

# Language Agnostic Botnet Detection based on ESOM and DNS

Christian Dietz and Rocco Mandrysch  
Urs Anliker and Gabi Dreo

Universität der Bundeswehr München and RUAG Defence

**botconf**  
30.11.2016



FZ *Forschungszentrum  
Cyber Defence*  
Universität der Bundeswehr München

**Together  
ahead. RUAG**

- 1 Introduction
- 2 Emergent Self-organizing Maps
- 3 Analysis and Results
- 4 Summary and Conclusion

# Introduction

## Motivation:

- Bots are using DNS protocol as a communication channel:
  - such as covert channel
  - to transport commands or information
- In many cases DNS names are generated via "Domain Generation Algorithm" (DGA)
- Example: Tinba, Pisloader, PadCrypt, ...

## Goal:

- Botnet detection via DGA identification with
  - Emergent Self-Organizing Maps (ESOMs)
  - language independent features from domain names

# Introduction

Our requirement is

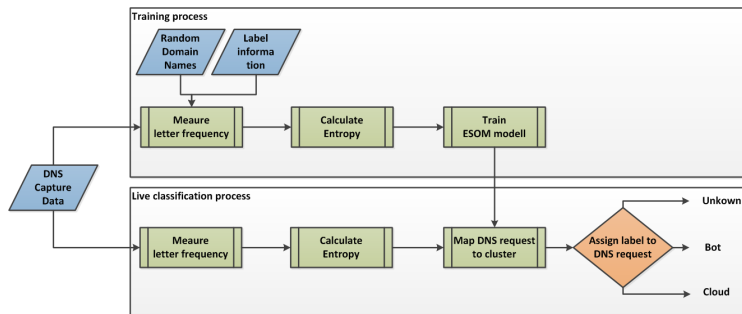
- to be language independent
- a distinction between cloud and malware DGAs

# Emergent Self-organizing Maps

- Self Organizing Map (SOM) is an artificial neural network (ANN), based on an unsupervised learning algorithm
  - NN are used to estimate or approximate function
  - unsupervised learning: type of algorithms that try to find correlations in and / or representations of the input data without any external informations (e.g. a classification) than the data itself
- Method of mapping high-dimensional input to low dimensional output such as (mostly) 2-D or 3-D
- Similar inputs are mapped to close locations on the low dimensional map such that the local topology is preserved.
- Large SOMs are called Emergent Self-Organizing Maps to emphasize the distinction.

# Analysis Strategy

- 1 Extract language features from domain names in data samples
- 2 Run ESOM training with compiled feature set
- 3 Categorize domain names from real live traffic with ESOM map



# Training data

- For training we used benign DNS and domain names from different malware DGAs

## Benign DNS

- is taken from real network data
  - contains only requests for
    - "standard domain names"
    - (auto generated) domain cloud names
- 
- DNS could be contaminated with auto generated domains.
  - An analysis with Alexa domain list and cloud domain names in separate training data samples is published in our paper.

# Training Data

- Malware DGAs with different characteristics among each other

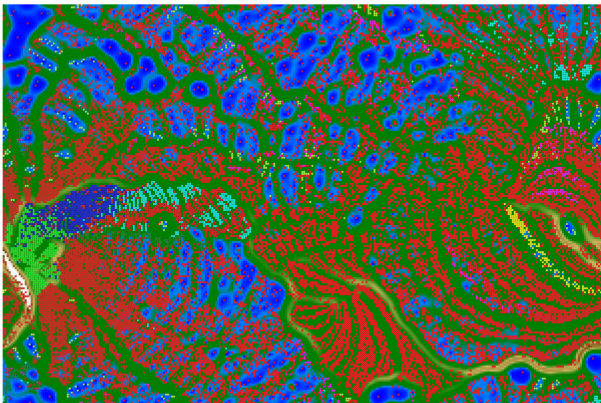
sample	characteristics
newgoz	doamin usually start with 1
ramnit	all letters except 'z'
shiotob	only digits 1,2,3,4,5,9
symmi	count sub-domains $> 2$ and all letters except 'zyj'
tinba	double counts of letters,e.g. 'jj'
murofet v3	numbers between 10 and 69
padcrypt	only letters: a,b,c,d,e,f,n,o,l,m,k



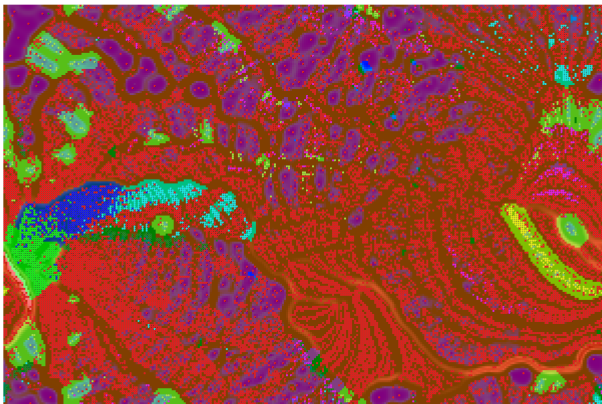
# Extracted features from domain names

language feature	third-level domain	second-level domain	top-level domain
letter frequency	×	×	×
digit frequency	×	×	×
vowel frequency	×	×	×
consonant frequency	×	×	×
Shannon entropy	×	×	×
bi-gram frequency	×	×	—
tri-gram frequency	×	×	—

- Calculated map with size  $200 \times 300$
- Every color stands for a data set
- Color gradient map: Earth
  - green (and then brown): low similarities between data points
  - blue: many similarities between data points



- Calculated surface of map
- Algorithm: gradient calculation between local minimas and maximas
- Color of area is dominated by color of most occurred data points
- Colored areas are used for domain name classification



- Compare results of every domain name in validation data set with every data point in map
- "False Positive" (FP) and "True Positive" (TP) classification rate for surface validation



# Calculating Detection Rate

- Compare results of every domain name in data set with every data point in map
- Rate  $\hat{=}$  domain names identified in a category divided by the total number of domains in a data sample

- Sample names in columns are the categories
- Rows are presenting the malware sample, which we want to detect

	benign dns	newgoz	padcrypt	ramnit	tinba
banjori	-	0.20	-	0.8	-
corebot	0.02	-	-	-	-
dircrypt	0.39	0.03	0.04	0.54	-
fobber	0.48	-	0.03	0.49	-
gozi	0.67	-	0.09	0.24	-
kraken	0.62	-	0.03	0.13	0.22
locky	0.36	-	0.03	0.13	0.47
necurs	0.5	0.04	0.02	0.07	0.36
nymain	0.66	-	0.03	0.17	0.14
proslikefan	0.63	-	0.05	0.09	0.24
pykspa	0.63	-	0.05	0.09	0.24
qadars	0.87	-	0.07	0.06	-
qakbot	0.23	0.1	0.03	0.53	0.09
ranbyus	0.44	0.01	-	0.19	0.35
simba	1	-	-	-	-
vawtrak	0.74	-	0.02	0.24	-

## Summary and Conclusion

- Presented a novel approach of detecting Botnets based on ESOMs.
- Method is language feature independent
- Our approach can classify between benign and botnet DGA domains.
- Next steps:
  - Evaluating uncertainties
  - Tune training parameters
  - Running long-term validation in real live networks

---

Thank you!

---