

# Detection of malicious domains through lexical analysis

International Conference On Cyber Security And Protection Of Digital Services

June 11, 2018

Egon Kidmose, Matija Stevanovic and Jens Myrup Pedersen  
egk@es.aau.dk, mst@es.aau.dk, jens@es.aau.dk

Department of Electronic Systems  
Aalborg University  
Denmark



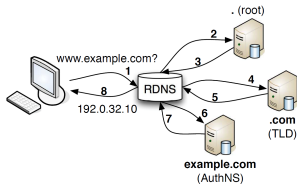
**AALBORG UNIVERSITY**  
DENMARK

## DNS: Domain Names System

- ▶ Used everywhere, by everybody on the Internet
- ▶ ... also criminals!
- ▶ Technology/service choke-point
- ▶ Only small fraction of Internet traffic

## Example domains

- ▶ `www.example.com`
- ▶ `google.com`
- ▶ `yf32d9ac7f0a9f463e8da4736b12d7044a.tk`



Source: Antonakakis et al. 2011.

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

### 2 Motivation

Assumption and Method

Data

Features

Results

Conclusion

# Abuse and Malicious Domains



- ▶ Command and control
- ▶ Phishing
- ▶ Spamming
- ▶ Typo-squatting
- ▶ Blending in, tunnelling
- ▶ Fast-flux, Double-flux
- ▶ Domain-flux (DGA)



Source: Haymarket Media, Inc. <sup>1</sup>

**Malicious domain:** Any domain that is used for criminal, malicious or otherwise nefarious activities.

---

<sup>1</sup>: <https://www.scmagazineuk.com/cyber-criminals-becoming-increasingly-professional/article/531709/>

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

### 3 Motivation

Assumption and Method

Data

Features

Results

Conclusion

# Assumption and Method



Assumption:

***“Malicious domain names can be detected by lexical features.”***

Method:



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

4 Assumption and Method

Data

Features

Results

Conclusion

Machine Learning: Supervised, 10-fold cross validation, 10 repeats, Random Forrest, Python, Scikit-learn.

Table I: DATA SETS OF MALICIOUS DOMAIN NAMES.

Data sets		Number of domains	
		Non-DGA	DGA
abuse.ch	ZeuS Tracker	437	-
	Palevo Tracker	14	-
	Ransomware Tracker	1248	-
malware-\ domains.com	Just domains	13073	-
	Zeus Gameover	-	190033
	Conficker	-	106101
	Pushdo	-	10951
	GOZ	-	7348
	Microsoft Botnet	22036	-
host-file.net	Ad/tracking servers (ATS)	47960	-
	Malware distribution (EMD)	137237	-
	Exploits sites (EXP)	17282	-
	Fraud sites (FSA)	134501	-
	Spamming sites (GRM)	674	-
	Spamming sites (HFS)	573	-
	Hijack sites (HJK)	74	-
	Misleading marketing (MMT)	5533	-
	Pharmacy activities (PHA)	23143	-
	Phishing sites (PSH)	133913	-
	WareZ distribution (WRZ)	3231	-
datadriven\ security.info	Cryptolocker	-	34319
	Goz	-	7347
	New Goz	-	10999
malware\ domainlist.com	Malware-related domains	1253	-
malc0de.com	Malware-related domains	208	-
malwarepatrol.net	Malicious URLs	35518	-
phistank.com	Phishing URLs	14807	-
AAU-STAR	Domains from malw.testng	27778	-
Total		620493	367098

Table II: DATA SETS OF BENIGN DOMAIN NAMES.

Data sets		Number of domains	
		Non-DGA	DGA
alexa.org	Most popular domains	971424	-
datadriven\ security.info	Legitimate DGA domains	-	133927

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

5 Data

Features

Results

Conclusion



## 1. Basic Domain Features

- ▶ Count, Categorical
- ▶ TLD and FQDN

## 2. Simple Lexical Features

- ▶ Length, count, ratio
- ▶ 2LD, character classes

## 3. Advanced Lexical Features

- ▶ Approximating language, recognising words
- ▶ Entropy of letter distribution
- ▶ N-gram analysis (alex.org, English)

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

6 Features

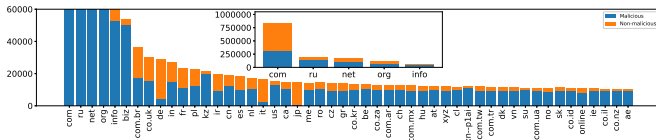
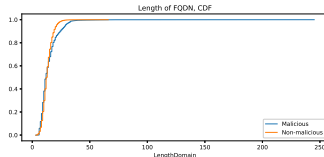
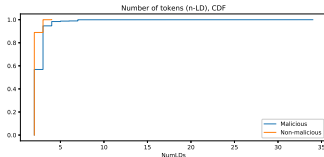
Results

Conclusion

---

Top-Level Domain (TLD): **www.example.com**.  
Second-Level Domain (2LD): **www.example.com**.  
Fully Qualified Domain Name: **www.example.com**

## Basic Domain Features



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

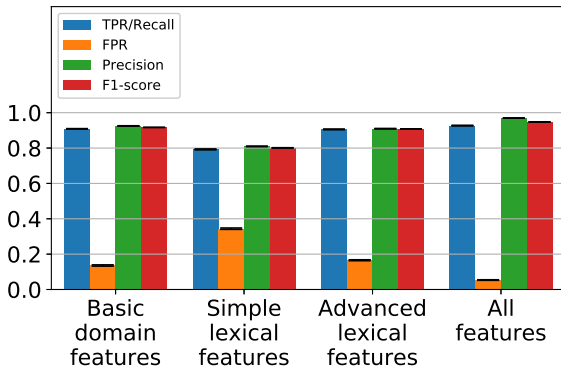
Features

7

Results

Conclusion

# Detection: Scenario I - Full Data Set



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

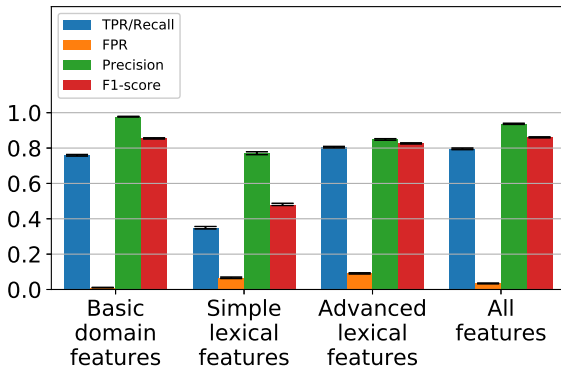
Features

8 Results

Conclusion



# Detection: Scenario II - Non-DGA



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

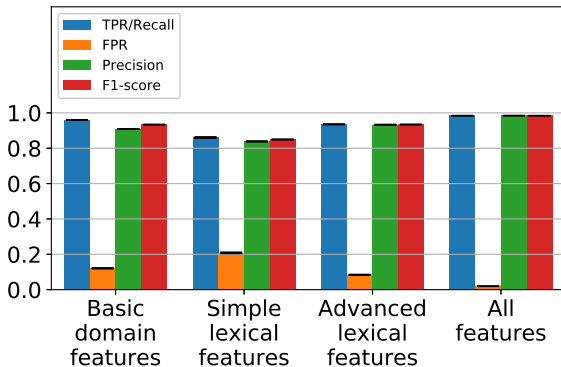
Data

Features

9 Results

Conclusion

# Detection: Scenario III - DGA



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

10 Results

Conclusion

Scenario III includes the benign non-DGA domains.

## Results

### Malicious domains can be detected by using lexical features only

- ▶ Detection of DGA-domains is particular promising
  - ▶ Precision: 0.984. Recall: 0.984. F1-score: 0.948.
- ▶ Detection of Non-DGA-domains requires further work
  - ▶ Precision: 0.937. Recall: 0.795. F1-score: 0.860.
- ▶ Using all features outperforms the individual sets

## General

- ▶ Vague distinction between DGA and Non-DGA
- ▶ Labelled data of high quality is hard to come by

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

11 Conclusion

# Conclusion



## Future work

- ▶ Explore other lexical features
- ▶ Explore effect of combining with non-lexical features
- ▶ Improve performance for non-DGA domains
- ▶ More substantial data sets
- ▶ Evaluate in a practical, online setting
- ▶ Apply unsupervised machine learning

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

12 Conclusion

# Backup slides



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

13 Conclusion



## Existing Solutions

- ▶ Reputation
- ▶ Traffic and activity (DNS, e-mail, ...)
- ▶ Resilience/anonymity techniques

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

14 Conclusion



## Basic Domain Features

1. Number of domain levels (n-LD).
2. Top level domain (TLD).
3. Length of Fully Qualified Domain Name (FQDN).

## Simple Lexical Features

1. Length of 2nd Level Domain (2-LD).
2. Ratio of consonants in the 2-LD.
3. Number of vowels in 2-LD.
4. Number of numeric characters in 2-LD.
5. Number of special characters in 2-LD.
6. Ratio of special characters in 2-LD.

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

15 Conclusion



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

16 Conclusion

## Advanced Lexical Features

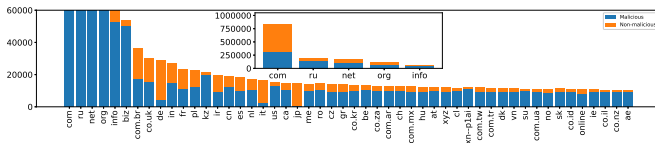
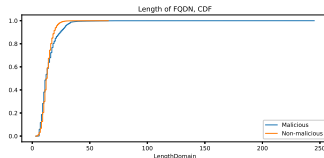
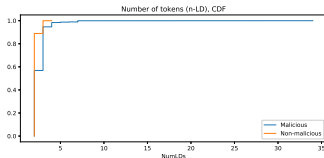
1. Language indicator (`langid.py`)
2. Number of English words in 2-LD.
3. Entropy of 2-LD.
4. N-gram analysis of 2-LD ([www.alexandria.org](http://www.alexandria.org))<sup>1</sup>.
5. N-gram analysis of 2-LD (English dictionary).

---

<sup>1</sup>{3,4,5}-grams,



## Basic Domain Features



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

17 Conclusion

## Simple Lexical Features

Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

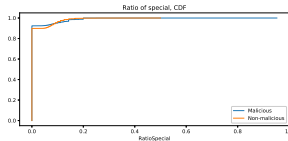
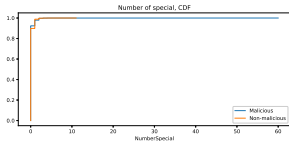
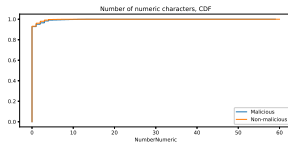
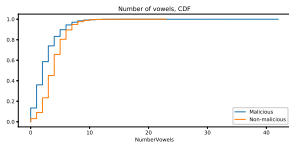
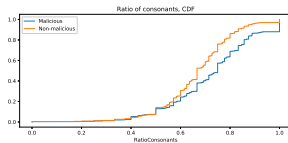
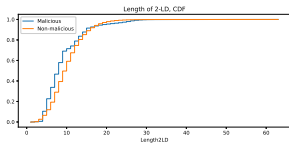
Assumption and Method

Data

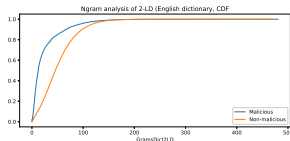
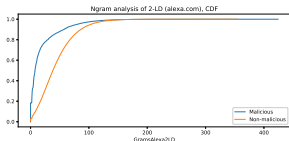
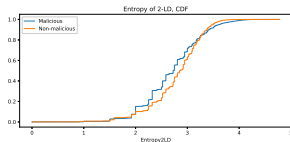
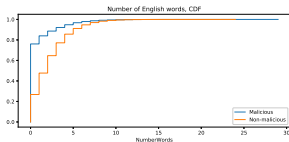
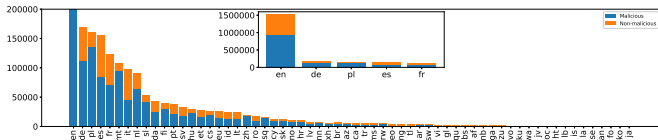
Features

Results

18 Conclusion



## Advanced Lexical Features



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

Conclusion

19

# Detection: Mean performance



Detection of malicious domains through lexical analysis

Egon Kidmose et al.

Motivation

Assumption and Method

Data

Features

Results

20 Conclusion

Scenario	Features	TPR/Recall	FPR	Precision	F1-score
I (Full)	1	9.09e-01	1.36e-01	9.25e-01	9.17e-01
	2	7.93e-01	3.44e-01	8.10e-01	8.01e-01
	3	9.06e-01	1.66e-01	9.10e-01	9.08e-01
	All	<b>9.27e-01</b>	<b>5.35e-02</b>	<b>9.70e-01</b>	<b>9.48e-01</b>
II (Non-DGA)	1	7.59e-01	<b>1.10e-02</b>	<b>9.78e-01</b>	8.55e-01
	2	3.50e-01	6.69e-02	7.69e-01	4.81e-01
	3	<b>8.05e-01</b>	9.16e-02	8.49e-01	8.26e-01
	All	7.95e-01	3.41e-02	9.37e-01	<b>8.60e-01</b>
III (DGA)	1	9.60e-01	1.21e-01	9.09e-01	9.34e-01
	2	8.61e-01	2.08e-01	8.38e-01	8.50e-01
	3	9.35e-01	8.40e-02	9.33e-01	9.34e-01
	All	<b>9.84e-01</b>	<b>1.98e-02</b>	<b>9.84e-01</b>	<b>9.84e-01</b>