change can be observed, there are no indicators based on DNS data that the domain is turning malicious. Since OpenINTEL has a resolution of one measurement per day, the DNS data that was analyzed is oblivious to quick changes in the DNS records. Therefore, if a domain has a malicious lifespan of less than one day (or even a few hours) and is pointed to a 'domain is blocked' page afterwards, the DNS records belonging to the 'domain is blocked' page are labeled `malicious` in the ground truth. This problem might be mitigated by investigating HTTP data, but for now this remains future work.





# Generic model design

In this chapter, the outcomes of the generic model design are displayed and discussed. First, in Section 5.1 the ground truth is shown. Afterwards, Section 5.2 provides an overview of the features that were selected and used. The generic model is then designed, validated and implemented in respectively Section 5.3, Section 5.4 and Section 5.5. Finally, the outcomes are discussed in Section 5.6.

## 5.1 Designing a ground truth

According to the approach described in Section 4.2, first all the .com domains having a substring match with both a trademark and a frequent combosquatting word were selected, which resulted in 285.327 domains. Next, the blacklists were retrieved and combined. After filtering only the .com domains out of the the blacklists, as displayed in Appendix Section A.3, 433.757 .com domains were present on the list. out of the 285.327 domains retrieved from OpenINTEL, 5274 were also present in the combined blacklist dataset. According to the approach another 2637 domains were added from the Top Alexa list, along with 2637 blacklisted domains not containing a trademark. In the end, this resulted in a labeled ground truth dataset of 10548, in which 5274 domains were labeled `combosquat` and 5274 domains were labeled `not-combosquat`.

## 5.2 Defining features

In this section, the features that were extracted are described. This is done according to the analogy by Zhauniarovich et al. [14]; *internal features* are features based on DNS data, while *contextual features* are features based on additional, non-DNS data.

Lexical feature	Type	Source
domainname_words	Ordinal	[7]
domainname_segments	Ordinal	[7]
domainname_popular_combosquatting_words	Ordinal	-
percentage_lms_of_total	Ordinal	[15], [18]
domainname_characters	Ordinal	[7], [18]
contains_minus_char	Boolean	[18]
contains_digit	Boolean	[18]

**Table 5.1:** Lexical features (total 7)

### 5.2.1 Defining internal features

For the purpose of this research the internal features are further split up in *lexical features*, which are extracted from the domain name itself, and *DNS features*, all other features extracted from the OpenINTEL project.

#### *Lexical features*

The lexical features that are used are displayed in Table 5.1. There are simple features (contains\_digit, contains\_minus\_char, domainname\_characters) and more complex features. The more complex features are explained below

domainname\_words and domainname\_segments are based on the *word segmentation algorithm*, originally proposed by [30] and also used in the paper by Kintis et al. [7]. Using this algorithm, the domain name is split into different sections based on their probability to be standalone sections. For example, 00fr-youtubevideos would result in the sections 00fr, youtube and videos. Following the classification by Kintis et al., if a section is present in one of multiple dictionaries [31]–[34] it is considered a *word*; otherwise it is considered a *segment*. In the example above, youtube and videos are considered words and 00fr is considered a segment. The two features count the number of these *words* and *segments* in a domain name.

The study of Kintis et al. furthermore provides a list of most frequent combosquatting words per category. All of these words are added to a dictionary set, and if a word matches one of these words the domainname\_popular\_combosquatting\_words value increases by 1. This dictionary of words is displayed in Appendix Section A.2. Finally, the percentage\_lms\_of\_total is calculated as the Longest Meaningful String of the total domain. The largest word in the domain name is selected and calculated as a percentage of the total domain name. In the example above, youtube is the longest meaningful string. Since it has length 7 and the total length of the domain name is 18, the value of percentage\_lms\_of\_total will be 39%.

DNS feature	Type	Source
number_of_A_records	Ordinal	-
number_of_AAAA_records	Ordinal	-
number_of_NS_records	Ordinal	-
number_of_MX_records	Ordinal	-
number_of_SOA_records	Ordinal	-
number_of_CNAME_records	Ordinal	-
number_of_DNSKEY_records	Ordinal	-
number_of_TXT_records	Ordinal	-
number_of_ipv4_addresses	Ordinal	[9], [15]
list_of_ipv4_addresses	Categorical	[15]
AS_number	Categorical	[7], [9], [14]
response_name_matches	Ordinal	[9]
country_code	Categorical	[9], [15]
soa_refresh	Ordinal	-
soa_retry	Ordinal	-
soa_minimum	Ordinal	[9]

**Table 5.2:** DNS features (total 16)*DNS features*

The DNS features are all extracted from the OpenINTEL project. In Section 3.1 background information is provided on the DNS. OpenINTEL stores many DNS record types and DNS fields<sup>1</sup>. Initially, all fields were taken into account when defining the features. Based on relevant literature, a selection of initial fields was made. In addition to that some experimental features were also added, such as the number of certain types of records. The fields were then converted into features, as shown in Table 5.2. It should be noted that this selection is broad; feature selection should later on filter out the less significant features so that only truly distinctive features remain.

Most of these features count occurrences of different types of RR's, IP addresses, AS numbers etc. The `response_name_matches` field indicates whether the query name & response name of a DNS query match.

**5.2.2 Defining contextual features**

Given the definition in [14], contextual features are obtained when DNS data is combined with external data sources. In this research, the external data sources are

<sup>1</sup><https://openintel.nl/background/dictionary/>

WHOIS feature	Type	Source
whois_registrar	Categorical	[15], [18]
whois_number_of_nameservers	Ordinal	-
whois_registrant_name	Categorical	-
whois_registrant_organization	Categorical	-
whois_registrant_country	Categorical	-

**Table 5.3:** WHOIS features (total 5)

CTL feature	Type	Source
ctl_number_of_current_certs	Ordinal	-
ctl_list_of_providers	Categorical	-

**Table 5.4:** CTL features (total 2)

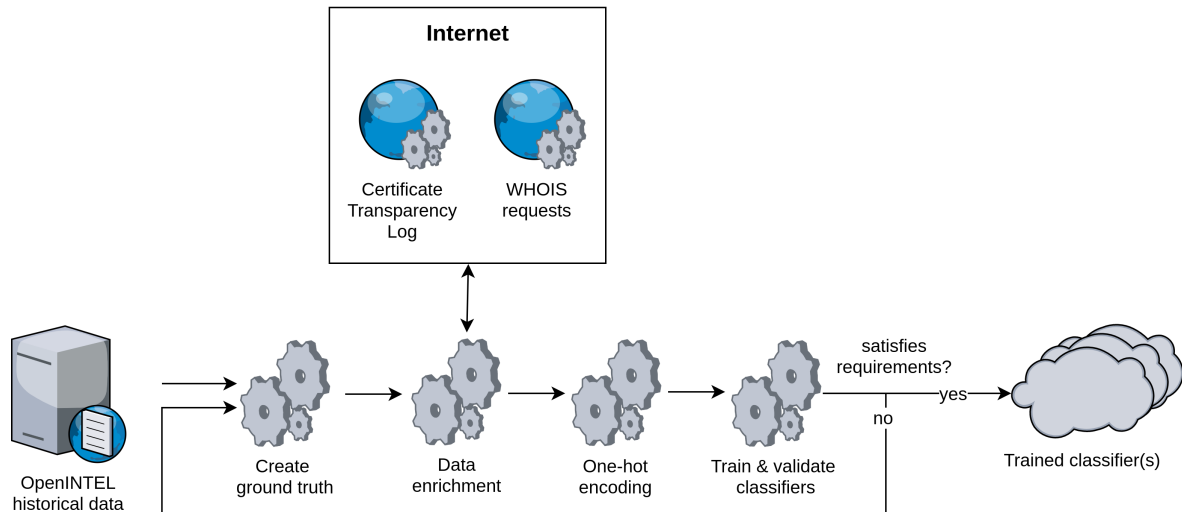
respectively the WHOIS-servers available on the internet and the Certificate Transparency Logs.

### *WHOIS information*

10 features were chosen from the total list of WHOIS-information response fields. Features such as `whois_zipcode` were left out, since they provide too specific information. Similarly some features were left out because they were based on time. For example, the `whois_created_day_ago` held the number of days that had passed since the domain was registered; a feature that is not useful to detect combosquatting domains at any fixed point in time. In the end 5 features were left out, leaving a total of 5 features as shown in Table 5.3.

### *Certificate Transparency Log*

The reports as described in the approach function as data input for the features. For a given domain, the CTL outputs a list of certificates that are currently issued to a domain. Here, the certificates issued to the 2LD and the `www` 3LD are taken into account. These certificates may be issued by multiple providers, for example Let's Encrypt Authority X3 and TERENA SSL High Assurance CA 3. Both the list of certificate providers, as well as the number of current certificates associated with a domain are used as features. If a certificate for a `www` 3LD is issued, the domain is stripped down to its 2LD equivalent. The used features are shown in Table 5.4.



**Figure 5.1:** Schematic overview of the CDM training phase

## 5.3 Prototype construction

After constructing a ground truth and having it enriched with the **30** features described in Table 5.1, Table 5.2, Table 5.3 and Table 5.4, a prototype was constructed in order to be able to validate the performance. The Python<sup>2</sup> programming language was chosen as the main language because of the author's proficiency with the language and the availability of data processing & machine learning libraries, such as `scikit-learn`<sup>3</sup>, `pandas`<sup>4</sup>, `numpy`<sup>5</sup> and `scipy`<sup>6</sup>. In Figure 5.1 and Figure 5.2 a schematic overview of both the training & test phase of the prototype is displayed.

As described earlier, the ground truth resulted in a list of 10548 domains Figure 4.5. Afterwards, the ground truth domains are enriched with the 32 selected features. This involves both invoking the OpenINTEL system to obtain the internal features, as well as querying two public systems to obtain the contextual features. After the data is enriched, a matrix of 10548 rows and 30 columns is present. In the next step, several ML classifiers are trained on the enriched data. ML classifiers can only handle ordinal or binary values. Therefore, the categorical values in the ground truth need to be converted into ordinal or binary values. A technique known as *one-hot encoding* is applied on the data to fulfill this need. One-hot encoding transforms categorical data into binary data, by creating a new column for every categorical value and setting a 1 if the row contains this category, or a 0 otherwise. An

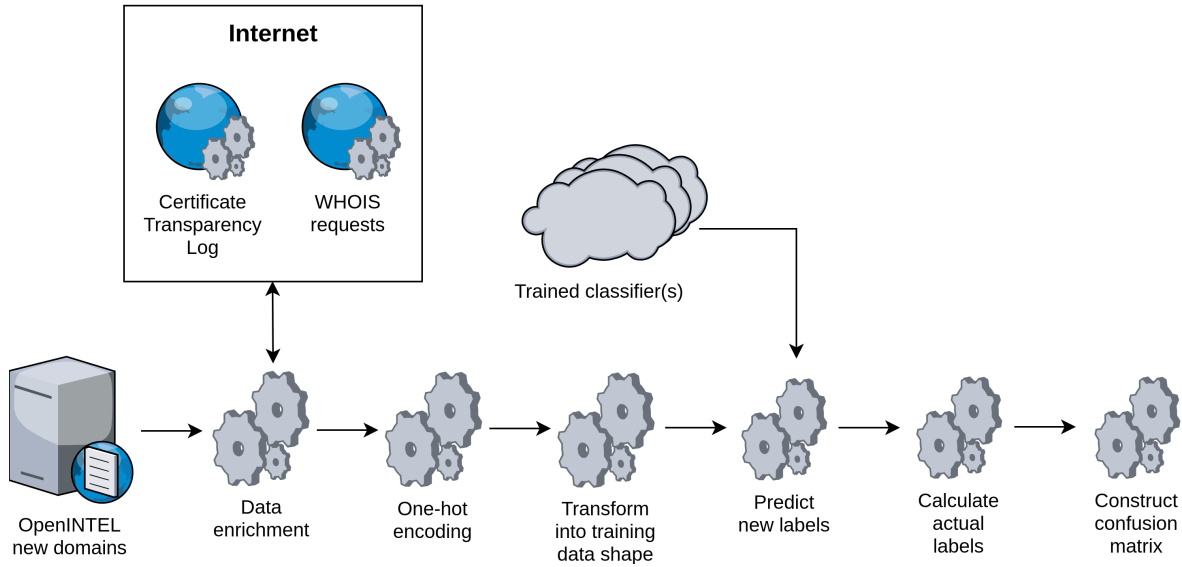
<sup>2</sup><https://www.python.org/>

<sup>3</sup><http://scikit-learn.org/stable/>

<sup>4</sup><https://pandas.pydata.org/>

<sup>5</sup><http://www.numpy.org/>

<sup>6</sup><https://www.scipy.org/>



**Figure 5.2:** Schematic overview of the CDM test phase

domainname	list_of_AS_numbers
example1.com	[12282, 28892]
example2.com	[3853]

**Table 5.5:** Before one-hot encoding

example transformation is shown in Table 5.5 and Table 5.6

In this example, the list in the `list_of_AS_numbers` column is transformed into multiple columns containing only binary information. All categorical features are transformed in this way, except for the `list_of_ipv4_addresses` column. Since there are  $2^{32}$  possible IPv4 addresses, performing one-hot encoding on this column would result in the same amount columns. This is undesirable, since each newly added column increases the memory usage during execution. Therefore, a more efficient way of using IPv4 addresses as features is proposed in the paper by Chiba et al. [35]. Following their transformation based on 'octets', this resulted in 1024 extra columns instead of the possible  $2^{32}$  columns. After the transformations, the matrix consists of 10548 rows and 6106 columns.

For efficiency purposes, a technique called *feature selection* is often applied to a enriched dataset. During feature selection, features that are not used in the clas-

domainname	AS_number_3853	AS_number_12282	AS_number_28892
example1.com	0	1	1
example2.com	1	0	0

**Table 5.6:** After one-hot encoding

sification process are omitted. In this research, feature selection is performed by training a *DecisionTreeClassifiers* with unlimited `max_depth` on the entire dataset and extracting the features that are used. Under the hood, the *DecisionTreeClassifier* tries to minimize the uncertainty based on the Gini impurity of the features. The Gini impurity is based on the probability that a domain is labeled incorrectly based on a certain feature. A high Gini value corresponds to a feature that is significant in the classification process; lower Gini values do not significantly contribute in the decision making process. By default, `sklearn`'s *DecisionTreeClassifier* constructs a tree and keeps adding new leafs to the tree until it reaches a point where it realizes that adding extra leafs is no longer beneficial to the classification process. By constructing such a tree and extracting the features that were used in the tree, not-significant features can be filtered out. Because a *DecisionTreeClassifier* is initialized with a random starting point, the selected features differ per execution round. On average, the total amount of features is reduced from 6106 to around 650, which greatly improves performance. This means that on average, the 650 selected features have the same significance in the classification process as the full set of 6106 features.

## 5.4 Prototype validation

According to the approach described in Section 4.2, the prototype is first validated in a model of its intended context by applying  $k$ -fold cross validation. The value of  $k$  can be arbitrarily chosen, however several sources show that a value of  $k = 10$  is often chosen as a default value<sup>7,8</sup>.

An average of the scores from the 10 iterations is calculated and shown in Table 5.7. These were the first results and no classifier satisfies the requirements; a minimum precision rate of 90% and a maximum *FP rate* of 5%. Our requirements state that if no classifier meets the requirements, the classifier with the lowest *FP rate* should be chosen. In this case, this was the *GaussianNB* classifier with a *FP rate* of 6%.

## 5.5 Real world validation

By retrieving all newly added domains of today, yesterday and the day before yesterday, a list of newly added domains can be obtained from OpenINTEL. The reason that the day before yesterday is included is to make sure that if a measurement error

---

<sup>7</sup><https://magoosh.com/data-science/k-fold-cross-validation/>

<sup>8</sup><https://machinelearningmastery.com/k-fold-cross-validation/>

Classifier	FP rate %	Accuracy %	Precision %
DecisionTreeClassifier	22	78	77
RandomForestClassifier	17	80	81
AdaBoostClassifier	18	80	80
KNeighborsClassifier	26	76	74
GaussianNB	6	51	54
BernoulliNB	26	73	73
MLPClassifier	31	67	69
SGDClassifier	76	48	48
GradientBoostingClassifier	23	81	77
ExtraTreesClassifier	20	81	79

**Table 5.7:** Classifiers and their scores

occured in OpenINTEL (and thus, one measurements for a certain domain is missing), it is not immediately considered a 'newly added domain'. The newly added domains were then enriched with features listed in Table 5.1, Table 5.3, Table 5.4 and Table 5.2. Afterwards, one-hot encoding was applied on the enriched data. However, the one-hot encoded newly added domains cannot be fed directly into the trained *GaussianNB* because of two reasons:

- 1) Categorical values present in the test data but not present in the training data are not known to the classifier.
- 2) Categorical values present in the training data but not present in the test data is missing.

The solution that has been chosen is to first iterate over the test data columns and remove the columns that are not present in the training data. Afterwards, when iterating over the training data columns, the columns that are not yet present in the test data are added and filled with zeroes. Finally, the prototype as described in Figure 5.2 was deployed in a real-world setting.

All newly added .com domains of 26-01-2019 were retrieved. This resulted in a list of 112.647 newly added domains. Subsequently, the trained *GaussianNB* classifier was used to predict the labels of each the new domains. Since the process of retrieving the WHOIS and CT log features is lengthy, the first 10.000 domains were used as a sample for the total set of 112.647 newly added domains.

Out of the total **10.000** domains, this resulted in **75** domains predicted as *combosquat*, and **9925** as *not-combosquat*. Now that the predicted labels are available, new 'actual' labels need to be calculated as well in order to calculate the confusion matrix. Since these domains were not part of the ground truth, the code that is shown in is shown in Appendix B.2 was used to calculate the actual labels. After-



wards, every predicted domain had a 'predicted' label as well as an 'actual' label assigned. These two values could then be used to create the confusion matrix, as shown in Table 5.8.

		Prediction outcome	
		cs	b
Actual value	cs'	TP 0	FN 15
	b'	FP 75	TN 9910
Total		75	9925

**Table 5.8:** Confusion matrix for the real-world validation

Based on the confusion matrix, the *precision*, *FP rate* and the *accuracy* can be calculated as explained in Subsection 4.2.3; these scores indicate how well the prediction on the new unseen data performed:

$$\begin{aligned}
 Precision &= \frac{0}{0 + 75} \\
 &= 0\%
 \end{aligned}$$

$$\begin{aligned}
 FP\ rate &= \frac{75}{75 + 9910} \\
 &= 0.75\%
 \end{aligned}$$

$$\begin{aligned}
 Accuracy &= \frac{0 + 9910}{0 + 15 + 75 + 9910} \\
 &= 99.1\%
 \end{aligned}$$

## 5.6 Discussion

Results indicated that the detection of new combosquat domains was not sufficient. As shown in the previous section, out of the 10.000 new .com domain names, 75 domain names were flagged as combosquatting domains. At first, the frequency of the detected domains is in line with distribution in the ground truth. 75 out of 10.000 means that 0.75% was flagged as combosquat. In the ground truth 285.327 domains

Domain name
sevenoakscommhomes.com
mybackyardrelaxation.com
smarthomegadgetguru.com
anthonyandmikayla.com
silvercloudinvestments.com
retroelectromotors.com
mizikmalemusic.com

**Table 5.9:** Selection of domain names that were falsely labeled as combosquat (False Positives)

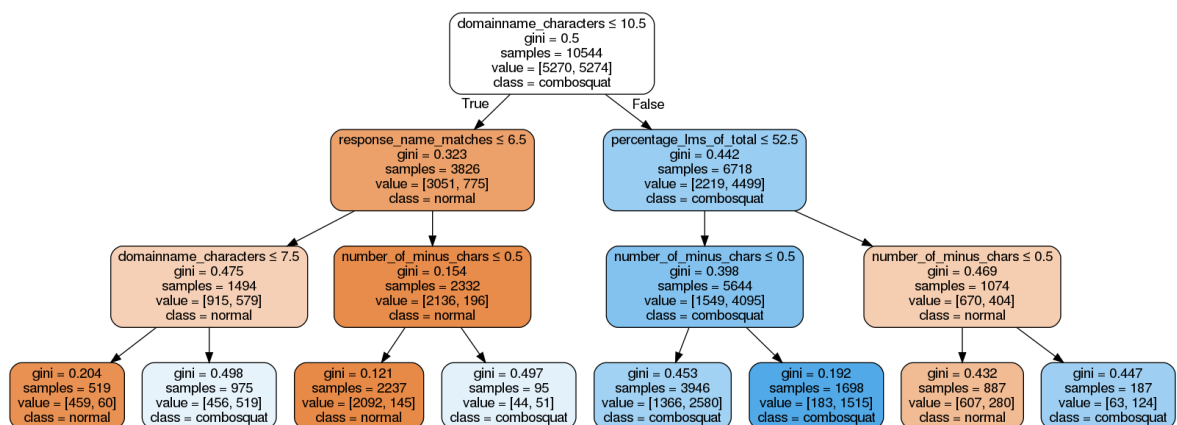
were used as combosquat domains on a total dataset of 137M .com domains; here they make up  $\frac{285327}{137000000} * 100 = 0.2\%$  of the total domains.

The *FP rate* and *accuracy* scores seem pretty satisfying at first. However, the *GaussianNB* classifier was unable to detect even one of the 75 actual combosquats in the 10.000 newly added domains. The *precision* score is therefore 0%, which is obviously not sufficient. When looking more in detail into the trained *GaussianNB* classifier, it becomes clear why this happened.

A selection of the domain names falsely flagged as combosquat is shown in Table 5.9. On the other hand, Table 5.10 shows the combosquat domains that were missed by the classifier. Since the *GaussianNB* classifier is trained based on the features selected by the *DecisionTreeClassifier* shown in Figure 5.3. the tree can be used to explain the results. For convenience, the `max_depth` has been set to 3, since this provides insight in the features with the highest Gini impurity value. The most significant feature is *domainname\_characters*, which represents the total length of the domainname. This means that the prediction is greatly based on this feature. Furthermore, it can be observed that the DNS features are not as important as the lexical features; features like *percentage\_lms\_of\_total* and *number\_of\_minus\_chars* are present, while only one DNS feature is present (*response\_name\_matches*). This means that combosquat domains cannot easily be fingerprinted by distinct DNS entries. The last observation is that the classifier seems unable to distinguish trademarks from regular words. When looking at Table 5.9 and Table 5.10, the classifier did not learn that 'samsung' is a trademark and 'backyard' is not.

Domain name
samsungblockchaincore.com
facebookdekatsuyaku.com
linkedinclassroom.com
shopifywebsitedesignerbuilder.com
godaddyholdings.com
ihategodaddy.com
mywalmartcoupons.com

**Table 5.10:** Selection of domain names that were falsely labeled as not-combosquat (False Negatives)



**Figure 5.3:** Decision tree of depth 3



# Domain lifecycle analysis & feature feature selection

This chapter will display and discuss the results of subquestions CTD2 and CTD3 in respectively Section 6.1 and Section 6.2.

## 6.1 Defining the killchain phases

According to the approach described in Section 4.3, the filtering was performed. 13693 blacklisted domains were available to be analyzed. In this stage a list of 106 trademarks was used and using a simple grouping function, the trademarks that were targeted most frequently could be analyzed. The top 10 of targeted trademarks is shown in Table 6.1.

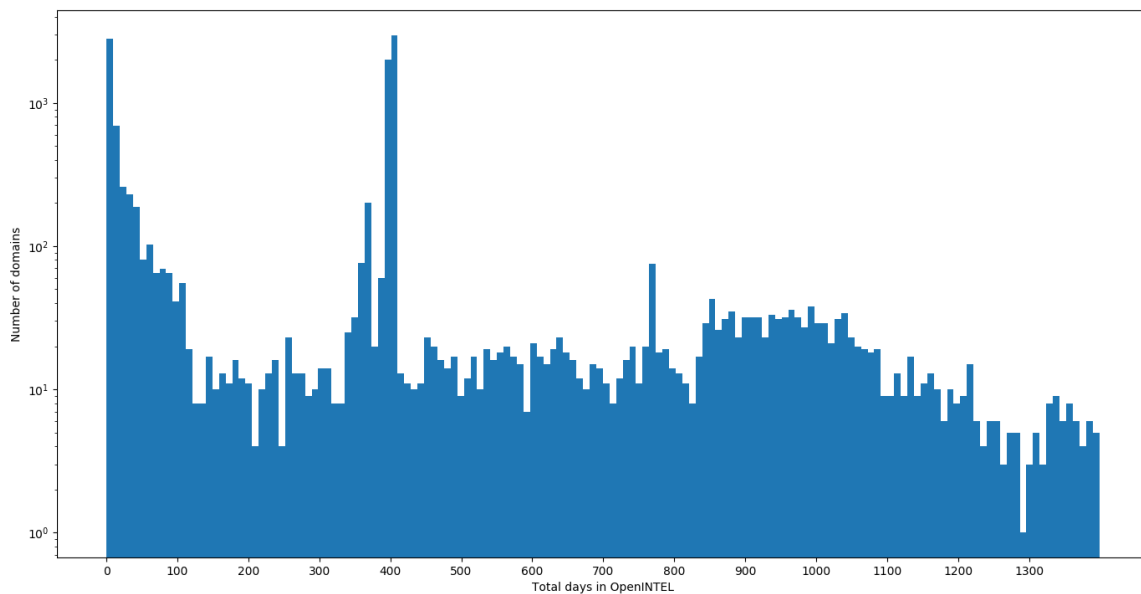
Rank	Trademark	Number of domains
1	Apple	8751
2	Paypal	1241
3	Microsoft	711
4	Netflix	592
5	Facebook	372
6	Amazon	323
7	Instagram	265
8	Google	213
9	Whatsapp	166
10	Wellsfargo	115

**Table 6.1:** Top 10 most targeted trademarks

The full list of trademarks and the corresponding frequencies is displayed in Section D.2. It can be observed that 'Apple' takes up the majority of malicious com-

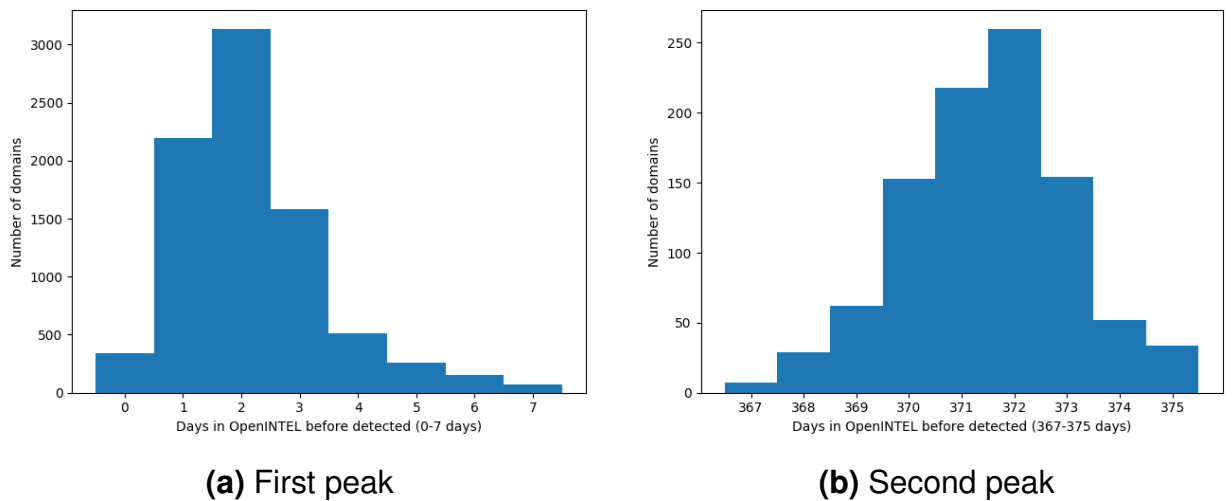
bosquat domains; as much as 8751 out of the 13693 domains. Following was 'Paypal' with 1241 hits. The top 5 is completed with 'Microsoft', 'Netflix' and 'Facebook', respectively with 711, 592 and 372 combosquat domains. Next, OpenINTEL was used to add more context to the domains. The next filtering process was applied, leaving 12115 domains for the analysis. These domains were then used to calculate four metrics: the total days of a domain in OpenINTEL, the days of a domain in OpenINTEL before it got blacklisted, the amount of domains that were still present in OpenINTEL after being removed from a blacklist and finally, the amount of days the domains was present in OpenINTEL after being blacklisted.

The first graph that could be created was the total days of a domain in OpenINTEL, as seen in Figure 6.1.

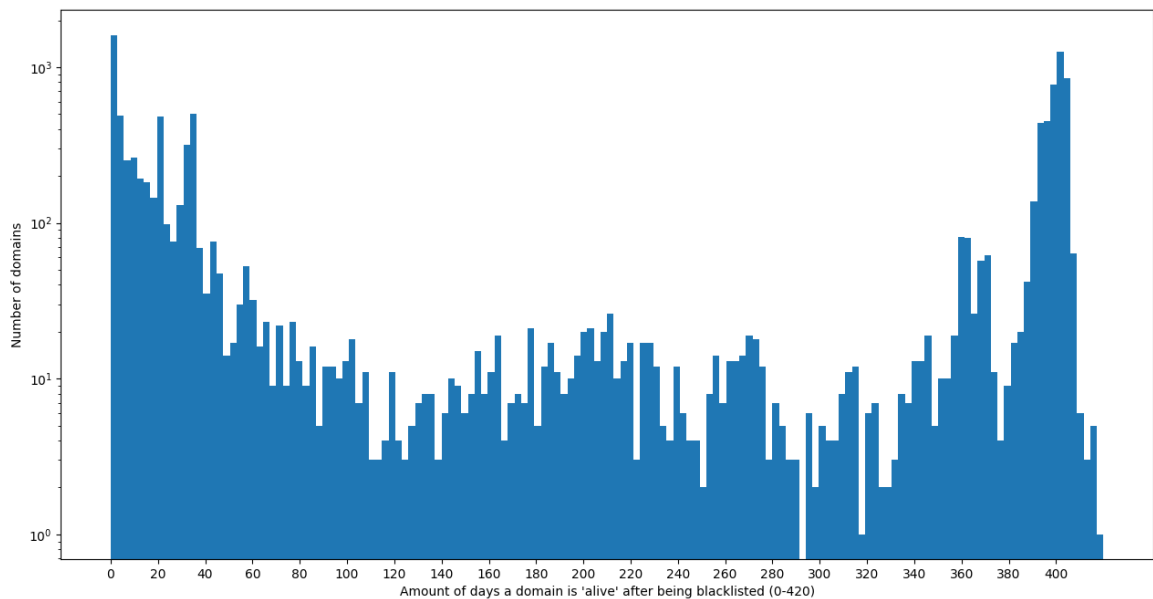


**Figure 6.1:** The total number of days a combosquat domain is present in OpenINTEL, with the  $y$ -axis in logarithmic scale.

When calculating and plotting the total days before a domain is detected two large peaks could be observed; the first peak lies at 2 days and the second peak lies at 372 days. These peaks are shown in Figure 6.2. The graph that shows all data can be found in Appendix D.1.



**Figure 6.2:** A graph showing the first and second peak from the graph showed in Appendix D.1



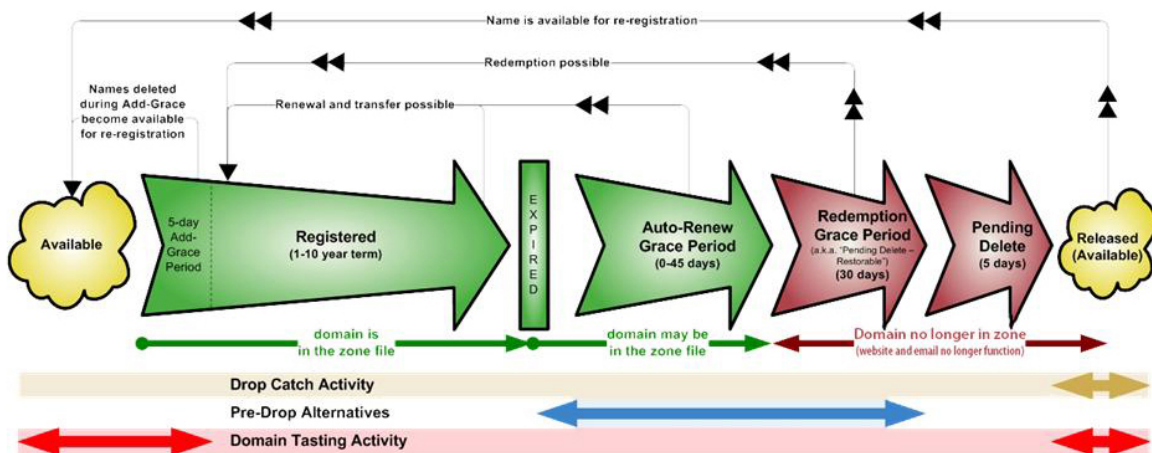
**Figure 6.3:** The number of days a domain is present in OpenINTEL after being detected, with the  $y$ -axis in logarithmic scale.

Next, it was calculated how many domains were still present in OpenINTEL after being removed from a blacklist. 86.7% of the domains were not present anymore, and 13.3% were still present in OpenINTEL after being removed from a blacklist. In this case, the reputation of the domain has improved in such a way that it got removed from the blacklist. Finally, Figure 6.3 shows the amount of days a domain is present in OpenINTEL after being detected.

Given the fact that most malicious combosquat domains are only active and detected after a few days, in order to define the killchain phases the domains needed to be measured more frequently. Since the historical DNS data was not suitable for this purpose, a additional small dataset consisting of active HTTP data was created using the Certstream<sup>1</sup> library which provide a real-time feed of newly signed SSL certificates. Since combosquat domains are used for phishing purposes, and phishing websites more frequently use SSL certificates to fake their legitimacy to users, it is expected that also SSL certificates for combosquat domains could be observed. A few combosquat domains were actively queried for their HTTP response during their lifecycle and the observations are summarized in Figure 6.4. This lifecycle sometimes only covered a few hours instead of a few days, which made it impossible for OpenINTEL to detect even the DNS changes. Note that this (extra) dataset was primarily used for exploring the possible stages the short-living domain could reside in. Although no conclusions can be based on this small dataset, future work could focus on the detection of malicious short-living combosquat domains primarily based on HTTP data.

### 6.1.1 Discussion

As a basis for the discussion, the domain lifecycle diagram provided by the ICANN<sup>2</sup> is provided in Figure 6.5.



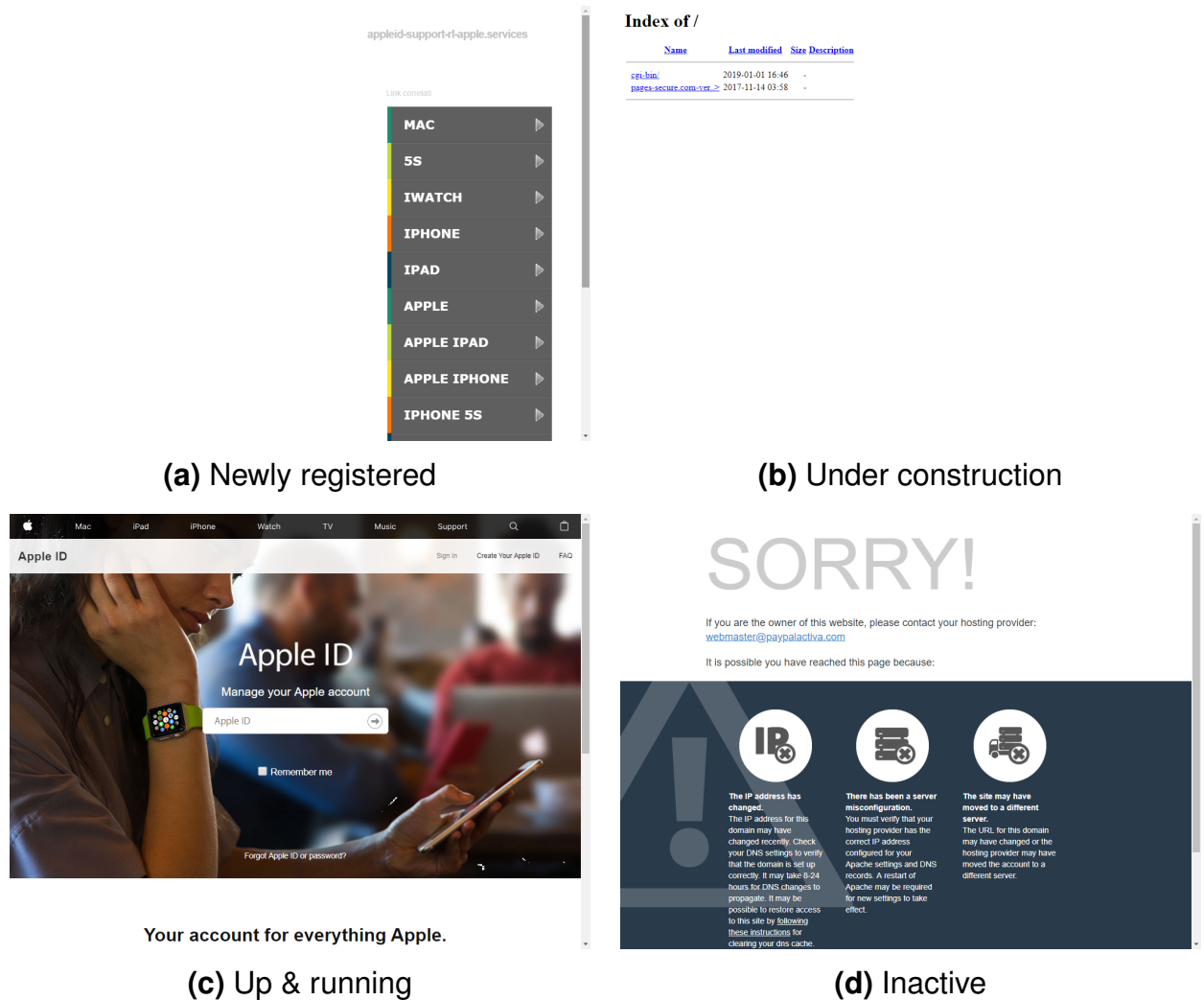
**Figure 6.5:** ICANN Domain lifecycle diagram

First, all the graphs that were created during the analysis are discussed. Figure 6.1 shows three interesting peaks. The first peak (around 0-7 days) are short-living domains; domains that are being registered and disappear within a few days.

<sup>1</sup><https://certstream.calidog.io>

<sup>2</sup>Internet Corporation for Assigned Names and Numbers





**Figure 6.4:** Screenshots of a short-living malicious combosquat domain

The second (small) peak is around 365, which correlated to a registration period of 1 year. However, the third peak lies around 405 days. This could be explained by summing the registration period (365 days) and the Auto-Renew Grace Period offered by a lot of registrars (40 days). In this Auto-Renew Grace Period, the domains is put 'on hold' to give the registrant some extra time to decide on continuing the registration or not, while the domain remains present in the zone file (see Figure 6.5). This means that no indication of an attackers' incentive to reuse or 'clean' a malicious domain can be seen; it seems that they just register combosquat domains in bulk and let them expire.

Figure 6.2 shows the two peaks in the total amount of days before domains are detected. It is interesting to see that the majority of detected domains are detected within ten days. Contrary to what Kintis et al. reported, the findings suggest that these types of domains are quickly blacklisted, often within necten days instead of the long-living domains reported by Kintis et. al. Figure 6.2b shows that the peak lies at 2 days. The second peak lies around 372. This peak has a similar shape as the peak around 2, however it has moved 370 days up front. This could be explained as a one-year registration period of 365 days + five additional days. It is assumed that these five days correspond to the Add-Grace Period (see Figure 6.5) of the registrar; within five days after registering a domain, the registration may be reversed by the registrar which in turn receives a full refund of the registry. After the registration is reversed, the domain becomes immediately available for re-registration. An explanation could be that the registration is reversed after five days, someone else registers that particular domain and it lives for another 365 days. Then, after 370 days, the domain expires and is drop-catched by a malicious registrant. This malicious registrant then uses the domain for malicious purposes, after which it gets blacklisted after an average 2 days.

Next are the amount of domains that are still present in OpenINTEL after being removed from a blacklist. A majority of 86.7% of the domains is not present in OpenINTEL after it has been removed from a blacklist, again indicating that there is no (economic) incentive for attackers to reuse the abused domains.

Lastly, Figure 6.3 shows the amount of days domains are alive after being blacklisted. This graph can be explained by taking Figure 6.1 and extracting the 2 days taken from Figure 6.2b in Figure 6.2. Therefore, this graph confirms the correlation between the different graphs.

Figure 6.4 shows the four phases that were frequently observed when actively obtaining HTTP information about short-living combosquat domains. Figure Figure 6.4a shows the domain being registered and parked; domains are not actively misused in this stage. Figure 6.4b shows the first signs of activity; a webserver is set-up and a folder containing malicious content is uploaded. Usually after this stage, it does

not take long before Figure 6.4c can be observed; the combosquat domains is up & running and its intentions are malicious. Finally, after some time the domain gets blacklisted / blocked / banned and the domain resolves to an error page, as shown in Figure 6.4d. Note that in this example, the domains is used for a phishing page, while several other types of abuse can also be present. Since users need to be directed to the phishing page, it is publicly promoted via email and/or the World Wide Web which possibly results in the fast detection.

Considering all of the above, OpenINTEL is not able to detect the really short-living combosquat domains (with a lifecycle of a few hours). Active HTTP would be needed for this purpose, which is unfortunately unavailable in this research. However, when examining Figure 6.2, some domains manage to remain undetected from a few days to a few hundred days. OpenINTEL is suitable to detect changes in DNS records related to a domain turning malicious (for example a change in A, AAAA or MX records).

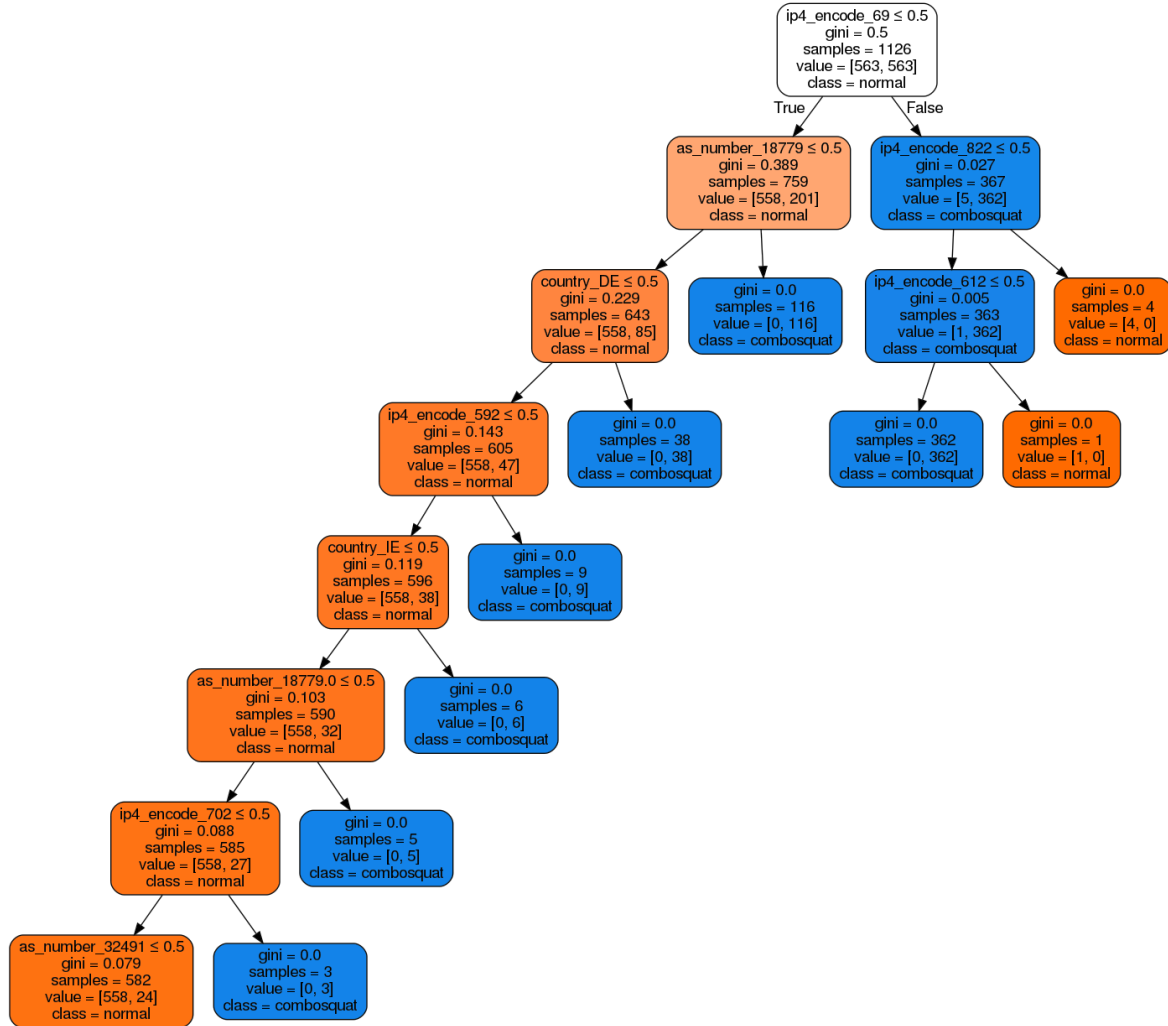
The killchain is used to define the different stages a combosquat domain can reside in. Based on the discussed figures and lifecycle phases, Table 6.2 was created. This table shows the first killchain stages in relation to combosquat domains, and whether it is possible to detect the transitions between killchain phases using either OpenINTEL or active HTTP measurements.

Killchain phase	Combosquat appliance	Detection with	
		OpenINTEL	HTTP
<i>Reconnaissance</i>	Attacker checks free domains	False	False
<i>Weaponization</i>	Attacker registers a combosquat domain	True	False
<i>Delivery</i>	Attacker configures malicious webserver	True*	True
<i>Exploitation</i>	Attacker directs users to domain	False	False
<i>Installation &amp; later</i>	User interacting with malicious page	False	True

**Table 6.2:** Killchain phases in combosquat perspective

The asterisk means that this is True when a change in DNS records is present. Usually, when registering a domain the registrar sets the DNS records to the registrar's defaults. Therefore, when an attacker for example changes the A or AAAA records to point to the malicious webserver, this can be observed in OpenINTEL. However, if the DNS records point to the malicious webserver from the start, OpenINTEL is not able to detect this change. Furthermore, if this change is made within one day after registration, OpenINTEL is not able to detect this because it only measures once per day.

## 6.2 Feature selection



**Figure 6.6:** Decision tree with the selected features

The ground truth has been created in the manner described in Section 4.4.  $M$  was defined as the day the domain first appeared on a blacklist,  $C$  as the day the change to the malicious DNS records was observed and  $B$  was the day before that change. The domains and the corresponding features from OpenINTEL were stored. For every domain, all dataframes were concatenated in chronological order. A hash was calculated over all dataframes and thus, over all features. First, the hash of day  $M - 1$  checked against the hash of  $M$ . If the hashed would differ,  $M - 1$  would be labeled as  $B$ . If the hashed matched, the iteration would continue until a day  $M - n$  would arise where the hash was not matching. The maximum value of  $n$  was set to 100 to increase the performance of the process; this means that a change in the DNS records had to be present in the 100 days before domain got blacklisted, which covers the first peak displayed in Figure 6.2 (a).

This resulted in a total of 5282 rows containing 'benign' and 'malicious' rows. This meant that 2641 domains were present that had the same features on  $C$  and  $M$ , and where the features of both  $B$  and  $M$  were available. This means that out of the 12115 domains, in 9474 cases no change was observed in the last 100 days, the features of  $C$  and  $M$  did not match, or the features of day  $B$  were not available.

The 5282 rows were then one-hot encoded, resulting in a ground truth dataframe with 5282 rows and 1533 columns. A *DecisionTreeClassifier* was then trained on the ground truth in the same way as described in Section 5.3. On average, this resulted in a decrease of columns from 1533 to 356. Finally, in Figure 6.6 the actual *DecisionTreeClassifier* that was used for the feature selection is shown.

### 6.2.1 Discussion

The top features outlined in Figure 6.6 are almost all based on the IP-addresses and AS numbers that correspond to the domains. The most significant features is *ipv\_encode\_69*, and the second-most significant features are *as\_number\_18779* and *ipv4\_encode\_822*. This means that no unique signature is observed for combosquat domains turning malicious except for a change in the IP address or AS range. This means that attackers change their DNS entries in order to point the domains to certain malicious IP addresses and AS ranges.



# Detection model design & validation

In this chapter, the results regarding the main research question are described. In Section 7.1 the model is designed and validated. Section 7.2 describes how the model was placed in the intended context and how the implementation was evaluated. Section 7.2 discusses the final outcomes of the combosquat detection model.

## 7.1 Treatment design and validation

The ground truth that was designed while answering CTD3 can again be used. This ground truth was used to extract the features as described in Figure 4.8.

A total of 10 classifiers were selected based on their usage in related work. For each of the 10 classifier, 10-fold cross validation was performed on the training i.e ground truth data. After training and validating the classifiers, the average scores of the 10 classifiers can be found in Table 7.1

Since the *RandomForestClassifier* has the highest *precision* and the lowest *FP rate*, this classifier is picked to be used in the real-world implementation.

## 7.2 Treatment implementation and evaluation

The trained *RandomForestClassifier* was used to detect malicious combosquat domains. The approach as described in Figure 4.3 was applied on the total set of combosquat domains on 30-07-2017. This date was chosen because it is in the middle of the total timespan of the dataset.

All combosquat domains on the date were extracted from OpenINTEL, resulting in a list of 173189 combosquat domains. For every domain in this list, the features were retrieved, one-hot encoding was applied and the columns of this test set were shaped to the training columns that were used by the trained *RandomForestClassifier*.

<b>Classifier</b>	<b>FP rate %</b>	<b>Accuracy %</b>	<b>Precision %</b>
DecisionTreeClassifier	24	70	72
RandomForestClassifier	22	70	73
AdaBoostClassifier	25	70	71
KNeighborsClassifier	25	68	70
GaussianNB	83	52	51
BernoulliNB	37	65	64
MLPClassifier	47	61	62
SGDClassifier	60	46	46
GradientBoostingClassifier	23	67	70
ExtraTreesClassifier	24	70	73

**Table 7.1:** Classifiers and their scores

After the predictions were made, the actual labels were also calculated for every domain and consecutively a confusion matrix was constructed.

		<b>Prediction outcome</b>	
		<b>m</b>	<b>b</b>
<b>Actual value</b>	<b>m'</b>	TP 2227	FN 5719
	<b>b'</b>	FP 30422	TN 134810
<b>Total</b>		32649	140529

**Table 7.2:** Confusion matrix for the real-world validation

Based on the confusion matrix, the *precision*, *FP rate* and the *accuracy* can be calculated as explained in Subsection 4.2.3:



$$\begin{aligned} Precision &= \frac{2227}{2227 + 30422} \\ &= 6.8\% \end{aligned}$$

$$\begin{aligned} FP\ rate &= \frac{30422}{30422 + 134810} \\ &= 18.41\% \end{aligned}$$

$$\begin{aligned} Accuracy &= \frac{2227 + 134810}{2227 + 5719 + 30422 + 134810} \\ &= 79.13\% \end{aligned}$$

The scores are insufficient to be used for efficient detection of combosquat domains turning malicious. The *Precision* score is too low to be of practical use, as well as the high *FP rate*.



# Conclusion

The chapter provides answers for the research questions that were formulated in Section 2.4. In Section 8.1 the conclusions from the subquestion are listed, and the main research question is answered. Finally, in Section 8.2 the recommendations for future work are discussed.

## 8.1 Conclusion

Based on the literature study there was a hypothesis that, with enough data in the form of active DNS measurements, a generic combosquat detection model could be designed that was able to warn customers in an early stage. To test the hypothesis, a research question was defined, that was further divided into multiple subquestions: *How to develop an active combosquatting detection model that meets the requirements set in Section 2.3, so that businesses can be warned in an earlier stage within a threat intelligence platform?*

The first sub question was about investigating whether a generic model was possible for the detection of combosquat domains using active DNS measurements. The first observation was that the key feature is `domainname_characters`, thus the total length of the domainname. Since the addition of words to a trademarks often results in a lengthy domain name, the classifier was mainly trained on this one feature. While this holds for combosquat domains, this also holds for other general purpose domain names, such as *smarthomegadgetguru.com*, not an unique domain name in itself. The second observation was that lexical features (based on the domain name itself) were more important than features selected from DNS data. Combosquatting domains do not significantly differ from benign domains and/or other malicious domains in a distinctive manner regarding the selected DNS features. The third observation was that based on the data & selected features, no clear distinction could

be made between a regular word and a trademark. Apple illustrates this problem clearly; it is both a well-known trademark, but also a regular word frequently used in domain names. These trademarks make it difficult for an automated approach to distinguish domains where Apple is used as a trademark and where not.

Through the observations it can be concluded that it is extremely difficult to construct a generic model for detecting combosquat domains without a predefined list of trademarks.

To answer the second subquestion, a temporal analysis was performed on the combined OpenINTEL & blacklist data. This resulted into new insights regarding combosquat domain usage by malicious users and the (economic) incentives behind it. These findings were then translated into actions that malicious users could perform in the separate killchain stages.

A combosquat domain can reside in all of the killchain phases, however the first five phases (Reconnaissance, Weaponization, Delivery, Exploitation and Installation) are most useful for detection purposes. Detection based on DNS data is under certain circumstances possible between the Weaponization & Delivery phase. Short-living domains are more difficult to detect with OpenINTEL because of the measurement frequency of one time per day; domains that are registered, turn malicious and are abandoned within a day remain undetected in this way. Detection of combosquat domains turning malicious based on OpenINTEL data should therefore be focused on domains living longer than 1 day. If one would want to detect also the short-living domains, other data sources such as active HTTP measurements should be added. Combosquat domains in the Weaponization phase are considered benign, while domains in the Delivery phase are considered malicious.

The last subquestion is answered by looking at the features selected in Section 6.2. These features were related to the IP-addresses and AS numbers of the domain. So based on changes in the IP-addresses and AS numbers of a domain, a benign combosquat domain turning malicious can be identified.

The results show that the detection of combosquat domains turning malicious based on active DNS measurements is not sufficient, when considering acceptable scores for the *False Positive rate*(5%) and the *precision*(90%). The *False Positive rate* of 18.41% is too high to be of practical use, similar to the low *precision* rate of 6.8%. Therefore, the first observation is that based on only active DNS measurements and

given the features that were used, it is not possible to detect when a combosquat domain transforms from an 'inactive' state to an 'active' state. The second observation is that when looking at the features that define a legitimate domain turning malicious, is that no unique indicator of change can be observed. As can be seen in the decision tree in Figure 6.6, the most relevant features consist of IP-addresses and AS numbers. This method of detection has been widely researched and is also actively being used in practice; IP-addresses and AS numbers are getting rated and blacklisted regularly. Thus, the classifier that uses IP-addresses and AS numbers is not considered 'new' and it matches the current detection methods.

## 8.2 Future work & recommendations

One of the main conclusions is that the detection of combosquat domains turning malicious is not sufficient when it is only based on active DNS measurements. This obviously does not mean that this is not possible at all. By using & combining other data sources, for example HTTP data, detection may be possible. Although a lot of research is also done in this field, it has not been applied to combosquat domains specifically.

It should also be noted that registrars at this point are not succeeding in declining combosquat registrations. On the one hand, it is because this is very hard to do, as this research shows. Another aspect is that it is against the nature of a registrant selling domain names is their core business and generates revenue. Limiting the number of registrations implies less revenue. Although currently there is legislation in place to prevent combosquatting abuse, it is obviously not enforced properly. An option might be to make the domain registration procedure more restricted by law; registrars may ask for more information about the registrants in order to verify the actual identity, such as passport numbers, legal entity numbers and more. Since every registrar is only responsible for a selection of TLD's, this may even be different per country.

To conclude, the detection results from this thesis can be used for future work, as an effective way of detecting combosquat domains is strongly desired.



# Bibliography

- [1] H. Al-Mohannadi, Q. Mirza, A. Namanya, I. Awan, A. Cullen, and J. Disso, "Cyber-attack modeling analysis techniques: An overview," in *2016 IEEE 4th International Conference on Future Internet of Things and Cloud Workshops (FiCloudW)*, Aug 2016, pp. 69–76.
- [2] P. Dreyer, T. Jones, K. Klima, J. Oberholtzer, A. Strong, J. W. Welburn, and Z. Winkelman, "Estimating the global cost of cyber risk," 2018.
- [3] "Fox-it | for a more secure society," <https://www.fox-it.com>, accessed: 2018-09-07.
- [4] Y. Ayrou, A. Raji, and M. Nassar, "Modelling cyber-attacks: a survey study," *Network Security*, vol. 2018, no. 3, pp. 13 – 19, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1353485818300254>
- [5] E. M. Hutchins, M. J. Cloppert, and R. M. Amin, "Intelligence-driven computer network defense informed by analysis of adversary campaigns and intrusion kill chains," *Leading Issues in Information Warfare & Security Research*, vol. 1, no. 1, p. 80, 2011.
- [6] W. Tounsi and H. Rais, "A survey on technical threat intelligence in the age of sophisticated cyber attacks," *Computers & Security*, vol. 72, pp. 212 – 233, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167404817301839>
- [7] P. Kintis, N. Miramirkhani, C. Lever, Y. Chen, R. Romero-Gómez, N. Pitropakis, N. Nikiforakis, and M. Antonakakis, "Hiding in plain sight: A longitudinal study of combosquatting abuse," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '17. New York, NY, USA: ACM, 2017, pp. 569–586. [Online]. Available: <http://doi.acm.org/10.1145/3133956.3134002>
- [8] "Combosquatting: The business of cybersquatting," Fairwind Partners LLC, Tech. Rep., 2008.

- [9] O. van der Toorn, R. van Rijswijk-Deij, B. Geesink, and A. Sperotto, "Melting the snow: Using active dns measurements to detect snowshoe spam domains," in *NOMS 2018 - 2018 IEEE/IFIP Network Operations and Management Symposium*, April 2018, pp. 1–9.
- [10] R. Wieringa, *Design science methodology for information systems and software engineering*. Springer, 2014, 10.1007/978-3-662-43839-8.
- [11] P. V. Mockapetris, "Domain names - concepts and facilities." *RFC*, vol. 1034, pp. 1–55, November 1987. [Online]. Available: <http://dblp.uni-trier.de/db/journals/rfc/rfc1000-1099.html>
- [12] E. Kidmose, E. Lansing, S. Brandbyge, and J. M. Pedersen, "Detection of malicious and abusive domain names," in *2018 1st International Conference on Data Intelligence and Security (ICDIS)*, April 2018, pp. 49–56.
- [13] "Sidn : Jouw wereld. ons domein." <https://www.sidn.nl/>, accessed: 2018-09-21.
- [14] Y. Zhauniarovich, I. Khalil, T. Yu, and M. Dacier, "A survey on malicious domains detection through DNS data analysis," *CoRR*, vol. abs/1805.08426, 2018. [Online]. Available: <http://arxiv.org/abs/1805.08426>
- [15] L. Bilge, S. Sen, D. Balzarotti, E. Kirda, and C. Kruegel, "Exposure: A passive dns analysis service to detect and report malicious domains," *ACM Trans. Inf. Syst. Secur.*, vol. 16, no. 4, pp. 14:1–14:28, Apr. 2014. [Online]. Available: <http://doi.acm.org/10.1145/2584679>
- [16] S. Schiavoni, F. Maggi, L. Cavallaro, and S. Zanero, "Phoenix: Dga-based botnet tracking and intelligence," in *Detection of Intrusions and Malware, and Vulnerability Assessment*, S. Dietrich, Ed. Cham: Springer International Publishing, 2014, pp. 192–211.
- [17] R. Sharifnya and M. Abadi, "Dfbotkiller: Domain-flux botnet detection based on the history of group activities and failures in dns traffic," *Digital Investigation*, vol. 12, pp. 15 – 26, 2015. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1742287614001182>
- [18] S. Hao, A. Kantchelian, B. Miller, V. Paxson, and N. Feamster, "Predator: Proactive recognition and elimination of domain abuse at time-of-registration," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '16. New York, NY, USA: ACM, 2016, pp. 1568–1579. [Online]. Available: <http://doi.acm.org/10.1145/2976749.2978317>



- [19] T. Vissers, J. Spooren, P. Agten, D. Jumpertz, P. Janssen, M. Van Wesemael, F. Piessens, W. Joosen, and L. Desmet, "Exploring the ecosystem of malicious domain registrations in the .eu tld," in *Research in Attacks, Intrusions, and Defenses*, M. Dacier, M. Bailey, M. Polychronakis, and M. Antonakakis, Eds. Cham: Springer International Publishing, 2017, pp. 472–493.
- [20] J. Szurdi and N. Christin, "Domain registration policy strategies and the fight against online crime," *WEIS*, June 2018.
- [21] Y.-M. Wang, D. Beck, J. Wang, C. Verbowski, and B. Daniels, "Strider typo-patrol: Discovery and analysis of systematic typo-squatting," in *Proceedings of the 2Nd Conference on Steps to Reducing Unwanted Traffic on the Internet - Volume 2*, ser. SRUTI'06. Berkeley, CA, USA: USENIX Association, 2006, pp. 5–5. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1251296.1251301>
- [22] P. Piredda, D. Ariu, B. Biggio, I. Corona, L. Piras, G. Giacinto, and F. Roli, "Deepsquatting: Learning-based typosquatting detection at deeper domain levels," in *AI\*IA*, 2017.
- [23] N. Nikiforakis, S. Van Acker, W. Meert, L. Desmet, F. Piessens, and W. Joosen, "Bitsquatting: Exploiting bit-flips for fun, or profit?" in *Proceedings of the 22Nd International Conference on World Wide Web*, ser. WWW '13. New York, NY, USA: ACM, 2013, pp. 989–998. [Online]. Available: <http://doi.acm.org/10.1145/2488388.2488474>
- [24] N. Nikiforakis, M. Balduzzi, L. Desmet, F. Piessens, and W. Joosen, "Sound-squatting: Uncovering the use of homophones in domain squatting," in *Information Security*, S. S. M. Chow, J. Camenisch, L. C. K. Hui, and S. M. Yiu, Eds. Cham: Springer International Publishing, 2014, pp. 291–308.
- [25] T. Holgers, D. E. Watson, and S. D. Gribble, "Cutting through the confusion: A measurement study of homograph attacks," in *Proceedings of the Annual Conference on USENIX '06 Annual Technical Conference*, ser. ATEC '06. Berkeley, CA, USA: USENIX Association, 2006, pp. 24–24. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1267359.1267383>
- [26] J.-W. Bullee, L. Montoya, M. Junger, and P. Hartel, "Spear phishing in organisations explained," *Information and Computer Security*, vol. 25, no. 5, pp. 593–613, 7 2017.
- [27] P. Lv, J. Ya, T. Liu, J. Shi, B. Fang, and Z. Gu, "You have more abbreviations than you know: A study of abbrevsquatting abuse," in *Computational Science – ICCS 2018*, Y. Shi, H. Fu, Y. Tian, V. V. Krzhizhanovskaya, M. H. Lees, J. Dongarra,

- and P. M. A. Sloot, Eds. Cham: Springer International Publishing, 2018, pp. 221–233.
- [28] “Alexa top 1m,” <http://s3.amazonaws.com/alexa-static/top-1m.csv.zip>, accessed: 2018-09-17.
- [29] “Certificate transparency,” <https://transparencyreport.google.com/https/certificates>, accessed: 2018-09-11.
- [30] T. Segaran and J. Hammerbacher, *Beautiful data: the stories behind elegant data solutions*. ” O’Reilly Media, Inc.”, 2009.
- [31] “Pyenchant tutorial,” [https://faculty.math.illinois.edu/~gfrancis/illimath/windows/aszgard\\_mini/movpy-2.0.0-py2.4.4/manuals/PyEnchant/PyEnchant%20Tutorial.htm](https://faculty.math.illinois.edu/~gfrancis/illimath/windows/aszgard_mini/movpy-2.0.0-py2.4.4/manuals/PyEnchant/PyEnchant%20Tutorial.htm), accessed: 01-11-2018.
- [32] “Slang dictionary - internet & text slang,” <http://www.noslang.com/dictionary/>, accessed: 01-11-2018.
- [33] “List of swear words, bad words, & curse words - starting with a,” <http://www.noswearing.com/dictionary>, accessed: 01-11-2018.
- [34] “Sowpods scrabble word list,” <https://www.wordgamedictionary.com/sowpods/>, accessed: 01-11-2018.
- [35] D. Chiba, K. Tobe, T. Mori, and S. Goto, “Detecting malicious websites by learning ip address features,” in *2012 IEEE/IPSJ 12th International Symposium on Applications and the Internet*, July 2012, pp. 29–39.

## **Appendix A**

### **Generic model data**

#### **A.1 Manual trademark selection**

n=107

soundcloud dropbox nbc sports nordstrom premierleague twitch alipay mailchimp skype aliexpress paypal siemens dailymail googleplus indiatimes tomshardware norton theguardian netflix quora playstation yandex ieee pinterest airbnb huffingtonpost flickr salesforce bankofamerica pastebin instagram tripadvisor foxnews wikipedia github nvidia apple americanexpress youtube fox-it twitter stackoverflow dailymotion office365 facebook spotify usatoday snapchat microsoftonline hewlett steam-powered duckduckgo homedepot hdfcbank thestartmagazine slideshare walmart bloomberg epicgames samsung amazonaws autodesk packard nytimes godaddy alibaba mediafire expedia wordpress linkedin Breitbart amazon tumblr marktplaats fujitsu elsevier filehippo ladbible google wells Fargo reddit nokia mozilla Symantec McAfee Microsoft Shopify WhatsApp Avast Utwente Gamepedia Verizon CloudFront 4shared Adobe StackExchange GitLab League of Legends Lenovo WeTransfer Wiktionary

#### **A.2 Frequent combosquatting words**

Taken from the paper by Kintis et al. [7]

car universal square villa cheap marketing search porno account print office content vacation official listen wire hot shipping worldwide county services pilgrim net free videos san shop plus sex business fuck health group maps online delivery apps phone channel play princess watch kindle support post home com wireless mobile island theme freight inn news foundation posting president vote yeah archive service photography gift glass store photos club south sale express trump tube life jobs energy mortgage mike file sucks world elect center ground galaxy views live update user xxx campaign garden stores time cards deals page media zine university blog

hotel trust login family movies movie buy just card head best line video stop lay gay  
land love google real music chill themes beach cars estate followers porn school  
city paris download games prices credit hotels mail bank price property tex truth  
new phones canada chemical movie photo converter investment

### A.3 Public blacklists

This table shows the blacklists that were used when answering CTD1.

Source
<a href="http://www.abuse.ch/">http://www.abuse.ch/</a>
<a href="http://www.malwaredomains.com/">http://www.malwaredomains.com/</a>
<a href="http://www.malwaredomains.com/">http://www.malwaredomains.com/</a>
<a href="http://dns-bh.sagadc.org/">http://dns-bh.sagadc.org/</a>
<a href="https://isc.sans.edu/suspicious_domains">https://isc.sans.edu/suspicious_domains</a>
<a href="http://www.urlvir.com/exporthosts/">http://www.urlvir.com/exporthosts/</a>
<a href="http://www.nothink.org/blacklist/blacklist_malware_dns.txt">http://www.nothink.org/blacklist/blacklist_malware_dns.txt.</a>
<a href="http://www.joewein.net/dl/bl/dom-bl.txt">http://www.joewein.net/dl/bl/dom-bl.txt.</a>

## Appendix B

# Python code snippets

This appendix holds several code snippets, used in the thesis. They can be used to clarify the methodology and to see what was actually calculated. Note that the imports are not shown, as well as code that is not useful to show, e.g. code that loads and saves DataFrames, trademarks, words and more.

## B.1 OpenINTEL queries

---

```
SELECT
lower(query_name) AS domainname,
count(case response_type when 'A' then 1 else null end) AS number_of_A_records,
count(case response_type when 'AAAA' then 1 else null end) AS number_of_AAAA_records,
count(case response_type when 'NS' then 1 else null end) AS number_of_NS_records,
count(case response_type when 'MX' then 1 else null end) AS number_of_MX_records,
count(case response_type when 'SOA' then 1 else null end) AS number_of_SOA_records,
count(case response_type when 'CNAME' then 1 else null end) AS number_of_CNAME_records,
count(case response_type when 'DNSKEY' then 1 else null end) AS number_of_DNSKEY_records,
count(case response_type when 'TXT' then 1 else null end) AS number_of_TXT_records,
count(ip4_address) AS number_of_ipv4_addresses,
concat("[", group_concat(ip4_address), "]") AS ipv4_addresses,
count(case query_name when response_name then 1 else null end) AS response_name_matches,
max(country) as country,
max('as') as as_number,
avg(soa_refresh) AS soa_refresh,
avg(soa_retry) AS soa_retry,
avg(soa_minimum) AS soa_minimum
FROM openintel.com_warehouse_parquet
WHERE year = [current_year] AND month = [current_month] AND day = [current_day]
AND lower(query_name) NOT LIKE '123-nonexistant-dnsjedi-456.%'
AND lower(query_name) NOT LIKE 'www.%'
AND regexp_like(lower(query_name), [manual_trademarks_regex])
AND regexp_like(lower(query_name), [frequent_combosquat_words_regex])
```

---

**GROUP BY** `lower(query_name)`

---

## B.2 Validation of real-world domains

Every domain has a predicted value True or False, provided by the trained classifier. The code below calculates the `actual` value for the domain. In the end, every domain has a `predicted` and `actual` boolean, which are then used to construct the confusion matrix.

---

```
'''
Function for validating whether the predicted domains are actual combosquat domains or not.
Consists of a few simple checks.
'''

def is_combosquat(domainname):
    contains_trademark = alexaregex.search(domainname)

    # Check for constraint 1) and 2)
    if contains_trademark:
        trademark = contains_trademark.group(0)
        # Rule out typosquatting, so a levenshtein distance of 1. This (
        # partially) checks constraint 5)
        if distance.levenshtein(trademark, domainname) == 1:
            return False

    segmented_domainname = wordsegment.segment(domainname)

    # Segmented_domainname contains a list of segments here. Now,
    # check for constraint 3)
    is_standalone_word = False

    for segment in segmented_domainname:
        if alexaregex.fullmatch(segment):
            is_standalone_word = True
            break

    if is_standalone_word:
        # Now we only have to check whether the two IP's are not
        # in the same range and the AS numbers do not match.
        # This checks constraint 4)
        return not ip_and_as_match(trademark, domainname)

    # If it's not a standalone word, dismiss it. E.g.
    # 'applejuice.com' should not be included.
    else:
```

```

        return False

    # If no trademark is present in the domainname, dismiss it already
    else:
        return False

'''
For a given trademark_domain (e.g. amazon) and new_domainname (amazon-secure-login),
this function checks whether the AS numbers match and the IP of the new_domainname
is in the /16 range of the IPv4 of the trademark_domain
'''
def ip_and_as_match(trademark_domain, new_domainname):
    originaldomain_as = domains_with_features.loc[domains_with_features['domainname']
        == trademark_domain + '.com.']['as_number'].values[0]
    originaldomain_ips = domains_with_features.loc[domains_with_features['domainname']
        == trademark_domain + '.com.']['ipv4_addresses'].values
    domainname_as = new_domains.loc[new_domains['domainname']
        == new_domainname + '.com.']['as_number'].values[0]
    domainname_ips = new_domains.loc[new_domains['domainname']
        == new_domainname + '.com.']['ipv4_addresses'].values

    ip_in_original_range = False
    as_numbers_match = originaldomain_as == domainname_as

    for originaldomain_ip in originaldomain_ips:
        network = ip_network(originaldomain_ip + "/16", strict=False)
        for domainname_ip in domainname_ips:
            if ip_address(domainname_ip) in network:
                # Here, the IP of the domain is in the original range
                ip_in_original_range = True

    return ip_in_original_range and as_numbers_match

```

---

## B.3 Running the prototype, training & test phase

No backup for today available, creating new one..

Today is: 2019-01-27 12:44:15.613630:

So, the ground truth is based on 2019-01-25 12:44:15.613631

Performing Kerberos authentication for OpenINTEL ...

found keytab: /home/jjansen/oi\_jjansen.keytab

OK

Querying OpenINTEL for combosquatting domains

Got query response, now writing to csv...

Written new info to combosquatdomains\_2512019.csv

Data transferred from voordeur!

```

Importing abuse.ch [Elapsed Time: 0:00:01] |#####| (Time: 0:00:01)
Importing hosts-file.net [Elapsed Time: 0:00:17] |#####| (Time: 0:00:17)
Importing nothink.org [Elapsed Time: 0:00:00] |#####| (Time: 0:00:00)
Importing malwaredomainlist delisted [Elapsed Time: 0:00:00] || (Time: 0:00:00)
Importing malwaredomainlist blacklist [Elapsed Time: 0:00:00] |#| (Time: 0:00:00)
Importing malwaredomains immortal [Elapsed Time: 0:00:00] |#| (Time: 0:00:00)
Importing malwaredomains default [Elapsed Time: 0:00:00] |##| (Time: 0:00:00)
Importing joewein [Elapsed Time: 0:00:00] |#####| (Time: 0:00:00)
Importing isc_sans_edu 1/3 [Elapsed Time: 0:00:00] |#####| (Time: 0:00:00)
Importing isc_sans_edu 2/3 [Elapsed Time: 0:00:00] |#####| (Time: 0:00:00)
Importing isc_sans_edu 3/3 [Elapsed Time: 0:00:00] |#####| (Time: 0:00:00)
Importing urlvir.com [Elapsed Time: 0:00:00] |#####| (Time: 0:00:00)

```

Out of the 285327 domains, a total of 5274 is present on passive blacklists!

Adding 2637 Alexa & Random malicious domains

Done

Starting adding lexical features..

Done

Start adding contextual features..

Done

Applying one hot encoding..

Done

Datagram shape: (10548, 6106)

is\_combosquatting

False 5274

True 5274

Name: is\_combosquatting, dtype: int64

Shape before features selection:

(10548, 6106)

Shape after feature selection:

(10548, 676)

Selected features:

Index(['number\_of\_a\_records', 'number\_of\_aaaa\_records', 'number\_of\_ns\_records',



```
'number_of_mx_records', 'number_of_txt_records',
'number_of_ipv4_addresses', 'response_name_matches', 'soa_refresh',
'soa_retry', 'soa_minimum',
...
'ip4_encode_995', 'ip4_encode_998', 'ip4_encode_999', 'ip4_encode_1005',
'ip4_encode_1006', 'ip4_encode_1008', 'ip4_encode_1009',
'ip4_encode_1011', 'ip4_encode_1012', 'ip4_encode_1017'],
dtype='object', length=676)

#####
DecisionTreeClassifier
#####

Average FP rate: 22 %
Average precision: 77 %
Average accuracy: 78 %
Raw FP: 118
Raw TP: 421
Raw TN: 408
Raw FN: 106

#####
RandomForestClassifier
#####

Average FP rate: 17 %
Average precision: 81 %
Average accuracy: 80 %
Raw FP: 92
Raw TP: 424
Raw TN: 434
Raw FN: 102

#####
AdaBoostClassifier
#####
```

```
Average FP rate: 18 %
Average precision: 80 %
Average accuracy: 80 %
Raw FP: 101
Raw TP: 434
Raw TN: 425
Raw FN: 92
```

```
#####
KNeighborsClassifier
#####
```

```
Average FP rate: 26 %
Average precision: 74 %
Average accuracy: 76 %
Raw FP: 144
Raw TP: 432
Raw TN: 382
Raw FN: 95
```

```
#####
GaussianNB
#####
```

```
Average FP rate: 6 %
Average precision: 54 %
Average accuracy: 51 %
Raw FP: 34
Raw TP: 54
Raw TN: 492
Raw FN: 472
```

```
#####
BernoulliNB
#####
```

Average FP rate: 26 %  
Average precision: 73 %  
Average accuracy: 73 %  
Raw FP: 141  
Raw TP: 392  
Raw TN: 385  
Raw FN: 135

#####  
MLPClassifier  
#####

Average FP rate: 31 %  
Average precision: 69 %  
Average accuracy: 67 %  
Raw FP: 166  
Raw TP: 357  
Raw TN: 360  
Raw FN: 170

#####  
SGDClassifier  
#####

Average FP rate: 76 %  
Average precision: 48 %  
Average accuracy: 48 %  
Raw FP: 406  
Raw TP: 393  
Raw TN: 120  
Raw FN: 133

#####  
GradientBoostingClassifier  
#####

```
Average FP rate: 23 %
Average precision: 77 %
Average accuracy: 81 %
Raw FP: 126
Raw TP: 459
Raw TN: 400
Raw FN: 68
```

```
#####
ExtraTreesClassifier
#####
```

```
Average FP rate: 20 %
Average precision: 79 %
Average accuracy: 81 %
Raw FP: 109
Raw TP: 446
Raw TN: 417
Raw FN: 81
```

Added all classifiers with a FP rate lower than 10 percent and a precision higher than

```
Starting the Combosquat Detection Model testing phase!
Loading new domains backup..
Done
Only processing first 1000 entries
Enriching new domains with lexical features
Start Alexa filtering
Done
Starting adding lexical features..
Done
Start adding contextual features
Done with the WHOIS requests, now heading to the CT Log features
Done with the contextual features
Applying one hot encoding..
Shape before one-hot encoding: (10000, 31)
Done
Removing columns not present in training columns..
```

Adding dummy test columns for missing training columns..

Done

Total new domains: 10000

75 combosquat domains found using GaussianNB

Done

Starting the Combosquat Detection Model validation phase!

Validation complete!

TP: 0

FP: 75

TN: 9910

FN: 15



## Appendix C

# Scraped blacklists

This table shows the blacklists that were included in the scraped blacklists, starting from **2016-07-08** and ending at **2019-01-11**.

Source
<a href="http://www.malwaredomainlist.com/hostslist/hosts.txt">http://www.malwaredomainlist.com/hostslist/hosts.txt</a>
<a href="http://www.malwaredomainlist.com/hostslist/delisted.txt">http://www.malwaredomainlist.com/hostslist/delisted.txt</a>
<a href="http://mirror1.malwaredomains.com/files/justdomains">http://mirror1.malwaredomains.com/files/justdomains</a>
<a href="http://www.joewein.net/dl/bl/dom-bl.txt">http://www.joewein.net/dl/bl/dom-bl.txt</a>
<a href="http://malc0de.com/bl/ZONES">http://malc0de.com/bl/ZONES</a>
<a href="https://zeustracker.abuse.ch/blocklist.php?download=domainblocklist">https://zeustracker.abuse.ch/blocklist.php?download=domainblocklist</a>
<a href="https://ransomwaretracker.abuse.ch/downloads/RW_DOMBL.txt">https://ransomwaretracker.abuse.ch/downloads/RW_DOMBL.txt</a>
<a href="https://hosts-file.net/hphosts-partial.txt">https://hosts-file.net/hphosts-partial.txt</a>
<a href="https://palevotracker.abuse.ch/blocklists.php?download=domainblocklist">https://palevotracker.abuse.ch/blocklists.php?download=domainblocklist</a> (until 06-12-2016)
<a href="https://feodoctracker.abuse.ch/blocklist/?download=domainblocklist">https://feodoctracker.abuse.ch/blocklist/?download=domainblocklist</a>
<a href="http://www.networksec.org/grabbho/block.txt">http://www.networksec.org/grabbho/block.txt</a>
<a href="https://openphish.com/feed.txt">https://openphish.com/feed.txt</a>
<a href="https://www.threatcrowd.org/feeds/domains.txt">https://www.threatcrowd.org/feeds/domains.txt</a>
<a href="https://urlhaus.abuse.ch/downloads/text/">https://urlhaus.abuse.ch/downloads/text/</a>
<a href="http://osint.bambenekconsulting.com/feeds/c2-dommasterlist.txt">http://osint.bambenekconsulting.com/feeds/c2-dommasterlist.txt</a>
<a href="http://vxvault.net/URL_List.php">http://vxvault.net/URL_List.php</a>

On a daily basis, the blacklists were fetched and stored in the following format:

domainname	source	date
------------	--------	------

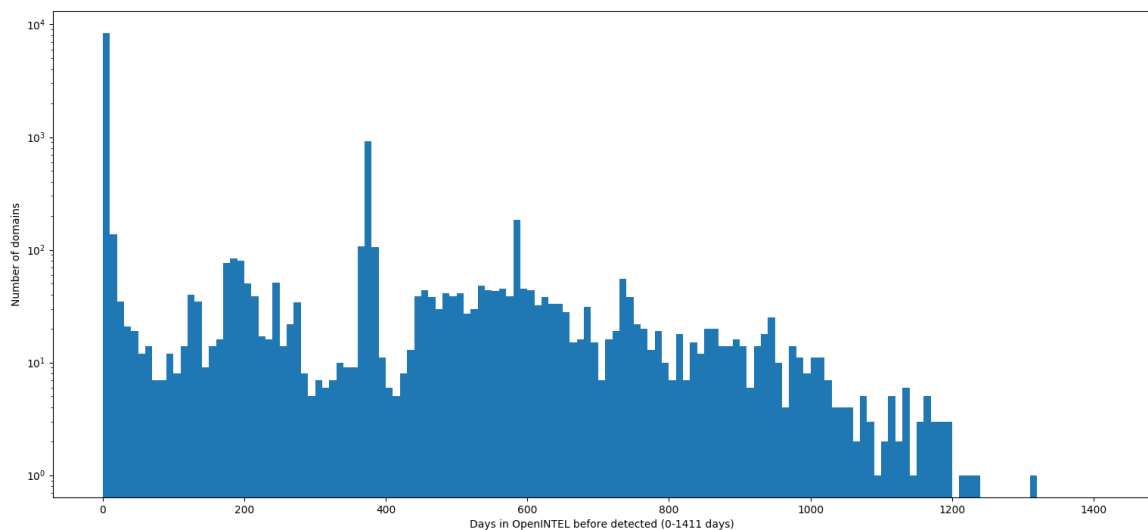
This resulted in a total of 870 files (8.3GB). Since the total days between the first and last date is 907 days, 37 days were missing due to switching to an other machine and/or temporary measurement failures. Out of the 870 files, 3 files turned out corrupt, leaving a total of **867** files to work with.





# Trademark distribution

## D.1 Total days before blacklisted



**Figure D.1:** The total number of days a combosquat domain is present in OpenINTEL before being listed on a blacklist, with the  $y$ -axis in logarithmic scale.

## D.2 Trademark frequency

